

Numerical Methods I

MATH-GA 2010.001/CSCI-GA 2420.001

Benjamin Peherstorfer
Courant Institute, NYU

Based on slides by G. Stadler and A. Donev

Outline

Organization

Conditioning of problems

Stability of algorithms

Representing real numbers

Organization

- ▶ **Time and location:** Mondays and Wednesdays 4:55–6:10PM, WWH 1302
- ▶ **Office hours:** Mondays, 6.10 – 7.10pm, stop by or make an appointment (please email). My office number is WWH #421
- ▶ **Course webpage:** https://docs.google.com/document/d/1VxdM4s-wiV-_C4uBDrP4ioi0dscfLwim8o1mokuj2oQ/edit?usp=sharing

You need to be logged in with your NYU-Google account to access it

- ▶ **Brightspace** <https://brightspace.nyu.edu/d2l/home/400947>

Required work and grading:

- ▶ 7 homework assignments (60% of your grade).

Required work and grading:

- ▶ 7 homework assignments (60% of your grade).
 - ▶ These will be mixed paper&pencil and computational/programming. You hand in solutions in LaTeX and Matlab code

Required work and grading:

- ▶ **7 homework assignments** (60% of your grade).
 - ▶ These will be mixed paper&pencil and computational/programming. You hand in solutions in LaTeX and Matlab code
 - ▶ Grader is Qi (Winston) Liang, lq504@nyu.edu

Required work and grading:

- ▶ **7 homework assignments** (60% of your grade).
 - ▶ These will be mixed paper&pencil and computational/programming. You hand in solutions in LaTeX and Matlab code
 - ▶ Grader is Qi (Winston) Liang, lq504@nyu.edu
 - ▶ Grader's office hour: tbd

Required work and grading:

- ▶ **7 homework assignments** (60% of your grade).
 - ▶ These will be mixed paper&pencil and computational/programming. You hand in solutions in LaTeX and Matlab code
 - ▶ Grader is Qi (Winston) Liang, lq504@nyu.edu
 - ▶ Grader's office hour: tbd
 - ▶ You are welcome to discuss with your colleagues, but *you've to write up your solution independently and write every line of code yourself.*

Required work and grading:

- ▶ **7 homework assignments** (60% of your grade).
 - ▶ These will be mixed paper&pencil and computational/programming. You hand in solutions in LaTeX and Matlab code
 - ▶ Grader is Qi (Winston) Liang, lq504@nyu.edu
 - ▶ Grader's office hour: tbd
 - ▶ You are welcome to discuss with your colleagues, but *you've to write up your solution independently and write every line of code yourself.*
 - ▶ The first homework assignment has been posted.

Required work and grading:

- ▶ **7 homework assignments** (60% of your grade).
 - ▶ These will be mixed paper&pencil and computational/programming. You hand in solutions in LaTeX and Matlab code
 - ▶ Grader is Qi (Winston) Liang, lq504@nyu.edu
 - ▶ Grader's office hour: tbd
 - ▶ You are welcome to discuss with your colleagues, but *you've to write up your solution independently and write every line of code yourself.*
 - ▶ The first homework assignment has been posted.
 - ▶ Follow the guidelines/rules provided with the homework to present and submit your solution.

Required work and grading:

- ▶ **7 homework assignments** (60% of your grade).
 - ▶ These will be mixed paper&pencil and computational/programming. You hand in solutions in LaTeX and Matlab code
 - ▶ Grader is Qi (Winston) Liang, lq504@nyu.edu
 - ▶ Grader's office hour: tbd
 - ▶ You are welcome to discuss with your colleagues, but *you've to write up your solution independently and write every line of code yourself.*
 - ▶ The first homework assignment has been posted.
 - ▶ Follow the guidelines/rules provided with the homework to present and submit your solution.
 - ▶ If you find errors/typos in the slides or homework assignments, please let me know.

Required work and grading:

- ▶ **7 homework assignments** (60% of your grade).
 - ▶ These will be mixed paper&pencil and computational/programming. You hand in solutions in LaTeX and Matlab code
 - ▶ Grader is Qi (Winston) Liang, lq504@nyu.edu
 - ▶ Grader's office hour: tbd
 - ▶ You are welcome to discuss with your colleagues, but *you've to write up your solution independently and write every line of code yourself.*
 - ▶ The first homework assignment has been posted.
 - ▶ Follow the guidelines/rules provided with the homework to present and submit your solution.
 - ▶ If you find errors/typos in the slides or homework assignments, please let me know.
 - ▶ Late homework policy: We understand that there are sometimes unexpected circumstances. You can hand in one homework late by 24h this semester if you send me and the grader an email at least 24h before the deadline that your work will be late.

Required work and grading:

- ▶ **7 homework assignments** (60% of your grade).
 - ▶ These will be mixed paper&pencil and computational/programming. You hand in solutions in LaTeX and Matlab code
 - ▶ Grader is Qi (Winston) Liang, lq504@nyu.edu
 - ▶ Grader's office hour: tbd
 - ▶ You are welcome to discuss with your colleagues, but *you've to write up your solution independently and write every line of code yourself.*
 - ▶ The first homework assignment has been posted.
 - ▶ Follow the guidelines/rules provided with the homework to present and submit your solution.
 - ▶ If you find errors/typos in the slides or homework assignments, please let me know.
 - ▶ Late homework policy: We understand that there are sometimes unexpected circumstances. You can hand in one homework late by 24h this semester if you send me and the grader an email at least 24h before the deadline that your work will be late.
- ▶ A **final** (40% of grade)

Required work and grading:

- ▶ **7 homework assignments** (60% of your grade).
 - ▶ These will be mixed paper&pencil and computational/programming. You hand in solutions in LaTeX and Matlab code
 - ▶ Grader is Qi (Winston) Liang, lq504@nyu.edu
 - ▶ Grader's office hour: tbd
 - ▶ You are welcome to discuss with your colleagues, but *you've to write up your solution independently and write every line of code yourself.*
 - ▶ The first homework assignment has been posted.
 - ▶ Follow the guidelines/rules provided with the homework to present and submit your solution.
 - ▶ If you find errors/typos in the slides or homework assignments, please let me know.
 - ▶ Late homework policy: We understand that there are sometimes unexpected circumstances. You can hand in one homework late by 24h this semester if you send me and the grader an email at least 24h before the deadline that your work will be late.
- ▶ A **final** (40% of grade)
 - ▶ Expect to write *some* code too

Organization issues

Prerequisites:

- ▶ Basic linear algebra; calculus; experience in Matlab (or Python or another programming language)

There is a part II of this class. . .

- ▶ . . . in the Spring semester. You should take both parts to get a reasonably complete overview of Numerical Methods.
- ▶ If you consider taking only one semester of Numerical Methods, I recommend taking Scientific Computing this semester instead of this class.

Topics covered in Numerical Methods I

Numerical Methods and their Analysis

- ▶ Stability; sources of errors; error propagation, representation of numbers in computers
- ▶ Numerical linear algebra: direct solution of sparse/dense linear systems; solution of least square systems; eigenvalue problems; iterative solution of linear systems
- ▶ Nonlinear systems; Newton's method; Nonlinear least squares
- ▶ Numerical optimization
- ▶ Interpolation and Approximation
- ▶ Numerical integration

Computing Issues

- ▶ What makes some computer codes faster than others?
- ▶ Where are numerical methods used, and what is their role in science research?
- ▶ How large/complicated problems can we solve today? Where are the challenges and limits of what we can do?

Topics of Numerical Methods II

Main topics covered in Numerical Methods II in the Spring semester

- ▶ Approximation of ordinary differential equations (ODEs)
- ▶ Approximation of partial differential equations (PDEs)
- ▶ Solvers for the resulting (high-dimensional) discrete problems

Programming

Programming the methods we discuss is an integral part of this course. To really understand methods & algorithms, one needs to implement them and experiment with them.

- ▶ Make sure you have access to **MATLAB** (CIMS, student license), you will need it for the first homework assignment.
- ▶ Alternatives to MATLAB: **Octave, Python or Julia**.
- ▶ We will talk about a few best coding practices, and how to present results.

Recommended textbooks/literature:

Text books:

- ▶ P. Deuffhard, A. Hohmann: *Numerical Analysis in Modern Scientific Computing. An Introduction*, 2nd edition, Springer, 2003.
- ▶ L. N. Trefethen, D. Bau: *Numerical Linear Algebra*, SIAM, 1997.
- ▶ A. Quarteroni, R. Sacco, F. Saleri: *Numerical Mathematics*, 2nd edition, Springer, 2007.
- ▶ M. Overton: *Numerical Computing with IEEE Floating Point Arithmetic*, SIAM, 2004.

Matlab/Programming:

- ▶ W. Gander, M. J. Gander, F. Kwok: *Scientific Computing - An Introduction Using Maple and MATLAB*. Texts in Computation Science and Engineering. Springer, 2014.
- ▶ C. Moler: *Numerical Computing with Matlab*, SIAM, 2007.

Numerical mathematics

Computer simulations have had a big influence on research and development; sometimes the ability to simulate phenomena is referred to as the **third pillar of science**.

Numerical mathematics is a part of mathematics that **develops, analyzes** and **applies** methods from scientific computing to

- ▶ analysis
- ▶ linear algebra
- ▶ optimization
- ▶ differential equations
- ▶ ...

It has applications accross many applied sciences, including:

- ▶ physics
- ▶ economics
- ▶ biology
- ▶ finance
- ▶ ...

Development of Numerical Methods at Courant

A few examples. . .

- ▶ Eigenvalue problems (Overton)
- ▶ Fast multipole method (Greengard, O'Neil, Zorin)
- ▶ Finite elements and contact problems (Panozzo, Zorin)
- ▶ Methods for studying dynamical systems, multiscale methods (Vanden-Eijnden)
- ▶ Methods for free boundary problems in fluid dynamics (Shelley)
- ▶ Scalable implicit solvers for viscous flows (Stadler)
- ▶ Sampling methods and Uncertainty Quantification (Goodman, Stadler, Peherstorfer)
- ▶ Scientific machine learning (Vanden-Eijnden, Stadler, Peherstorfer)
- ▶ . . .

Applications of Numerical Methods at Courant

A few examples. . .

- ▶ Simulation and analysis of natural and artificial heart valves (Peskin)
- ▶ Simulation of plate tectonics and mantle convection (Stadler)
- ▶ The physics of cell's interiors and their motion (Shelley)
- ▶ Optimal complexity wave simulations (Greengard)
- ▶ Simulation of blood cells-resolving blood flow (Zorin)
- ▶ Plasmas (Stadler, Kaptanoglu)
- ▶ . . .


Seminars

Computational Mathematics and Scientific Computing seminar

- ▶ Fridays at 10:00, WWH 1302
- ▶ Talks about current research
- ▶ <https://cims.nyu.edu/dynamic/calendars/seminars/computational-mathematics-and-scientific-computing-seminar/>

The screenshot shows the top navigation bar of the Courant Institute of Mathematical Sciences website. It includes the NYU logo, the text "COURANT INSTITUTE OF MATHEMATICAL SCIENCES", a search bar, and a "Courant Login" link. Below the navigation bar is a dark menu with links for "Institute", "Academics", "Research", "People", "Calendars", "Resources", "About Us", and "Giving". The main content area features the title "Computational Mathematics and Scientific Computing Seminar" in a large, bold font. Below the title, there is a paragraph of text explaining the seminar's format and contact information for organizers Georg Stadler and Benjamin Peherstorfer. A link is provided for subscribing to the seminar mailing list. At the bottom, there is a dropdown menu for "Fall 2023" and a section titled "Upcoming Events" with a list item for "Friday, September 15, 2023".

NYU | COURANT INSTITUTE OF MATHEMATICAL SCIENCES

Search  [Courant Login](#)

[Institute](#) ▾ [Academics](#) ▾ [Research](#) ▾ [People](#) ▾ [Calendars](#) ▾ [Resources](#) ▾ [About Us](#) ▾ [Giving](#)

Computational Mathematics and Scientific Computing Seminar

The Computational Mathematics and Scientific Computing seminar will be in person at the usual time on Fridays at 10am unless otherwise noted. In rare cases we have zoom talks and then the zoom link to join the seminar will be sent to the seminar mailing list. Contact the organizers [Georg Stadler](#) and [Benjamin Peherstorfer](#) if you haven't received the Zoom link.

To subscribe to the CMSC seminar mailing list, please see [here](#).

Seminar Organizer(s): Georg Stadler and Benjamin Peherstorfer

[Fall 2023](#) ▾

Upcoming Events

- Friday, September 15, 2023

Modeling and Simulation meeting

- ▶ Thursdays at 12:30, WWH 1302
- ▶ Student-driven meeting on topics related to computational mathematics
- ▶ <https://math.nyu.edu/dynamic/research/pages/research-and-training-group-mathematical-modeling-and-simulation/activities/group-meeting/>

Modeling and Simulation meeting

- ▶ Thursdays at 12:30, WWH 1302
- ▶ Student-driven meeting on topics related to computational mathematics
- ▶ <https://math.nyu.edu/dynamic/research/pages/research-and-training-group-mathematical-modeling-and-simulation/activities/group-meeting/>

Mathematics colloquium

- ▶ Mondays at 3:45, WWH 1302
- ▶ <https://math.nyu.edu/dynamic/calendars/seminars/mathematics-colloquium/>

Modeling and Simulation meeting

- ▶ Thursdays at 12:30, WWH 1302
- ▶ Student-driven meeting on topics related to computational mathematics
- ▶ <https://math.nyu.edu/dynamic/research/pages/research-and-training-group-mathematical-modeling-and-simulation/activities/group-meeting/>

Mathematics colloquium

- ▶ Mondays at 3:45, WWH 1302
- ▶ <https://math.nyu.edu/dynamic/calendars/seminars/mathematics-colloquium/>

Math and data

- ▶ Thursdays at 2.00, Auditorium Hall 150, Center for Data Science, NYU, 60 5th ave.
- ▶ Interface of Applied Mathematics, Statistics and Machine Learning
- ▶ <https://mad.cds.nyu.edu/seminar/>

Outline

Organization

Conditioning of problems

Stability of algorithms

Representing real numbers

Condition of a problem

- ▶ Consider a generic problem: given F and data/input x , find output y such that

$$F(x, y) = 0$$

- ▶ Let's assume there is a unique solution so that we can write

$$y = f(x),$$

for a function f in the following

- ▶ Well-posed: Unique solution + If we perturb the input x a little bit, the solution y gets perturbed by a small amount.
- ▶ Otherwise, the problem is ill-posed; no numerical method can help with that.
(What should we do in such a situation?)

Condition of a problem

- ▶ Consider a generic problem: given F and data/input x , find output y such that

$$F(x, y) = 0$$

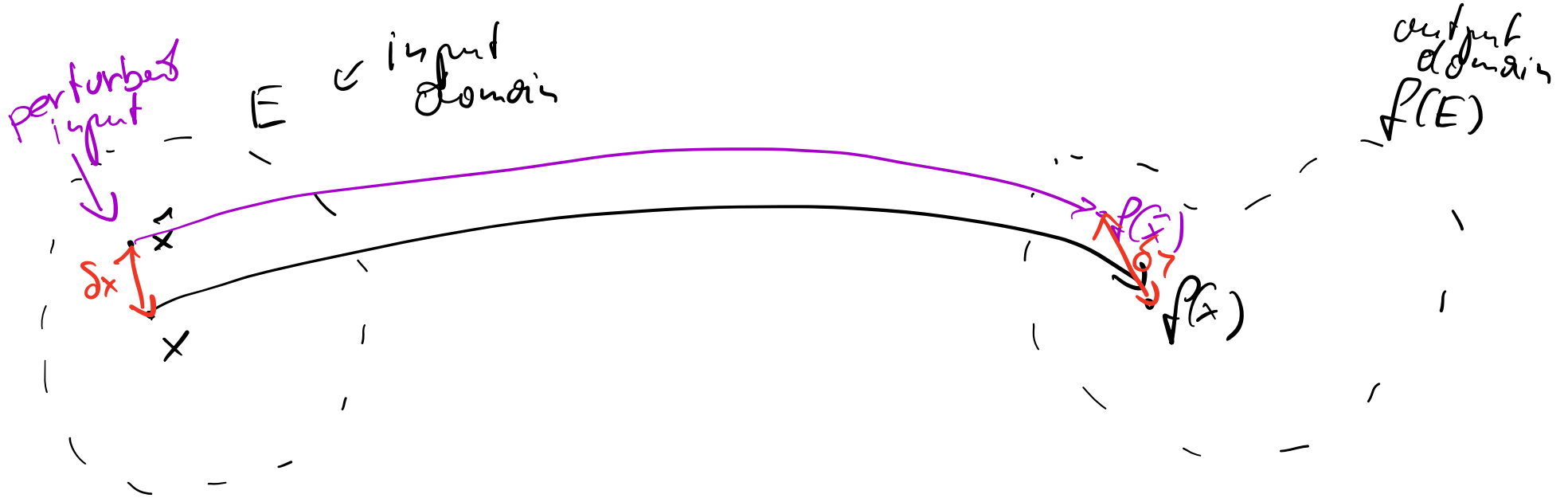
- ▶ Let's assume there is a unique solution so that we can write

$$y = f(x),$$

for a function f in the following

- ▶ Well-posed: Unique solution + If we perturb the input x a little bit, the solution y gets perturbed by a small amount.
- ▶ Otherwise, the problem is ill-posed; no numerical method can help with that.
(What should we do in such a situation? \rightsquigarrow change the problem)

Condition of a problem (visualization)

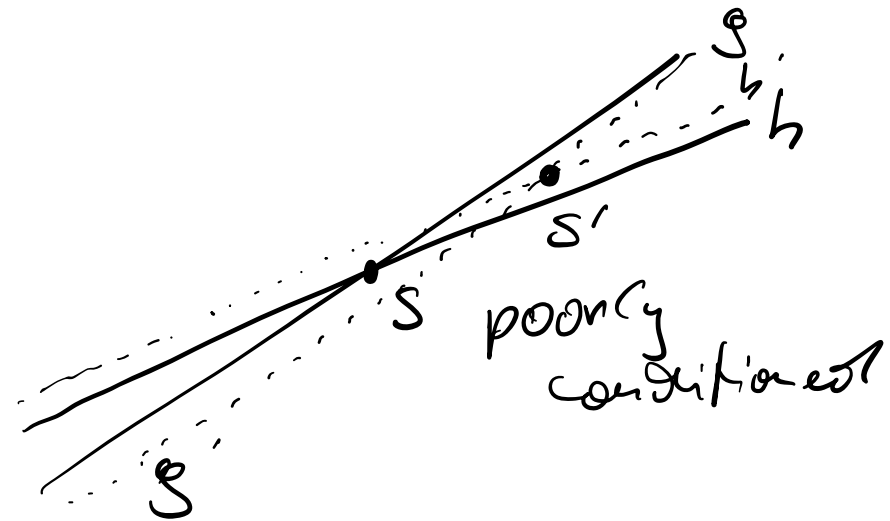
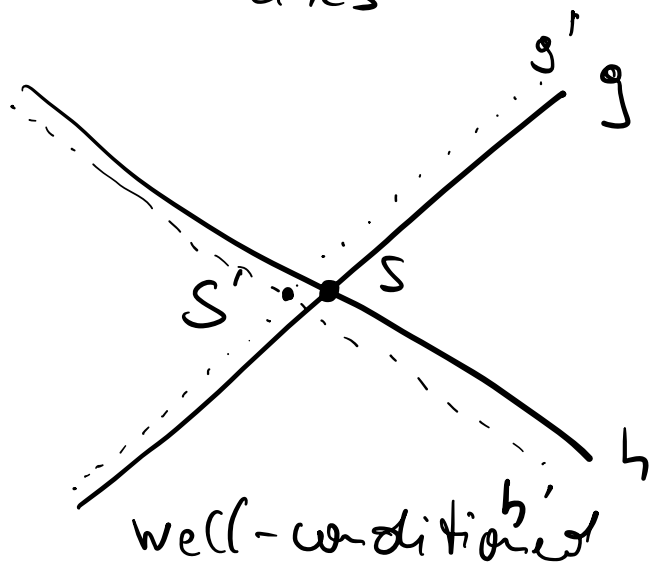


Notice: $\delta y = f(\bar{x}) - f(x)$
is independent of any algorithm
that we use. There is no numerics here!

Condition of a problem (intersecting lines)

Consider the problem of determining the intersection point of two lines

$f: (g, h) \mapsto s \leftarrow \text{intersection point}$
↑↑
lines



Condition of a problem (cont'd)

- ▶ Terms such as “little bit” and a “small amount” already point to that we need to measure something
- ▶ Therefore, we assume the map f is given as

$$f : U \subset \mathbb{R}^n \rightarrow \mathbb{R}^m$$

and we are interested in the norm $\| \cdot \|$

- ▶ The input error is then

$$\|x - \hat{x}\| \leq \delta \text{ (absolute)} \quad \|x - \hat{x}\| \leq \delta \|x\| \text{ (relative)}$$

- ▶ Correspondingly we measure the output error $f(x) - f(\hat{x})$ in $\| \cdot \|$ (we could also have looked at a componentwise error)

Condition of a problem (cont'd)

- ▶ Absolute condition number at x is

$$\kappa_{\text{abs}} = \lim_{\delta \rightarrow 0} \sup_{\|x - \hat{x}\| \leq \delta} \frac{\|f(x) - f(\hat{x})\|}{\|x - \hat{x}\|}$$

- ▶ Relative condition number at x is

$$\kappa_{\text{rel}} = \lim_{\delta \rightarrow 0} \sup_{\|x - \hat{x}\| \leq \delta} \frac{\|f(x) - f(\hat{x})\| / \|f(x)\|}{\|x - \hat{x}\| / \|x\|}$$

- ▶ If f is differentiable in x , then

$$\kappa_{\text{abs}} = \|f'(x)\| \quad \kappa_{\text{rel}} = \frac{\|x\|}{\|f(x)\|} \|f'(x)\|,$$

where $\|f'(x)\|$ is the norm of the Jacobian $f'(x)$ in the operator norm

$$\|A\| = \sup_{x \neq 0} \frac{\|Ax\|}{\|x\|} = \sup_{\|x\|=1} \|Ax\|$$

- ▶ Another way to interpret the condition number at x is via the bounds

$$\|f(\hat{x}) - f(x)\| \leq \kappa_{\text{abs}} \|\hat{x} - x\|$$

and

$$\frac{\|f(\hat{x}) - f(x)\|}{\|f(x)\|} \leq \kappa_{\text{rel}} \frac{\|\hat{x} - x\|}{\|x\|},$$

for infinitesimal δ (or $\hat{x} \rightarrow x$)

Condition of a problem (cont'd)

▶ If $\kappa_{\text{rel}} \sim 1$,

Condition of a problem (cont'd)

- ▶ If $\kappa_{\text{rel}} \sim 1$, then the problem is well conditioned: If the relative error in the data/input is small, then the relative error in the answer/output is similarly small
- ▶ If $\kappa_{\text{rel}} \gg 1$,

Condition of a problem (cont'd)

- ▶ If $\kappa_{\text{rel}} \sim 1$, then the problem is well conditioned: If the relative error in the data/input is small, then the relative error in the answer/output is similarly small
- ▶ If $\kappa_{\text{rel}} \gg 1$, then the problem is poorly conditioned: Small relative input error can lead to large relative output error
- ▶ If κ_{rel} (and κ_{abs}) do not exist, then the problem is ill conditioned.
- ▶ What is poorly conditioned depends on desired accuracy: if the input accuracy is low but we expect a high output accuracy, then problems are quickly poorly conditioned. If we are happy with a less accurate output, we might consider the problem still well conditioned.
- ▶ Sometimes, the possibly large error in the output does not matter and so we can solve poorly conditioned problems (think of early design stages, rapid prototyping, etc); but we should be very much aware of the condition of the problem.

Condition of a problem: Example

Condition number of "addition"

linear map

$$f: \mathbb{R}^2 \rightarrow \mathbb{R}, \quad (a, b) \mapsto f(a, b) = a + b$$

derivative

$$f'(a, b) = [1, 1] \in \mathbb{R}^{1 \times 2}$$

choose 1-norm $\|\cdot\|_1$ on \mathbb{R}^2

$$\left\| \begin{bmatrix} a \\ b \end{bmatrix} \right\|_1 = |a| + |b|$$

Induced operator norm is max column sum

$$\|f'(a,b)\|_1 = \|[1, 1]\|_1 = 1$$

$$\kappa_{\text{obs}} = \|f'(a,b)\|_1 = 1$$

$$\begin{aligned}\kappa_{\text{rel}} &= \frac{\|x\|}{\|f(x)\|} \|f'(x)\| = \\ &= \frac{|a| + |b|}{|a + b|} \cdot 1\end{aligned}$$

a, b have the same sign

if $a \cdot b \geq 0$, $\kappa_{\text{rel}} = 1 \rightarrow$ perfect

a, b don't have same sign

$a \cdot b \leq 0 \rightarrow$ subtraction

set $a > 0$, $b = -(a + \epsilon)$, $0 < \epsilon \ll 1$

$$\kappa_{\text{rel}} = \frac{a + a + \epsilon}{|a - (a + \epsilon)|} = \frac{2a + \epsilon}{\epsilon} \gg 1$$

\leadsto poorly conditioned

Condition number of Mat Vec: $x \mapsto Ax$

$$f(x) = Ax$$

$$f'(x) = A$$

$$\kappa_{\text{abs}} = \|f'(x)\| = \|A\| = \sup_{\|x\|=1} \|Ax\|$$

$$\kappa_{\text{rel}} = \frac{\|x\|}{\|Ax\|} \|A\|$$

Norm plays a critical role!

Norm of $\|A\|$ is important

Example: take $\|\cdot\|_2$ on input x , then operator norm

$$\|A\|_2 = \rho_{\text{max}}(A)$$

↑ largest singular value

Singular values tell us by how much a vector gets 'stretched' by matrix

↑
perturbation

Numerical Methods I

MATH-GA 2010.001/CSCI-GA 2420.001

Benjamin Peherstorfer
Courant Institute, NYU

Based on slides by G. Stadler and A. Donev

Today

Last time

- ▶ Condition of problems

Today

- ▶ More on condition of problems
- ▶ Stability of algorithms
- ▶ Matlab recap

Announcements

- ▶ Homework 1 was posted last week; is due in two weeks Mon, Sep 23 *before class*

Recap: Condition of a problem

- ▶ Consider a generic problem: given F and data/input x , find output y such that

$$F(x, y) = 0$$

- ▶ Let's assume there is a unique solution so that we can write

$$y = f(x),$$

for a function f in the following

- ▶ Well-posed: Unique solution + If we perturb the input x a little bit, the solution y gets perturbed by a small amount.
- ▶ Otherwise, the problem is ill-posed; no numerical method can help with that.
(What should we do in such a situation?)

Recap: Condition of a problem

- ▶ Consider a generic problem: given F and data/input x , find output y such that

$$F(x, y) = 0$$

- ▶ Let's assume there is a unique solution so that we can write

$$y = f(x),$$

for a function f in the following

- ▶ Well-posed: Unique solution + If we perturb the input x a little bit, the solution y gets perturbed by a small amount.
- ▶ Otherwise, the problem is ill-posed; no numerical method can help with that. (What should we do in such a situation? \rightsquigarrow change the problem)

Recap: Condition of a problem (cont'd)

- ▶ Absolute condition number at x is

$$\kappa_{\text{abs}} = \lim_{\delta \rightarrow 0} \sup_{\|x - \hat{x}\| \leq \delta} \frac{\|f(x) - f(\hat{x})\|}{\|x - \hat{x}\|}$$

- ▶ Relative condition number at x is

$$\kappa_{\text{rel}} = \lim_{\delta \rightarrow 0} \sup_{\|x - \hat{x}\| \leq \delta} \frac{\|f(x) - f(\hat{x})\| / \|f(x)\|}{\|x - \hat{x}\| / \|x\|}$$

- ▶ If f is differentiable in x , then

$$\kappa_{\text{abs}} = \|f'(x)\| \quad \kappa_{\text{rel}} = \frac{\|x\|}{\|f(x)\|} \|f'(x)\|,$$

where $\|f'(x)\|$ is the norm of the Jacobian $f'(x)$ in the operator norm

$$\|A\| = \sup_{x \neq 0} \frac{\|Ax\|}{\|x\|} = \sup_{\|x\|=1} \|Ax\|$$

Recap: Condition of a problem (cont'd)

- ▶ If $\kappa_{\text{rel}} \sim 1$,

Recap: Condition of a problem (cont'd)

- ▶ If $\kappa_{\text{rel}} \sim 1$, then the problem is well conditioned: If the relative error in the data/input is small, then the relative error in the answer/output is similarly small
- ▶ If $\kappa_{\text{rel}} \gg 1$,

Recap: Condition of a problem (cont'd)

- ▶ If $\kappa_{\text{rel}} \sim 1$, then the problem is well conditioned: If the relative error in the data/input is small, then the relative error in the answer/output is similarly small
- ▶ If $\kappa_{\text{rel}} \gg 1$, then the problem is poorly conditioned: Small relative input error can lead to large relative output error
- ▶ If κ_{rel} (and κ_{abs}) do not exist, then the problem is ill conditioned.
- ▶ What is poorly conditioned depends on desired accuracy: if the input accuracy is low but we expect a high output accuracy, then problems are quickly poorly conditioned. If we are happy with a less accurate output, we might consider the problem still well conditioned.
- ▶ Sometimes, the possibly large error in the output does not matter and so we can solve poorly conditioned problems (think of early design stages, rapid prototyping, etc); but we should be very much aware of the condition of the problem.

Condition of a problem: Example

Condition number of MatVec: $x \mapsto Ax$

$$f(x) = Ax$$

$$f'(x) = A$$

$$\kappa_{\text{abs}} = \|f'(x)\| = \|A\|$$

$$\kappa_{\text{rel}} = \frac{\|x\|}{\|Ax\|} \|A\|$$

norm plays a critical role!

if we take $\|\cdot\|_2$ on input x , then operator norm

$$\|A\|_2 = \sigma_{\max}(A)$$

Singular values tell us by how much
a vector gets 'stretched' by matrix
↑ perturbation

Condition number of solving linear systems

$$Ax = b$$

(1) Problem: fix A , consider $b \mapsto A^{-1}b$

$$f: \mathbb{R}^n \rightarrow \mathbb{R}^n, \quad f(b) = A^{-1}b$$

differentiable

$$\kappa_{\text{abs}} = \|f'(b)\| = \|A^{-1}\|$$

$$\kappa_{\text{rel}} = \frac{\|b\|}{\|A^{-1}b\|} \|A^{-1}\|$$

if $A^{-1}b = x$

$$\kappa_{\text{rel}} = \frac{\|Ax\|}{\|x\|} \|A^{-1}\| \leq \frac{\|A\| \|x\|}{\|x\|} \cdot \|A^{-1}\|$$

$$= \|A\| \|A^{-1}\| = \kappa(A)$$

a condition number of matrix A

If we take $\|\cdot\|_2$ on input b , then operator norm

$$\|A^{-1}\|_2 = \sigma_{\max}(A^{-1}) = \frac{1}{\sigma_{\min}(A)}$$

then $\kappa(A) = \frac{\sigma_{\max}(A)}{\sigma_{\min}(A)}$

(2) Problem: fix b , consider $A \mapsto A^{-1}b$

$$f(A) = A^{-1}b$$

is differentiable

\leadsto feedback by
Deuschel

$$\|f'(A)\| \leq \|A^{-1}\| \|x\|$$

and $x = A^{-1}b$

$$\kappa_{\text{abs}} = \|f'(A)\| \leq \|A^{-1}\| \|x\|$$

$$\kappa_{\text{rel}} = \frac{\|A\|}{\|f(A)\|} \|f'(A)\|$$

$$\leq \frac{\|A\|}{\|x\|} \|A^{-1}\| \|x\| = \|A\| \cdot \|A^{-1}\| = \kappa(A)$$

Condition number $\kappa(A)$ of matrix A , tells us how well system of linear equations $Ax=b$ is conditioned.

Outline

Organization

Conditioning of problems

Stability of algorithms

Representing real numbers

Stability of algorithms

Is $\tilde{f}(x)$, computed with an algorithm \tilde{f} , a good approximation of $f(x)$?

We are happy if the error due to the algorithm

$$\tilde{f}(x) - f(x)$$

lies within reasonable bounds of the error due to the input

$$f(\tilde{x}) - f(x)$$

Stability: Stability

We say that an algorithm \tilde{f} for a problem f is stable if for each $x \in E$ the error

$$\frac{\|\tilde{f}(x) - f(\tilde{x})\|}{\|f(\tilde{x})\|}$$

is small for \tilde{x} with small

$$\frac{\|\tilde{x} - x\|}{\|x\|}$$

A stable algorithm gives nearly the right answer ($\tilde{f}(x)$) to nearly the right question ($f(\tilde{x})$).

In forward error analysis one tries to establish stability by showing error bounds on the result in each operation in the algorithm in order to bound the error in the end result

Stability: Backward stability

Backward stability: Pass the errors of the algorithm back and interpret as input errors.

An algorithm \tilde{f} for a problem f is backward stable if for each $x \in X$ we have $\tilde{f}(x) = f(\tilde{x})$ for an \tilde{x} with

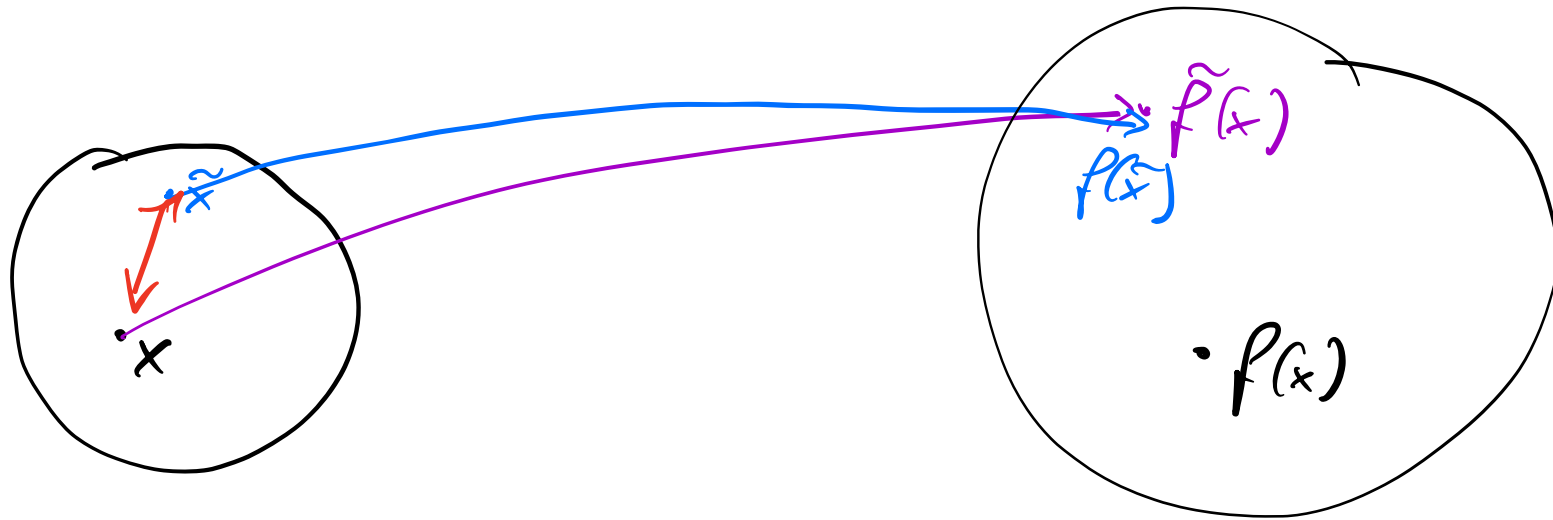
$$\frac{\|\tilde{x} - x\|}{\|x\|}$$

small

This is a tightening of the definition of stability of the previous slide:

A backward stable algorithm gives exactly the right answer to nearly the right question.

In backward error analysis one calculates, for a given output, how much one would need to perturb the input in order for the answer to be exact.



algorithmic error is interpreted as input error

Errors and error analyses

Relative errors:

$$\frac{\|x - x_n\|}{\|x\|}$$

Absolute error:

$$\|x - x_n\|$$

- ▶ Used for theoretical arguments
- ▶ In numerical practice: exact solution is not available, so these errors must be approximated.

A priori analysis is performed before a specific solution is computed. Typically, the analysis is performed for a large class of possible inputs.

A posteriori analysis bounds the error for a specific numerical solution \hat{x} (computed with a specific numerical method), and uses, e.g., residuals for the a posteriori analysis.

Computational errors

Numerical algorithms try to control or minimize, rather than eliminate, the various computational errors:

- ▶ **Approximation error** due to replacing the computational problem with an easier-to-solve approximation. Also called discretization error for ODEs/PDEs.
- ▶ **Truncation error** due to replacing limits and infinite sequences and sums by a finite number of steps. Closely related to approximation error.
- ▶ **Roundoff error** due to finite representation of real numbers and arithmetic on the computer, $x \neq \hat{x}$.
- ▶ **Propagated error** due to errors in the data from user input or previous calculations in iterative methods.
- ▶ **Statistical error** in stochastic calculations such as Monte Carlo calculations.

Intuition: Stability, Consistency, Convergence

Instead of solving $F(x, y) = 0$ directly, many numerical methods generate a solution sequentially

$$\bar{F}(x_i, x_{i-1}) = 0, \quad i = 1, 2, 3, \dots,$$

with $x_0 = x$ and sequence (x_i) converging to y

Additionally, we use a numerical method \hat{F}_n instead of \bar{F}

$$\hat{F}_n(\hat{x}_i, \hat{x}_{i-1}) = 0, \quad i = 1, 2, 3, \dots,$$

with method \hat{F}_n depending on a parameter n : Increasing n typically means investing more computational time for a hopefully more accurate result

Consistent: A numerical method is consistent if the local error made at each step vanishes for $n \rightarrow \infty$

$$\hat{F}_n(x_i, x_{i-1}) \rightarrow \bar{F}(x_i, x_{i-1}) \quad (n \rightarrow \infty)$$

This is one of the most basic requirements that we have on a numerical approach. If it is not consistent, it means we can invest more computational time (more effort) and certainly won't get lower errors.

Stability: Because we use \hat{F}_n instead of \bar{F} , in each iteration we make a local error (see above). We have \hat{x}_i at iteration i rather than x_i . Stability means here that the local error can be amplified only by a constant that is independent of n .

Convergence: If the numerical error can be made arbitrarily small by increasing the computational effort $n \rightarrow \infty$

consistency + stability \rightarrow convergence

A concrete and formal description of these concepts for finite difference approximations can be found in Chapter 2 of LeVeque's textbook on finite difference methods.

Speed of convergence

Let $x_n \rightarrow x$ in a normed space X , $\|\cdot\|$ for $n \rightarrow \infty$.

$$\lim_{n \rightarrow \infty} \frac{\|x - x_{n+1}\|}{\|x - x_n\|^q} < C$$

with $C > 0$ and $q \geq 1$

- ▶ Linear convergence: $q = 1$ and $C < 1$

$$\|x - x_{n+1}\| \leq C \|x - x_n\|$$

- ▶ Quadratic convergence: $q = 2$

$$\|x - x_{n+1}\| \leq C \|x - x_n\|^2$$

Beyond convergence

- ▶ An algorithm will produce the correct answer if it is convergent, but...
- ▶ Not all convergent methods are equal. We can differentiate them further based on:

Beyond convergence

- ▶ An algorithm will produce the correct answer if it is convergent, but...
- ▶ Not all convergent methods are equal. We can differentiate them further based on:
- ▶ **Accuracy** How much computational work do you need to expand to get an answer to a desired relative error?

Beyond convergence

- ▶ An algorithm will produce the correct answer if it is convergent, but...
- ▶ Not all convergent methods are equal. We can differentiate them further based on:
- ▶ **Accuracy** How much computational work do you need to expand to get an answer to a desired relative error?
- ▶ **Robustness** Does the algorithm work (equally) well for all (reasonable) input data d ?

Beyond convergence

- ▶ An algorithm will produce the correct answer if it is convergent, but...
- ▶ Not all convergent methods are equal. We can differentiate them further based on:
- ▶ **Accuracy** How much computational work do you need to expand to get an answer to a desired relative error?
- ▶ **Robustness** Does the algorithm work (equally) well for all (reasonable) input data d ?
- ▶ **Efficiency** How fast does the implementation produce the answer? This depends on the algorithm, on the computer, the programming language, the programmer, etc. (more next class)

Beyond convergence

- ▶ An algorithm will produce the correct answer if it is convergent, but...
- ▶ Not all convergent methods are equal. We can differentiate them further based on:
- ▶ **Accuracy** How much computational work do you need to expand to get an answer to a desired relative error?
- ▶ **Robustness** Does the algorithm work (equally) well for all (reasonable) input data d ?
- ▶ **Efficiency** How fast does the implementation produce the answer? This depends on the algorithm, on the computer, the programming language, the programmer, etc. (more next class)
- ▶ **Difficulty** How easy is it to implement and apply in practice? Do I need to spend 5 years of my time to implement it or can I code it up in 2 lines of code?

Matlab peculiarities [Following slides: A. Donev]

- MATLAB is an **interpreted language**, meaning that commands are interpreted and executed as encountered. MATLAB caches some stuff though...
- Many of MATLAB's **intrinsic routines** are however compiled and optimized and often based on well-known libraries (BLAS, LAPACK, FFTW, etc.).
- Variables in scripts/workspace are global and persist throughout an interactive session (use *whos* for info and *clear* to clear workspace).
- Every variable in MATLAB is, unless specifically arranged otherwise, a matrix, **double precision float** if numerical.
- Vectors (column or row) are also matrices for which one of the dimensions is 1.
- **Complex arithmetic** and complex matrices are used where necessary.

Matrices [Slide: A. Donev]

```
>> format compact; format long
>> x=-1; % A scalar that is really a 1x1 matrix
>> whos('x')
  Name      Size      Bytes  Class      Attributes
  x         1x1         8     double

```



```
>> y=sqrt(x) % Requires complex arithmetic
y =
      0 + 1.0000000000000000i
>> whos('y')
  Name      Size      Bytes  Class      Attributes
  y         1x1        16     double     complex

```



```
>> size(x)
ans =      1      1
>> x(1)
ans =     -1
>> x(1,1)
ans =     -1
>> x(3)=1;
>> x
x =     -1      0      1

```

Vectorization/Optimization [Slide: A. Donev]

- MATLAB uses **dynamic memory management** (including garbage collection), and matrices are re-allocated as needed when new elements are added.
- It is however much better to **pre-allocate space** ahead of time using, for example, *zeros*.
- The **colon notation** is very important in accessing array sections, and x is different from $x(:)$.
- **Avoid for loops** unless necessary: Use array notation and intrinsic functions instead.
- To see how much CPU (computing) time a section of code took, use *tic* and *toc* (but beware of timing small sections of code).
- MATLAB has built-in **profiling tools** (*help profile*).

Pre-allocation [Slide: A. Donev]

```
format compact; format long
clear; % Clear all variables from memory

N=100000; % The number of iterations

% Try commenting this line out:
f=zeros(1,N); % Pre-allocate f

tic;
f(1)=1;
for i=2:N
    f(i)=f(i-1)+i;
end
elapsed=toc;

fprintf('The result is f(%d)=%g, computed in %g s\n', ...
        N, f(N), elapsed);
```

Vectorization [Slide: A. Donev]

```
function vect(vectorize)
    N=1000000; % The number of elements
    x=linspace(0,1,N); % Grid of N equi-spaced points

    tic;
    if(vectorize) % Vectorized
        x=sqrt(x);
    else % Non-vectorized
        for i=1:N
            x(i)=sqrt(x(i));
        end
    end
    elapsed=toc;

    fprintf('CPU_time_for_N=%d_is_%g_s\n', N, elapsed);
end
```

Matlab examples [Slide: A. Donev]

```
>> fibb % Without pre-allocating
```

```
The result is  $f(100000)=5.00005e+09$ , computed in 6.53603 s
```

```
>> fibb % Pre-allocating
```

```
The result is  $f(100000)=5.00005e+09$ , computed in 0.000998 s
```

```
>> vect(0) % Non-vectorized
```

```
CPU time for N=1000000 is 0.074986 s
```

```
>> vect(1) % Vectorized — don't trust the actual number
```

```
CPU time for N=1000000 is 0.002058 s
```

Vectorization/Optimization [Slide: A. Donev]

- Recall that everything in MATLAB is a double-precision matrix, called **array**.
- Row vectors are just matrices with first dimension 1. Column vectors have row dimension 1. Scalars are 1×1 matrices.
- The syntax x' can be used to construct the **conjugate transpose** of a matrix.
- The **colon notation** can be used to select a subset of the elements of an array, called an **array section**.
- The default arithmetic operators, $+$, $-$, $*$, $/$ and \wedge are **matrix addition/subtraction/multiplication**, linear solver and matrix power.
- If you prepend a **dot before an operator** you get an **element-wise operator** which works for arrays of the same shape.

Matrices [Slide: A. Donev]

```
>> x=[1 2 3; 4 5 6] % Construct a matrix
x =
     1     2     3
     4     5     6
```

```
>> size(x) % Shape of the matrix x
ans =
     2     3
```

```
>> y=x(:) % All elements of y
y =
     1     4     2     5     3     6
```

```
>> size(y)
ans =
     6     1
```

```
>> x(1,1:3)
ans =
     1     2     3
```

```
>> x(1:2:6)
ans =
     1     2     3
```

Matrices [Slide: A. Donev]

```
>> sum(x)
```

```
ans =  
    5    7    9
```

```
>> sum(x(:))
```

```
ans =  
    21
```

```
>> z=1i; % Imaginary unit
```

```
>> y=x+z
```

```
y =  
    1.0000 + 1.0000i    2.0000 + 1.0000i    3.0000 + 1.0000i  
    4.0000 + 1.0000i    5.0000 + 1.0000i    6.0000 + 1.0000i
```

```
>> y'
```

```
ans =  
    1.0000 - 1.0000i    4.0000 - 1.0000i  
    2.0000 - 1.0000i    5.0000 - 1.0000i  
    3.0000 - 1.0000i    6.0000 - 1.0000i
```

Matrices [Slide: A. Donev]

```
>> x*y
??? Error using ==> mtimes
Inner matrix dimensions must agree.
```

```
>> x.*y
ans =
    1.0000 + 1.0000i    4.0000 + 2.0000i    9.0000 + 3.0000i
   16.0000 + 4.0000i   25.0000 + 5.0000i   36.0000 + 6.0000i
```

```
>> x*y'
ans =
   14.0000 - 6.0000i   32.0000 - 6.0000i
   32.0000 -15.0000i   77.0000 -15.0000i
```

```
>> x'*y
ans =
   17.0000 + 5.0000i   22.0000 + 5.0000i   27.0000 + 5.0000i
   22.0000 + 7.0000i   29.0000 + 7.0000i   36.0000 + 7.0000i
   27.0000 + 9.0000i   36.0000 + 9.0000i   45.0000 + 9.0000i
```

Coding guidelines [Slide: A. Donev]

- Learn to reference the **MATLAB help**: Including reading the examples and “fine print” near the end, not just the simple usage.
- **Indentation, comments, and variable naming** make a big difference! Code should be readable by others.
- Spending a few extra moments on the code will pay off when using it.
- Spend some time learning how to **plot in MATLAB**, and in particular, how to plot with different symbols, lines and colors using *plot*, *loglog*, *semilogx*, *semilogy*.
- Learn how to **annotate plots**: *xlim*, *ylim*, *axis*, *xlabel*, *title*, *legend*. The intrinsics *num2str* or *sprintf* can be used to create strings with embedded parameters.
- Finer controls over fonts, line widths, etc., are provided by the intrinsic function *set*...including using the LaTeX interpreter to typeset mathematical notation in figures.

Today

Last time

- ▶ Condition of problems
- ▶ Stability of algorithms

Today

- ▶ Float-point numbers in IEEE format
- ▶ Rounding, propagation of errors, and cancellation
- ▶ Truncation errors

Announcements

- ▶ Homework 1 was posted last week; is due next week Mon, Sep 23 *before class*

Recap: Condition of a problem

- ▶ Terms such as “little bit” and a “small amount” already point to that we need to measure something
- ▶ Therefore, we assume the map f is given as

$$f : U \subset \mathbb{R}^n \rightarrow \mathbb{R}^m$$

and we are interested in the norm $\| \cdot \|$

- ▶ The input error is then

$$\|x - \hat{x}\| \leq \delta \text{ (absolute)} \quad \|x - \hat{x}\| \leq \delta \|x\| \text{ (relative)}$$

- ▶ Correspondingly we measure the output error $f(x) - f(\hat{x})$ in $\| \cdot \|$ (we could also have looked at a componentwise error)

Recap: Condition of a problem (cont'd)

- ▶ Absolute condition number at x is

$$\kappa_{\text{abs}} = \lim_{\delta \rightarrow 0} \sup_{\|x - \hat{x}\| \leq \delta} \frac{\|f(x) - f(\hat{x})\|}{\|x - \hat{x}\|}$$

- ▶ Relative condition number at x is

$$\kappa_{\text{rel}} = \lim_{\delta \rightarrow 0} \sup_{\|x - \hat{x}\| \leq \delta} \frac{\|f(x) - f(\hat{x})\| / \|f(x)\|}{\|x - \hat{x}\| / \|x\|}$$

- ▶ If f is differentiable in x , then

$$\kappa_{\text{abs}} = \|f'(x)\| \quad \kappa_{\text{rel}} = \frac{\|x\|}{\|f(x)\|} \|f'(x)\|,$$

where $\|f'(x)\|$ is the norm of the Jacobian $f'(x)$ in the operator norm

$$\|A\| = \sup_{x \neq 0} \frac{\|Ax\|}{\|x\|} = \sup_{\|x\|=1} \|Ax\|$$

Recap: Condition of a problem (cont'd)

- ▶ If $\kappa_{\text{rel}} \sim 1$,

Recap: Condition of a problem (cont'd)

- ▶ If $\kappa_{\text{rel}} \sim 1$, then the problem is well conditioned: If the relative error in the data/input is small, then the relative error in the answer/output is similarly small
- ▶ If $\kappa_{\text{rel}} \gg 1$,

Recap: Condition of a problem (cont'd)

- ▶ If $\kappa_{\text{rel}} \sim 1$, then the problem is well conditioned: If the relative error in the data/input is small, then the relative error in the answer/output is similarly small
- ▶ If $\kappa_{\text{rel}} \gg 1$, then the problem is poorly conditioned: Small relative input error can lead to large relative output error
- ▶ If κ_{rel} (and κ_{abs}) do not exist, then the problem is ill conditioned.
- ▶ What is poorly conditioned depends on desired accuracy: if the input accuracy is low but we expect a high output accuracy, then problems are quickly poorly conditioned. If we are happy with a less accurate output, we might consider the problem still well conditioned.
- ▶ Sometimes, the possibly large error in the output does not matter and so we can solve poorly conditioned problems (think of early design stages, rapid prototyping, etc); but we should be very much aware of the condition of the problem.

Recap: Condition number of a matrix

Consider a matrix $A \in \mathbb{R}^{n \times n}$. Its condition number is

$$\kappa(A) = \|A\| \|A^{-1}\|$$

Widely used is the $\|\cdot\|_2$ norm and then

$$\kappa_2(A) = \|A\|_2 \|A^{-1}\|_2 = \frac{\sigma_{\max}(A)}{\sigma_{\min}(A)}$$

with the maximal and minimal singular value $\sigma_{\max}(A)$ and $\sigma_{\min}(A)$ of A

Consider a system of linear equations $Ax = b$. Then, the problems $A \mapsto A^{-1}b$ and $b \mapsto A^{-1}b$ have relative condition numbers

$$\kappa_{\text{rel}} \leq \kappa(A)$$

Recap: Backward stability

Backward stability: Pass the errors of the algorithm back and interpret as input errors.

An algorithm \tilde{f} for a problem f is backward stable if for each $x \in X$ we have $\tilde{f}(x) = f(\tilde{x})$ for an \tilde{x} with

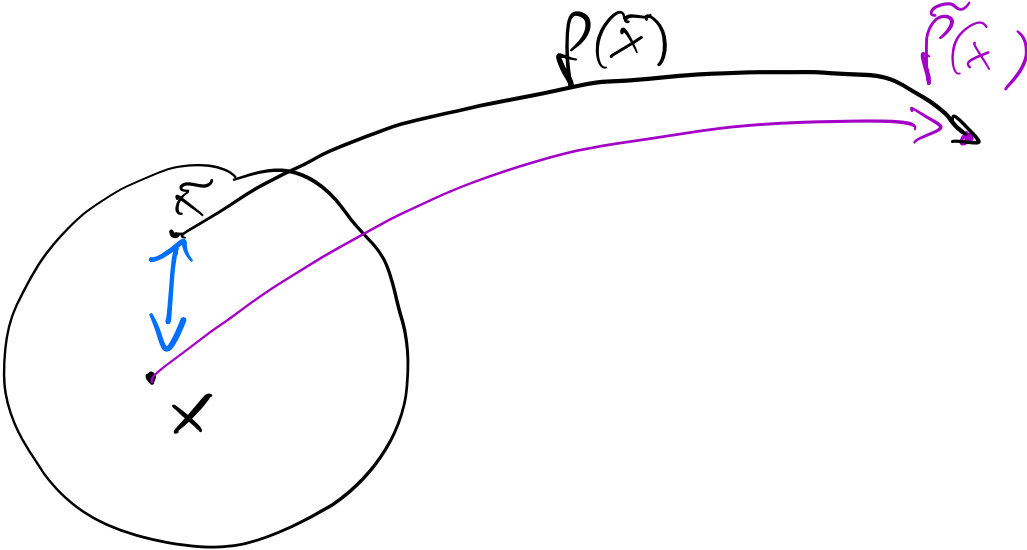
$$\frac{\|\tilde{x} - x\|}{\|x\|}$$

small

A backward stable algorithm gives exactly the right answer to nearly the right question.

In backward error analysis one calculates, for a given output, how much one would need to perturb the input in order for the answer to be exact.

Recap: Backward stability (cont'd)



Representing real numbers

Representing real numbers

- ▶ Computers represent everything using bit strings, i.e., integers in base 2. A finite number of integers can thus be exactly represented. But not real numbers! This leads to roundoff errors.
- ▶ Assume we have N digits to represent real numbers on a computer that can represent integers using a given number system, say decimal for human purposes.
- ▶ Fixed-point representation of numbers

$$x = (-1)^s \cdot [a_{N-2}a_{N-3} \cdots a_k \cdot a_{k-1} \cdots a_0]$$

has a problem of representing either small or larger numbers because the decimal point $.$ is fixed at position k

What could we do?

Floating-point numbers

- ▶ Instead, let's use floating-point representation

$$x = (-1)^s \cdot [0.a_1a_2 \cdots a_t] \cdot \beta^e = (-1)^s \cdot m \cdot \beta^{e-t}$$

similar to the common scientific number representation

$$0.1156 \cdot 10^1 = 1156 \cdot 10^{-3} \quad t = 4$$

- ▶ A floating-point number in base β is represented using one sign bit $s = 0$ or 1 , a t -digit integer mantissa

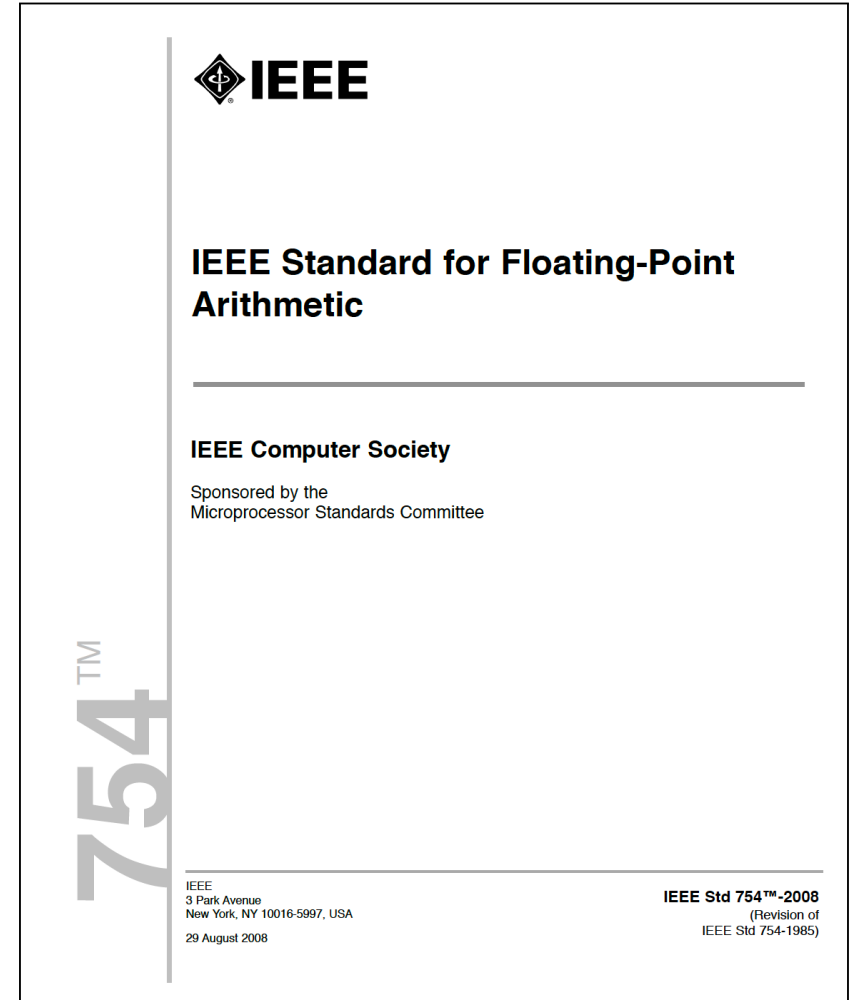
$$0 \leq m = [a_1a_2 \cdots a_t] \leq \beta^t - 1$$

and an integer exponent $L \leq e \leq U$

- ▶ Computers today use binary numbers and so $\beta = 2$

IEEE 754 standard

- ▶ Formats for representing and encoding real numbers using bit strings (single and double precision).
- ▶ Rounding algorithms for performing accurate arithmetic operations (e.g., addition, subtraction, division, multiplication) and conversions (e.g., single to double precision).
- ▶ Exception handling for special situations (e.g., division by zero and overflow).



Single precision IEEE floating-point numbers have the standardized storage format:

sign + power + fraction

with

$$N_s + N_p + N_f = 1 + 8 + 23 = 32 \text{ bits}$$

and are interpreted as

$$x = (-1)^s \cdot 2^{p-127} \cdot (1.f)_2$$

- ▶ Sign $s = 1$ for negative numbers
- ▶ Power $1 \leq p \leq 254$ determines the exponent
- ▶ Fractional part of the mantissa f
- ▶ `single` in Matlab, `float` in C/C++, `REAL` in Fortran

$$(138)_{10} = (10001010)_2$$

$$x = (-1)^0 2^{(10001010)_2 - 127} \cdot \underbrace{(1.01011)_2}_f$$

$$x = [s | p | f] = [0 | 10001010 | 010110 \dots 0]$$

$$= \underline{(452c000)}_{16}$$

IEEE representation example

Take the number $x = 2752 = 0.2752 \cdot 10^4$. Converting 2752 to the binary number system

$$\begin{aligned}x &= 2^{11} + 2^9 + 2^7 + 2^6 = (101011000000)_2 = 2^{11} \cdot (1.01011)_2 \\ &= (-1)^0 2^{138-127} \cdot (1.01011)_2 = (-1)^0 2^{(10001010)_2-127} \cdot (1.01011)_2\end{aligned}$$

On the computer:

$$\begin{aligned}x &= [s \quad | \quad p \quad | \quad f] \\ &= [0 \quad | \quad 100,0101,0 \quad | \quad 010,1100,0000,0000,0000,0000] \\ &= (452c0000)_{16}\end{aligned}$$

```
format hex;  
>> a=single(2.752E3)  
a =  
    452c0000
```

$$\begin{aligned}&(0100\ 0101\ 0010\ 1100\ 0000\ 0000\ 0000\ 0000)_2 \\ &=(452c0000)_{16}\end{aligned}$$

Double precision IEEE numbers

Double precision IEEE numbers (default in Matlab, `double` in C/C++) follow the same principle but use 64 bits to give higher precision and range

$$N_s + N_p + N_f = 1 + 11 + 52 = 64 \text{ bits}$$

$$x = (-1)^s \cdot 2^{p-1023} \cdot (1.f)_2$$

Even higher (extended) precision formats are not really standardized or widely implemented/used.

There is also software-emulated variable precision arithmetic in, e.g., Maple

Extremal exponent values

The extremal exponent values have special meaning (here single precision)

value	power p	fraction f
± 0	0	0
$\pm\infty$	255	0
Not a number (NaN)	255	> 0

Important facts about floating-point numbers

- ▶ Not all real numbers x can be represented exactly as a floating-point number. Instead, they must be *rounded* to the nearest floating point number $\hat{x} = \text{fl}(x)$
- ▶ Floating-point numbers have a *relative rounding error* that is smaller than the *machine precision* or *roundoff-unit* u

$$\frac{|\hat{x} - x|}{|x|} \leq u = 2^{-(N_f+1)} = \begin{cases} 2^{-24} \sim 6.0 \cdot 10^{-8}, & \text{for single precision} \\ 2^{-53} \sim 1.1 \cdot 10^{-16}, & \text{for double precision.} \end{cases}$$

- ▶ Often the machine precision/roundoff-unit is denoted as ϵ
- ▶ **The rule of thumb is that single precision gives 7-8 digits of precision and double 16 digits.**
- ▶ There is a smallest and largest possible number due to limit for the exponent.

Two axioms

Ignoring over- and underflow, we assume the following two “axioms” to hold for computers we work with:

1. For all $x \in \mathbb{R}$, there exists ϵ with $|\epsilon| \leq u$ (roundoff unit) such that

$$\text{fl}(x) = x(1 + \epsilon),$$

where $\text{fl}(\cdot)$ rounds to the the closest floating point approximation.

2. Consider two floating point numbers x, y . The floating-point operation \circledast (=add, sub, mult, div) of $*$ (=add, sub, mult, div) satisfies

$$x \circledast y = \text{fl}(x * y)$$

Axiom 1 and 2 imply that for two floating-point numbers x, y , there exists ϵ with $|\epsilon| \leq u$ such that

$$x \circledast y = (x * y)(1 + \epsilon).$$

Floating-point exceptions

Computing with floating point values may lead to exceptions, which may halt the program:

Floating-point exceptions

Computing with floating point values may lead to exceptions, which may halt the program:

- ▶ **Divide-by-zero:** if the result is $\pm\infty$, e.g., $1/0$

Floating-point exceptions

Computing with floating point values may lead to exceptions, which may halt the program:

- ▶ **Divide-by-zero:** if the result is $\pm\infty$, e.g., $1/0$
- ▶ **Invalid:** If the result is a *NaN*, e.g., taking $\sqrt{-1}$ (note that Matlab supports complex numbers...)

```
1: >>> x = math.sqrt(-1)
2: Traceback (most recent call last):
3:   File "<stdin>", line 1, in <module>
4: ValueError: math domain error
```

Floating-point exceptions

Computing with floating point values may lead to exceptions, which may halt the program:

- ▶ **Divide-by-zero:** if the result is $\pm\infty$, e.g., $1/0$
- ▶ **Invalid:** If the result is a *NaN*, e.g., taking $\sqrt{-1}$ (note that Matlab supports complex numbers...)

```
1: >>> x = math.sqrt(-1)
2: Traceback (most recent call last):
3:   File "<stdin>", line 1, in <module>
4: ValueError: math domain error
```

- ▶ **Overflow:** If the result is too large to be represented, e.g., adding two numbers, each on the order of *realmax*

Floating-point exceptions

Computing with floating point values may lead to exceptions, which may halt the program:

- ▶ **Divide-by-zero:** if the result is $\pm\infty$, e.g., $1/0$
- ▶ **Invalid:** If the result is a *NaN*, e.g., taking $\sqrt{-1}$ (note that Matlab supports complex numbers...)

```
1: >>> x = math.sqrt(-1)
2: Traceback (most recent call last):
3:   File "<stdin>", line 1, in <module>
4: ValueError: math domain error
```

- ▶ **Overflow:** If the result is too large to be represented, e.g., adding two numbers, each on the order of *realmax*
- ▶ **Underflow:** If the result is too small to be represented, e.g., dividing a number close to *realmin* by a large number.

Avoiding overflow

Numerical software needs to be careful about avoiding exceptions:

Mathematically equivalent expressions are not necessarily computationally equivalent!

- ▶ For example, computing $\sqrt{x^2 + y^2}$ may lead to overflow in computing $x^2 + y^2$ even though the result does not overflow
- ▶ Matlab's hypot function guards against this:

$$\sqrt{x^2 + y^2} = |x| \sqrt{1 + \left(\frac{y}{x}\right)^2} \text{ ensuring that } |x| > |y|$$

works correctly

- ▶ These kind of careful constructions may have higher computational cost (more CPU operations) or make roundoff errors worse.

Floating-point in practice

- Most scientific software **uses double precision** to avoid range and accuracy issues with single precision (better be safe than sorry). Single precision may offer speed/memory/vectorization advantages however (e.g. GPU computing).
- **Do not compare floating point numbers** (especially for loop termination), or more generally, do not rely on logic from pure mathematics.
- Optimization, especially in compiled languages, can rearrange terms or perform operations using **unpredictable** alternate forms (e.g., wider internal registers).
Using parenthesis helps , e.g. $(x + y) - z$ instead of $x + y - z$, but does not eliminate the problem.
- Library functions such as \sin and \ln will typically be computed almost to full machine accuracy, but do not rely on that for special/complex functions.

Propagation of errors

- ▶ Assume that we are calculating something with numbers that are not exact, e.g., a rounded floating-point number \hat{x} versus the exact real number x .
- ▶ For IEEE representations, recall that

$$\frac{|\hat{x} - x|}{|x|} \leq u = 2^{-(N_f+1)} = \begin{cases} 2^{-24} \sim 6.0 \cdot 10^{-8}, & \text{for single precision} \\ 2^{-53} \sim 1.1 \cdot 10^{-16}, & \text{for double precision.} \end{cases}$$

- ▶ In general, the *absolute error* $\delta x = \hat{x} - x$ may have contributions from each of the different types of error (roundoff, truncation, propagated, statistical).
- ▶ Assume we have an estimate or bound for the relative error

$$\left| \frac{\delta x}{x} \right| \lesssim \epsilon_x \ll 1$$

based on some analysis, e.g., for roundoff error the IEEE standard determines $\epsilon_x = u$ (roundoff-unit)

Propagation of errors

How does the relative error change (propagate) during numerical calculations?

Propagation of errors

multiplication / division

$$\begin{aligned} E_{xy} &= \left| \frac{(x + \delta x)(y + \delta y) - xy}{xy} \right| \\ &= \left| \frac{\cancel{xy} + \delta xy + \cancel{xy} + \frac{\delta x \delta y}{xy} - \cancel{xy}}{\cancel{xy}} \right| \\ &= \left| \frac{\delta x}{x} + \frac{\delta y}{y} + \frac{\delta x \delta y}{xy} \right| \end{aligned}$$

Now use

$$\left| \frac{\delta x}{x} \right| \ll 1, \quad \left| \frac{\delta y}{y} \right| \ll 1$$

$$\left| \frac{\delta x \delta y}{xy} \right| \leq \max \left\{ \left| \frac{\delta x}{x} \right|, \left| \frac{\delta y}{y} \right| \right\}$$

Then

$$\left| \frac{\delta x}{x} + \frac{\delta y}{y} + \frac{\delta x \delta y}{xy} \right| \leq 2(\epsilon_x + \epsilon_y)$$

$$\begin{aligned} &\lesssim \epsilon_x + \epsilon_y \\ &\quad \uparrow \text{up to constant} \end{aligned}$$

\Rightarrow rel. errors are added

\Rightarrow safe because accurate input gives output of similar accuracy

Addition / subtraction

this is more dangerous

\Rightarrow lead to catastrophic errors

Example

$$x = 1, \quad y = 0.999, \quad \delta_x = 10^{-6}, \quad \delta_y = 9 \times 10^{-7}$$

$$\left| \frac{\delta_x}{x} \right| = \underline{10^{-6}}, \quad \left| \frac{\delta_y}{y} \right| = \left| \frac{9 \times 10^{-7}}{9.99 \times 10^{-1}} \right| \approx \underline{10^{-6}}$$

$$\left| \frac{(x + \delta_x) - (y + \delta_y) - (x - y)}{x - y} \right|$$

$$= \left| \frac{\delta_x - \delta_y}{x - y} \right| = \frac{10^{-7}}{10^{-3}} = 10^{-4} \gg 10^{-6}$$

\Rightarrow errors can propagate quickly in +/-

Propagation of errors: Numerical experiment [From A. Donev]

Harmonic sum

$$H(N) = \sum_{i=1}^N \frac{1}{i}$$

```
1: function nhsum = harmonic(N)
2: nhsum = 0;
3: for i = 1:N
4:     nhsum = nhsum + 1.0/i;
5: end
6: end
```

What are the numerical issues of this implementation?

Propagation of errors: Numerical experiment [From A. Donev]

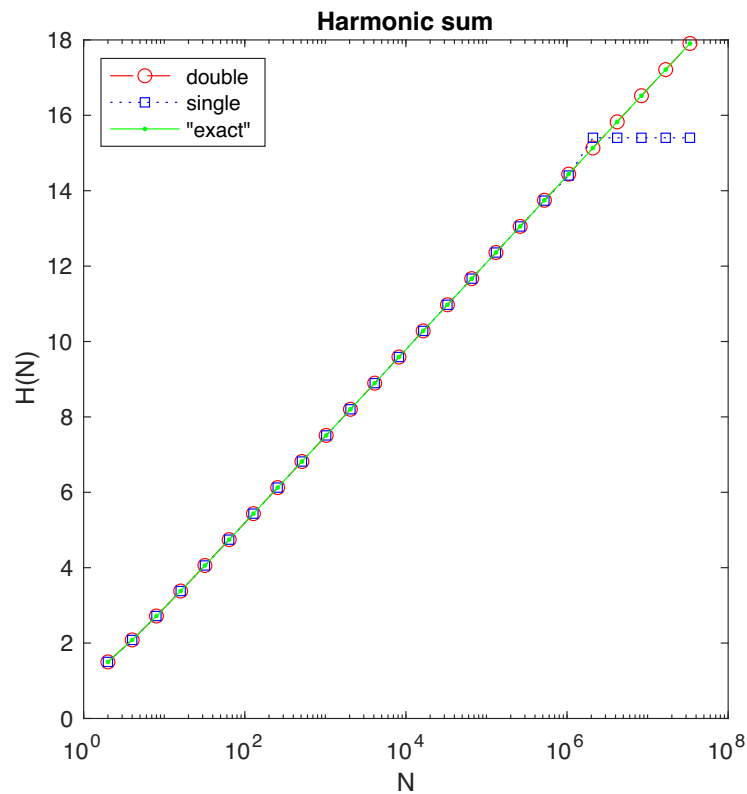
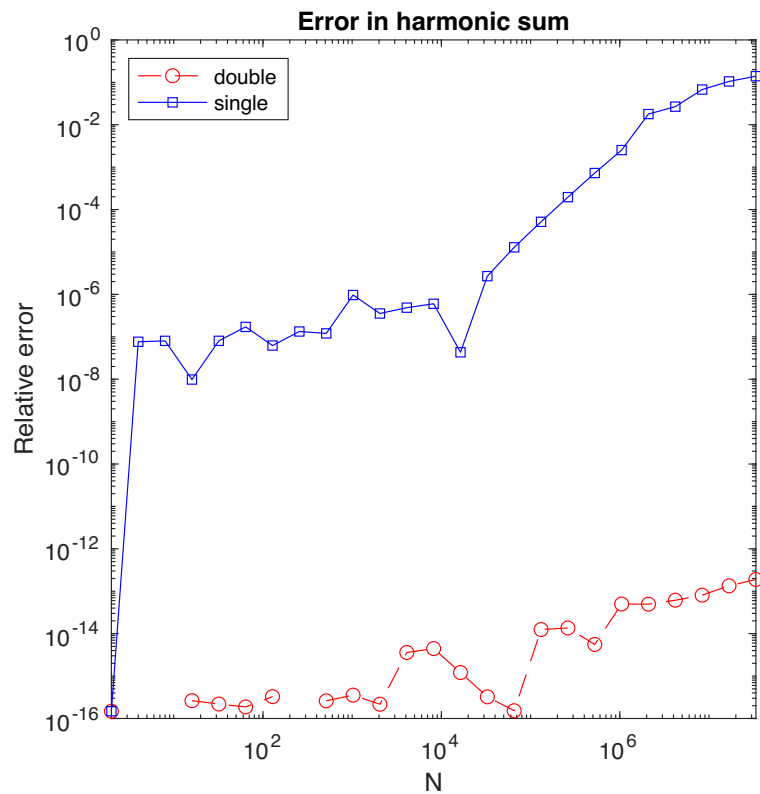
Harmonic sum

$$H(N) = \sum_{i=1}^N \frac{1}{i}$$

```
1: function nhsum = harmonic(N)
2: nhsum = 0;
3: for i = 1:N
4:     nhsum = nhsum + 1.0/i;
5: end
6: end
```

What are the numerical issues of this implementation?

↪ Adds very small number $1/i$ to potentially large number `nhsum`

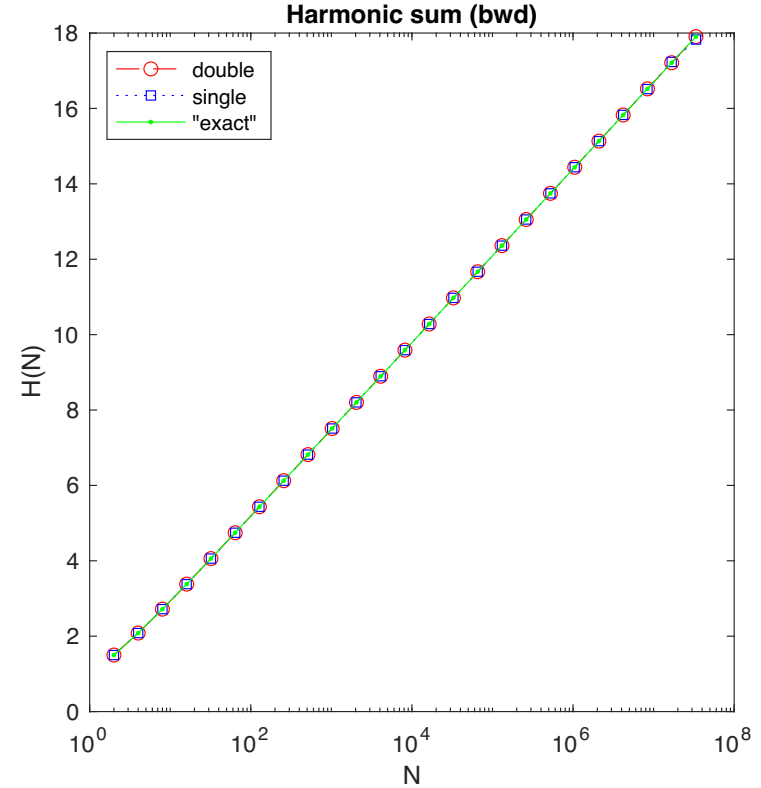
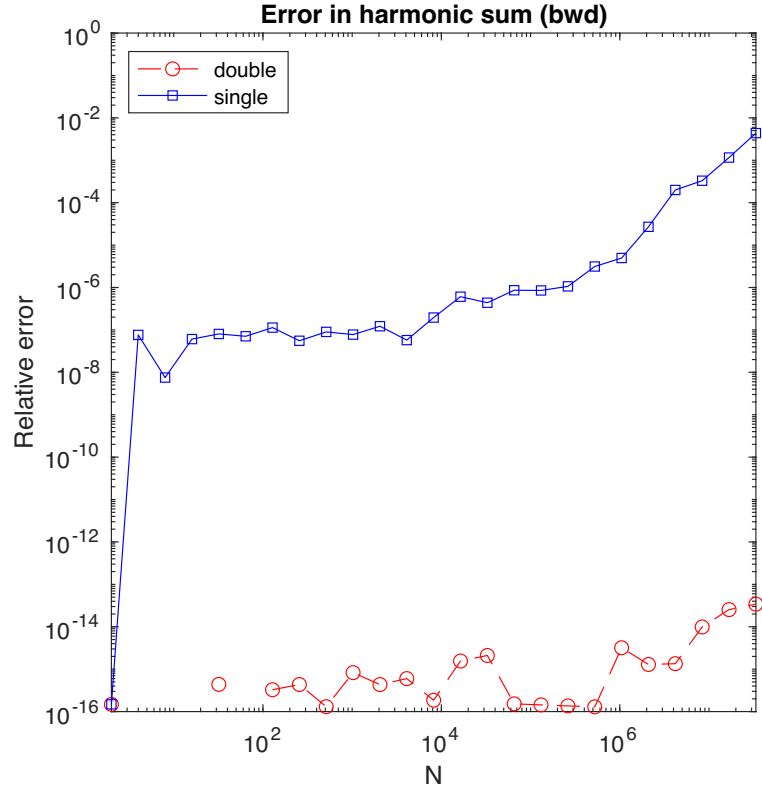


What can we do about it?

Implementation with backward summation

```
1: function nhsum = harmonicBwd(N)
2: nhsum = 0;
3: for i = N:-1:1
4:     nhsum = nhsum + 1.0/i;
5: end
6: end
```

Better, because adds small numbers to small numbers and larger numbers to large numbers.



Numerical Methods I

MATH-GA 2010.001/CSCI-GA 2420.001

Benjamin Peherstorfer
Courant Institute, NYU

Based on slides by G. Stadler and A. Donev

Today

Last time

- ▶ Float-point numbers in IEEE format
- ▶ Rounding
- ▶ Propagation of errors

Today

- ▶ Cancellation
- ▶ Truncation errors
- ▶ Solving linear systems

Announcements

- ▶ Homework 1 was posted last week; is due next week Mon, Sep 23 *before class*

→ Groder office hour: Fri, 3-4pm, WWH 303

Recap: Floating-point numbers

- ▶ Let's use floating-point representation

$$x = (-1)^s \cdot [0.a_1a_2 \cdots a_t] \cdot \beta^e = (-1)^s \cdot m \cdot \beta^{e-t}$$

similar to the common scientific number representation

$$0.1156 \cdot 10^1 = 1156 \cdot 10^{-3} \quad t = 4$$

- ▶ A floating-point number in base β is represented using one sign bit $s = 0$ or 1 , a t -digit integer mantissa

$$0 \leq m = [a_1a_2 \cdots a_t] \leq \beta^t - 1$$

and an integer exponent $L \leq e \leq U$

- ▶ Computers today use binary numbers and so $\beta = 2$

Recap: Single precision

Single precision IEEE floating-point numbers have the standardized storage format:

sign + power + fraction

with

$$N_s + N_p + N_f = 1 + 8 + 23 = 32 \text{ bits}$$

and are interpreted as

$$x = (-1)^s \cdot 2^{p-127} \cdot (1.f)_2$$

- ▶ Sign $s = 1$ for negative numbers
- ▶ Power $1 \leq p \leq 254$ determines the exponent
- ▶ Fractional part of the mantissa f
- ▶ `single` in Matlab, `float` in C/C++, `REAL` in Fortran

Recap: Important facts about floating-point numbers

- ▶ Not all real numbers x can be represented exactly as a floating-point number. Instead, they must be *rounded* to the nearest floating point number $\hat{x} = \text{fl}(x)$
- ▶ Floating-point numbers have a *relative rounding error* that is smaller than the *machine precision* or *roundoff-unit* u

$$\frac{|\hat{x} - x|}{|x|} \leq u = 2^{-(N_f+1)} = \begin{cases} 2^{-24} \sim 6.0 \cdot 10^{-8}, & \text{for single precision} \\ 2^{-53} \sim 1.1 \cdot 10^{-16}, & \text{for double precision.} \end{cases}$$

- ▶ Often the machine precision/roundoff-unit is denoted as ϵ
- ▶ **The rule of thumb is that single precision gives 7-8 digits of precision and double 16 digits.**
- ▶ There is a smallest and largest possible number due to limit for the exponent.

Recap: Two axioms

Ignoring over- and underflow, we assume the following two “axioms” to hold for computers we work with:

1. For all $x \in \mathbb{R}$, there exists ϵ with $|\epsilon| \leq u$ (roundoff unit) such that

$$\text{fl}(x) = x(1 + \epsilon),$$

where $\text{fl}(\cdot)$ rounds to the the closest floating point approximation.

2. Consider two floating point numbers x, y . The floating-point operation \circledast (=add, sub, mult, div) of $*$ (=add, sub, mult, div) satisfies

$$x \circledast y = \text{fl}(x * y)$$

Axiom 1 and 2 imply that for two floating-point numbers x, y , there exists ϵ with $|\epsilon| \leq u$ such that

$$x \circledast y = (x * y)(1 + \epsilon).$$

Numerical cancellation

If x and y are close to each other, then $x - y$ can have reduced accuracy due to catastrophic cancellation.

Consider computing the smaller root of the quadratic equation

$$x^2 - 2x + c = 0$$

for $|c| \ll 1$ and focus on propagation/accumulation of roundoff errors.

$$x^2 - 2x + c = 0, \quad |c| \ll 1$$

Solution: $x = 1 - \sqrt{1-c}$

Before we start: Taylor

$$f(z) = \sqrt{1+z}, \quad f'(z) = \frac{1}{2\sqrt{1+z}}$$

$$f(z) = f(0) + \frac{f'(0)}{1!} (z-0) + \frac{f''(0)}{2!} (z-0)^2 + \dots$$

truncate: $\sqrt{1+z} \doteq 1 + \frac{1}{2\sqrt{1+0}} z = 1 + \frac{z}{2}$

Thus, for $|c|$ small, get

$$x = 1 - \sqrt{1-c} \doteq 1 - (1 - \frac{c}{2}) = \frac{c}{2}$$

Case 1: $|c| < |u| \ll \text{round-off unit}$

$$f(1-c) = 1$$

no point in moving forward

Case 2: $|u| < |c| \ll 1$

Calculate $1-c$ in floating point

$$f(1) \ominus f(c) = f(1 - f(c))$$

$$= (f(1) - f(c))(1 + \epsilon) \quad |\epsilon| \leq u$$

$$= \underbrace{(1(1 + \epsilon_1) - c(1 + \epsilon_c))}_{(1-c) + (\epsilon_1 - c\epsilon_c)} (1 + \epsilon), \quad |\epsilon_1| \leq u, |\epsilon_c| \leq u$$

$$f(1) \ominus f(c) - (1-c) = \underbrace{(1-c)}_{\text{order 1}} \epsilon + \underbrace{(\epsilon_1 - c\epsilon_c)}_{\text{order } |u|} \underbrace{(1+\epsilon)}_{\text{order } |u|}$$

$\left| \frac{\delta(1-c)}{1-c} \right|$ also order $|c|$ because $1-c$ close to 1

Assume $\text{sqrt}(c)$ computes root to within relative accuracy of round-off unit ν

$$\begin{aligned}\sqrt{x+\delta x} &= \sqrt{x} \left(1 + \frac{\delta x}{x}\right)^{\frac{1}{2}} \\ &\approx \sqrt{x} \left(1 + \frac{\delta x}{2x}\right)\end{aligned}$$

So far shows that $\sqrt{1-c}$ has abs and relative error of order $|c|$

We have $x = 1 - \sqrt{1-c}$

\ominus (ready know

$$f(1) \ominus f(y) - (1-y) = \delta(1-y)$$

is of order $|1-y| \cdot |c| + |c|$

Relative error:

$$\left| \frac{\delta(1-y)}{1-y} \right| \approx \frac{|1-y| \cdot |c| + |c|}{|1-y|} = |c| + \frac{|c|}{|1-y|}$$

It is higher than $|c|$ for $1-y \approx 0$

To avoid cancellation, we should not directly implement $1 - \sqrt{1 - c}$

Rather, we can take the Taylor approximate $x \approx \frac{c}{2}$, which provides a good approximation for small c .

Even better, we could use the **mathematically equivalent but numerically preferred form**:

$$1 - \sqrt{1 - c} = \frac{c}{1 + \sqrt{1 - c}}$$

which does not suffer from cancellation problems as c becomes smaller.

(Notice that $1 - \sqrt{\dots}$ is avoided and therefore the cancellation problem shown by our analysis is avoided. We showed that $\sqrt{1 - c}$ is safe.)

To avoid cancellation, we should not directly implement $1 - \sqrt{1 - c}$

Rather, we can take the Taylor approximate $x \approx \frac{c}{2}$, which provides a good approximation for small c .

Even better, we could use the **mathematically equivalent but numerically preferred form**:

$$1 - \sqrt{1 - c} = \frac{c}{1 + \sqrt{1 - c}}$$

which does not suffer from cancellation problems as c becomes smaller.

(Notice that $1 - \sqrt{\dots}$ is avoided and therefore the cancellation problem shown by our analysis is avoided. We showed that $\sqrt{1 - c}$ is safe.)

```
1: >>> c = 1e-10 # solution roughly 5.000000000125 x 10^-11
2: >>> 1 - math.sqrt(1 - c)
3: 5.000000413701855e-11
4: >>> c/(1 + math.sqrt(1 - c))
5: 5.000000000125e-11
```

Big \mathcal{O} notation

Useful to compare growth of functions.

We write $f \in \mathcal{O}(g)(x \rightarrow \infty)$ if there exists constant $C > 0$ such that for an x_0 the following holds

$$\forall x \geq x_0 : |f(x)| \leq C|g(x)|$$

We also write $f \in \mathcal{O}(g)(x \rightarrow 0)$ if there exists a constant $C > 0$ such that for an $x_0 > 0$ the following holds

$$\forall |x| \leq x_0 : |f(x)| \leq C|g(x)|$$

In many cases we do not write explicitly whether we mean $x \rightarrow \infty$ or $x \rightarrow 0$ because it is clear from the context.

Question: In practice, would you prefer an algorithm with costs growing as $c_1(x) \in \mathcal{O}(x)$ or $c_2(x) \in \mathcal{O}(x^2)$? Why?

Question: In practice, would you prefer an algorithm with costs growing as $c_1(x) \in \mathcal{O}(x)$ or $c_2(x) \in \mathcal{O}(x^2)$? Why?

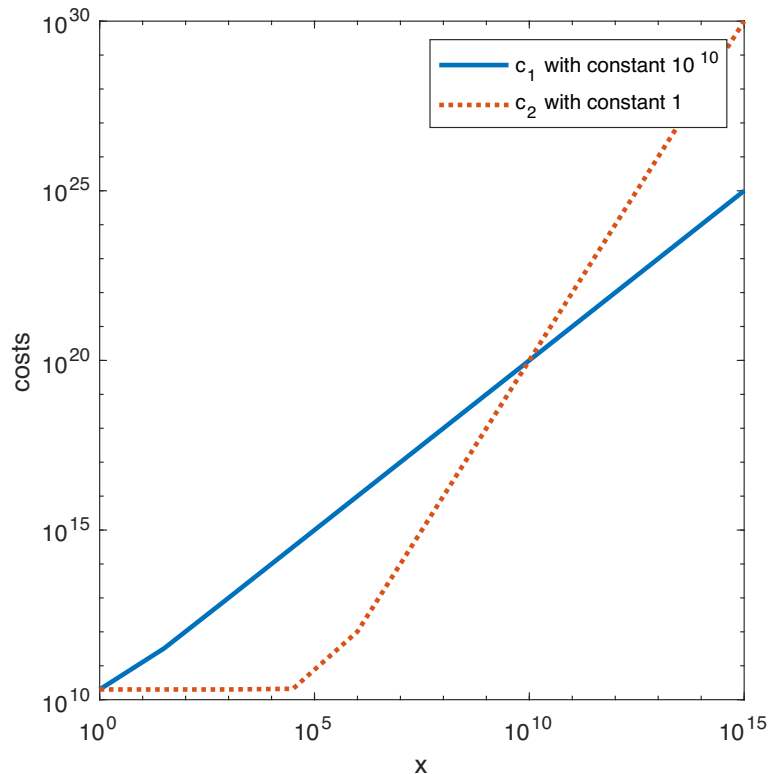
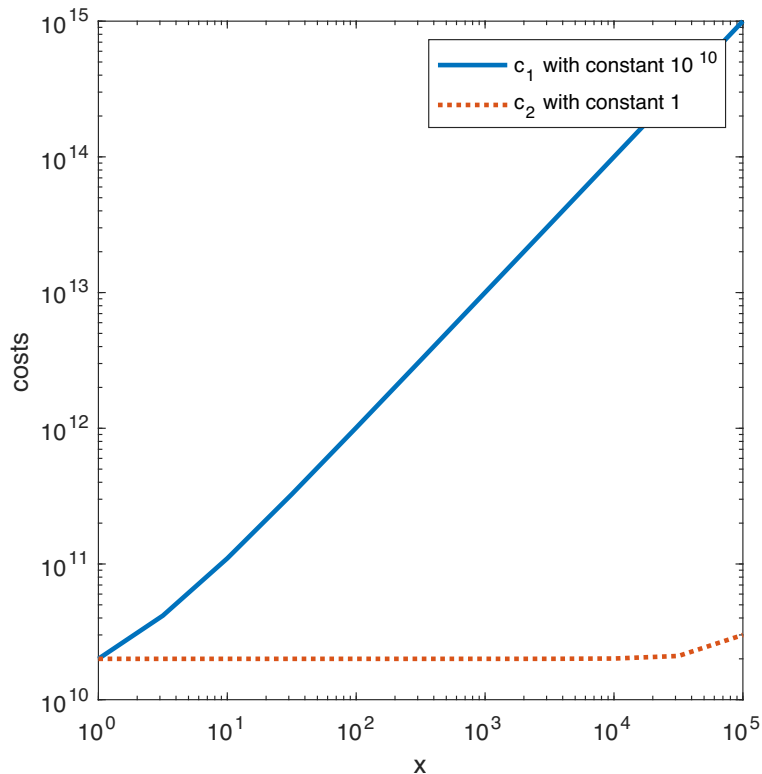
Answer: It depends on the hidden constants C and x_0 . If c_1 and c_2 have roughly the same constants, then probably c_1 .

However, if the constant for c_1 is $x_0 = 10^{10}$ and the constant for c_2 is $x_0 = 1$, then in most practical situations we prefer c_2 because we most likely will never reach the asymptotics of $x > x_0$ for c_1 in practice!

Warning: The Big \mathcal{O} notation tells us something about the asymptotics. The constants x_0 and C that are hidden in $\mathcal{O}(\cdot)$ do matter in practice!

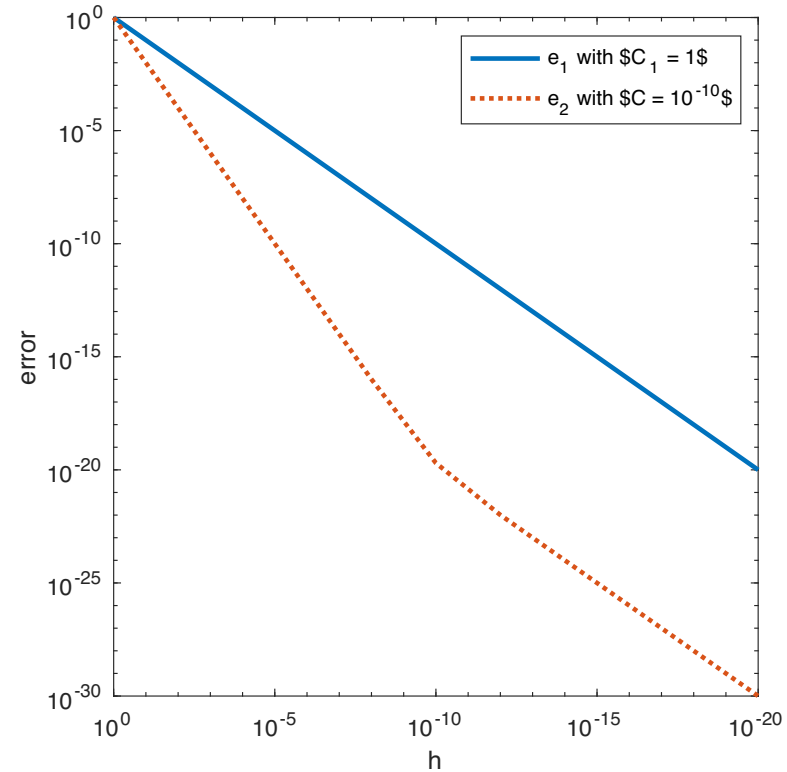
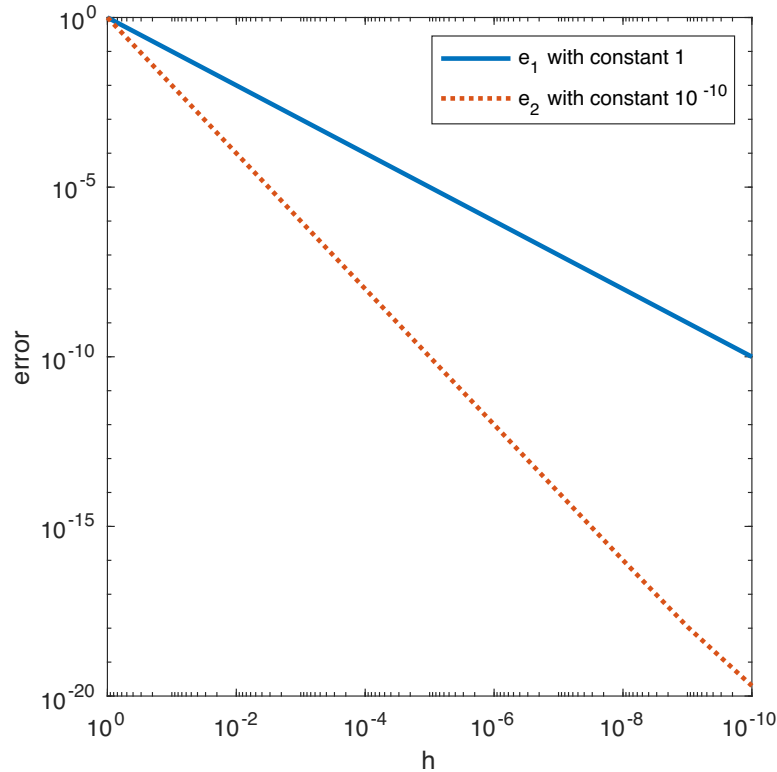
Costs (i.e., $x \rightarrow \infty$) of $c_1(x) = 10^{10} + 10^{10}x$ and $c_2(x) = 2 \times 10^{10} + x^2$. Then, $c_1 \in \mathcal{O}(x)$ and $c_2 \in \mathcal{O}(x^2)$ for $x \rightarrow \infty$. I.e., asymptotically the costs of c_2 grow faster than c_1 .

Costs (i.e., $x \rightarrow \infty$) of $c_1(x) = 10^{10} + 10^{10}x$ and $c_2(x) = 2 \times 10^{10} + x^2$. Then, $c_1 \in \mathcal{O}(x)$ and $c_2 \in \mathcal{O}(x^2)$ for $x \rightarrow \infty$. I.e., asymptotically the costs of c_2 grow faster than c_1 .



Warning: The Big \mathcal{O} notation tells us something about the asymptotics. The constants x_0 and C that are hidden in $\mathcal{O}(\cdot)$ do matter in practice!

Set now error: $e_1(h) = h$ and $e_2(h) = 10^{-10}h + h^2$. Then, $e_1 \in \mathcal{O}(h)$ and $e_2 \in \mathcal{O}(h)$ for $h \rightarrow 0$.



Warning: The Big \mathcal{O} notation tells us something about the asymptotics. The constants h_0 and C that are hidden in $\mathcal{O}(\cdot)$ do matter in practice!

Revisiting stability

Recall that we said: An algorithm \tilde{f} for a problem f is backward stable if for each $x \in X$ we have $\tilde{f}(x) = f(\tilde{x})$ for an \tilde{x} with

$$\frac{\|\tilde{x} - x\|}{\|x\|}$$

small.

We now can be more precise: An algorithm \tilde{f} for a problem f is backward stable if for each $x \in X$ we have $\tilde{f}(x) = f(\tilde{x})$ for an \tilde{x} with

$$\frac{\|\tilde{x} - x\|}{\|x\|} \in \mathcal{O}(u),$$

where u is the roundoff unit

- ▶ Recall that, loosely speaking, the symbol $\mathcal{O}(u)$ means “on the order of the roundoff unit.”
- ▶ By allowing $u \rightarrow 0$ (which is implied here by the \mathcal{O}), we consider an idealization of a computer (in practice, u is fixed). So what we mean is that the error should decrease in proportion to u or faster.

Suppose a backward stable algorithm is applied to solve a problem $f : X \rightarrow Y$ with relative condition number κ . Then, the relative errors satisfy

$$\frac{\|\tilde{f}(x) - f(x)\|}{\|f(x)\|} \in \mathcal{O}(\kappa(x)u).$$

Proof board

Let \tilde{f} be a backward stable algorithm of $f: X \rightarrow Y$ with rel condition κ , then

$$\frac{\| \tilde{f}(x) - f(x) \|}{\| f(x) \|} \in \mathcal{O}(\kappa(x) \nu)$$

where ν is the round off unit.

Backward stability tells us $\tilde{f}(x) = f(\tilde{x})$ for $\tilde{x} \in X$ satisfying

$$\frac{\| \tilde{x} - x \|}{\| x \|} \in \mathcal{O}(\nu)$$

Condition number

$$\kappa = \lim_{\delta \rightarrow 0} \sup_{\|x - \tilde{x}\| \leq \delta} \frac{\|x - \tilde{x}\|}{\|x\|}$$

$$\frac{\| f(x) - f(\tilde{x}) \|}{\| f(x) \|} \bigg/ \frac{\| x - \tilde{x} \|}{\| x \|}$$

thus with $\tilde{f}(x) = f(\tilde{x})$ we have

$$\frac{\| f(x) - \tilde{f}(x) \|}{\| f(x) \|} \leq \kappa(x) \frac{\| x - \tilde{x} \|}{\| x \|} \quad \nu \rightarrow 0$$

$$\in \mathcal{O}(\kappa(x) \nu)$$

Local truncation error

- ▶ **Approximation error** comes about when we replace a mathematical problem with some easier to solve approximation.
- ▶ This error is separate from and in addition to any numerical algorithm or computation used to actually solve the approximation itself, such as roundoff or propagated error.
- ▶ Truncation error is a common type of approximation error that comes from replacing infinitesimally small quantities with finite step size and truncating infinite sequences/series with finite ones.
- ▶ This is the most important type of error in methods for numerical interpolation, integration, solving differential equations, and others.

Local truncation error (cont'd)

- ▶ Analysis of local truncation error is almost always based on using Taylor series to approximate a function about a given point x

$$f(x + h) = \sum_{n=0}^{\infty} \frac{h^n}{n!} f^{(n)}(x) = f(x) + hf'(x) + \frac{h^2}{2} f''(x) + \dots ,$$

where we call h the step size

- ▶ We cannot do a series (infinite number of terms) numerically, so we truncate

$$f(x + h) \approx F_p(x, h) = \sum_{n=0}^p \frac{h^n}{n!} f^{(n)}(x)$$

- ▶ Question: What is the truncation error in this approximation? \rightsquigarrow **This kind of error estimate is one of the most commonly used in numerical analysis.**

- ▶ The remainder theorem of calculus provides a formula for the error: If the derivatives of f up to order $p + 1$ exist and are continuous in the interval $(x, x + h)$, then there is a $\xi \in [x, x + h]$ so that

$$f(x + h) - F_p(x, h) = \frac{h^{p+1}}{(p + 1)!} f^{(p+1)}(\xi)$$

- ▶ If we set

$$C = \frac{1}{(p + 1)!} \max_{y \in [x, x+h]} |f^{(p+1)}(y)|$$

then

$$|f(x + h) - F_p(x, h)| \leq Ch^{p+1}$$

- ▶ Intuition: If h is small and $f^{(p+1)}$ smooth, then $x \approx \arg \max_{y \in [x, x+h]} |f^{(p+1)}(y)|$ and the remainder term is nearly equal to the first neglected term

$$f(x + h) - F_p(x, h) \approx \frac{h^{p+1}}{(p + 1)!} f^{(p+1)}(x)$$

For example, let $e(h) = |f(x + h) - F_p(x, h)|$ and let the remainder theorem from the previous slide apply, then

$$e(h) \in \mathcal{O}(h^{p+1}), \quad h \rightarrow 0$$

Conclusions and summary

- ▶ No numerical method can compensate for a poorly conditioned problem. But not every numerical method will be a good one for a well conditioned problem.
- ▶ A numerical method needs to control the various computational errors (approximation, truncation, roundoff, propagated, statistical) while balancing computational cost.
- ▶ A numerical method must be consistent and stable in order to converge to the correct answer.
- ▶ The IEEE standard standardizes the single and double precision floating-point formats, their arithmetic, and exceptions. It is widely implemented.
- ▶ Numerical overflow, underflow and cancellation need to be carefully considered and avoided: Mathematically equivalent forms are not numerically equivalent.

Linear systems

Linear systems of equations

It is said that 70-80% of computational mathematics research involves solving systems of m linear equations in n unknowns

$$\sum_{j=1}^n a_{ij}x_j = b_i, \quad i = 1, \dots, m.$$

Linear systems arise directly from discrete models (e.g., in machine learning). Or through representing some abstract linear operator (such as a differential operator) in a finite basis as when numerically solving partial differential equations.

The common abstract way of writing systems of linear equations is

$$\mathbf{Ax} = \mathbf{b},$$

with matrix $\mathbf{A} \in \mathbb{R}^{m \times n}$, right-hand side $\mathbf{b} \in \mathbb{R}^m$, and solution $\mathbf{x} \in \mathbb{R}^n$

The goal is to calculate solution \mathbf{x} given data \mathbf{A}, \mathbf{b} in a numerically stable and computationally efficient way.

The matrix inverse

- ▶ A square matrix $\mathbf{A} \in \mathbb{R}^{n \times n}$ is invertible or nonsingular if there exists a matrix inverse $\mathbf{A}^{-1} = \mathbf{B} \in \mathbb{R}^{n \times n}$ such that

$$\mathbf{AB} = \mathbf{BA} = \mathbf{I},$$

where \mathbf{I} is the identity matrix.

- ▶ Matrix norm induced by a given vector norm

$$\|\mathbf{A}\| = \sup_{\mathbf{x} \neq \mathbf{0}} \frac{\|\mathbf{Ax}\|}{\|\mathbf{x}\|} \implies \|\mathbf{Ax}\| \leq \|\mathbf{A}\| \|\mathbf{x}\|$$

with sub-multiplicativity: $\|\mathbf{AB}\| \leq \|\mathbf{A}\| \|\mathbf{B}\|$

- ▶ Special case of interest: The 2-norm or spectral norm: $\|\mathbf{A}\|_2 = \sigma_1$ (largest singular value)
- ▶ The Euclidean or Frobenius norm is *not* an induced norm

$$\|\mathbf{A}\|_F = \sqrt{\sum_{i,j} |a_{ij}|^2}$$

but still is sub-multiplicative: $\|\mathbf{AB}\|_F \leq \|\mathbf{A}\|_F \|\mathbf{B}\|_F$

Condition of solving system of linear equations

Recall that we derived the condition number $\kappa(\mathbf{A}) = \|\mathbf{A}\| \|\mathbf{A}^{-1}\|$ of a matrix \mathbf{A}

Condition of solving system of linear equations (cont'd)

Now consider the general perturbations of the data

$$(\mathbf{A} + \delta\mathbf{A})(\mathbf{x} + \delta\mathbf{x}) = \mathbf{b} + \delta\mathbf{b}$$

One obtains the condition (proof in Quarteroni et al., Sec. 3.1)

$$\frac{\|\delta\mathbf{x}\|}{\|\mathbf{x}\|} \leq \frac{\kappa(\mathbf{A})}{1 - \kappa(\mathbf{A}) \frac{\|\delta\mathbf{A}\|}{\|\mathbf{A}\|}} \left(\frac{\|\delta\mathbf{b}\|}{\|\mathbf{b}\|} + \frac{\|\delta\mathbf{A}\|}{\|\mathbf{A}\|} \right)$$

Important practical estimate: Roundoff error in the data, with rounding unit u (recall $\approx 10^{-16}$ for double precision), produces a relative error

$$\frac{\|\delta\mathbf{x}\|_\infty}{\|\mathbf{x}\|_\infty} \leq \frac{2}{1 - \kappa(\mathbf{A})u} u\kappa(\mathbf{A})$$

\implies makes *no* sense to try to numerically solve systems with $\kappa(\mathbf{A}) > 10^{16}$ in double precision

Numerical solution of linear systems

There are many numerical methods for solving a system of linear equations

The most appropriate method depends on the properties of \mathbf{A}

- ▶ General dense matrices, where the entries in \mathbf{A} are mostly non-zero and nothing special is known \rightsquigarrow we focus on Gaussian elimination today
- ▶ General sparse matrices, where only a small fraction of $a_{ij} \neq 0$ (sparse typically means that $\mathcal{O}(n)$ entries are non-zero in an $n \times n$ matrix)
- ▶ Symmetric and positive-definite matrices
- ▶ Special structured sparse matrices, often arising from specific physical properties of the underlying system

It is also important to consider how many times a linear system with the same or related matrix or right-hand side needs to be solved.

Numerical Methods I

MATH-GA 2010.001/CSCI-GA 2420.001

Benjamin Peherstorfer
Courant Institute, NYU

Based on slides by G. Stadler and A. Donev

Today

Last time

- ▶ Cancellation
- ▶ Truncation errors
- ▶ Solving linear systems

Today

- ▶ Solving linear systems

Announcements

- ▶ Homework 1 was posted last week; is due next week Mon, Sep 23 *before class*

Recap: Linear systems of equations

It is said that 70-80% of computational mathematics research involves solving systems of m linear equations in n unknowns

$$\sum_{j=1}^n a_{ij}x_j = b_i, \quad i = 1, \dots, m.$$

Linear systems arise directly from discrete models (e.g., in machine learning). Or through representing some abstract linear operator (such as a differential operator) in a finite basis as when numerically solving partial differential equations.

The common abstract way of writing systems of linear equations is

$$\mathbf{Ax} = \mathbf{b},$$

with matrix $\mathbf{A} \in \mathbb{R}^{m \times n}$, right-hand side $\mathbf{b} \in \mathbb{R}^m$, and solution $\mathbf{x} \in \mathbb{R}^n$

The goal is to calculate solution \mathbf{x} given data \mathbf{A}, \mathbf{b} in a numerically stable and computationally efficient way.

Recap: The matrix inverse

- ▶ A square matrix $\mathbf{A} \in \mathbb{R}^{n \times n}$ is invertible or nonsingular if there exists a matrix inverse $\mathbf{A}^{-1} = \mathbf{B} \in \mathbb{R}^{n \times n}$ such that

$$\mathbf{AB} = \mathbf{BA} = \mathbf{I},$$

where \mathbf{I} is the identity matrix.

- ▶ Matrix norm induced by a given vector norm

$$\|\mathbf{A}\| = \sup_{\mathbf{x} \neq \mathbf{0}} \frac{\|\mathbf{Ax}\|}{\|\mathbf{x}\|} \implies \|\mathbf{Ax}\| \leq \|\mathbf{A}\| \|\mathbf{x}\|$$

with sub-multiplicativity: $\|\mathbf{AB}\| \leq \|\mathbf{A}\| \|\mathbf{B}\|$

- ▶ Special case of interest: The 2-norm or spectral norm: $\|\mathbf{A}\|_2 = \sigma_1$ (largest singular value)
- ▶ The Euclidean or Frobenius norm is *not* an induced norm

$$\|\mathbf{A}\|_F = \sqrt{\sum_{i,j} |a_{ij}|^2}$$

but still is sub-multiplicative: $\|\mathbf{AB}\|_F \leq \|\mathbf{A}\|_F \|\mathbf{B}\|_F$

Recap: Condition of solving system of linear equations (cont'd)

Now consider the general perturbations of the data

$$(\mathbf{A} + \delta\mathbf{A})(\mathbf{x} + \delta\mathbf{x}) = \mathbf{b} + \delta\mathbf{b}$$

One obtains the condition (proof in Quarteroni et al., Sec. 3.1)

$$\frac{\|\delta\mathbf{x}\|}{\|\mathbf{x}\|} \leq \frac{\kappa(\mathbf{A})}{1 - \kappa(\mathbf{A}) \frac{\|\delta\mathbf{A}\|}{\|\mathbf{A}\|}} \left(\frac{\|\delta\mathbf{b}\|}{\|\mathbf{b}\|} + \frac{\|\delta\mathbf{A}\|}{\|\mathbf{A}\|} \right)$$

Important practical estimate: Roundoff error in the data, with rounding unit u (recall $\approx 10^{-16}$ for double precision), produces a relative error

$$\frac{\|\delta\mathbf{x}\|_{\infty}}{\|\mathbf{x}\|_{\infty}} \leq \frac{2}{1 - \kappa(\mathbf{A})u} u\kappa(\mathbf{A})$$

\implies makes *no* sense to try to numerically solve systems with $\kappa(\mathbf{A}) > 10^{16}$ in double precision

Recap: Numerical solution of linear systems

There are many numerical methods for solving a system of linear equations

The most appropriate method depends on the properties of \mathbf{A}

- ▶ General dense matrices, where the entries in \mathbf{A} are mostly non-zero and nothing special is known \rightsquigarrow we focus on Gaussian elimination today
- ▶ General sparse matrices, where only a small fraction of $a_{ij} \neq 0$ (sparse typically means that $\mathcal{O}(n)$ entries are non-zero in an $n \times n$ matrix)
- ▶ Symmetric and positive-definite matrices
- ▶ Special structured sparse matrices, often arising from specific physical properties of the underlying system

It is also important to consider how many times a linear system with the same or related matrix or right-hand side needs to be solved.

Gauss elimination and LU factorization

$$\begin{array}{|c|c|c|} \hline a_{11}^{(1)} & a_{12}^{(1)} & a_{13}^{(1)} \\ \hline a_{21}^{(1)} & a_{22}^{(1)} & a_{23}^{(1)} \\ \hline a_{31}^{(1)} & a_{32}^{(1)} & a_{33}^{(1)} \\ \hline \end{array} \left[\begin{array}{c} x_1 \\ \hline x_2 \\ \hline x_3 \end{array} \right] = \left[\begin{array}{c} b_1^{(1)} \\ b_2^{(1)} \\ b_3^{(1)} \end{array} \right]$$

↓ eliminate

multipliers

$$l_{21} = a_{21}^{(1)} / a_{11}^{(1)}$$

$$l_{31} = a_{31}^{(1)} / a_{11}^{(1)}$$

$$\begin{bmatrix} a_{11}^{(1)} & a_{12}^{(1)} & a_{13}^{(1)} \\ a_{21}^{(1)} - l_{21} a_{11}^{(1)} = 0 & a_{22}^{(1)} - l_{21} a_{12}^{(1)} & a_{23}^{(1)} - l_{21} a_{13}^{(1)} \\ 0 & a_{32}^{(1)} - l_{31} a_{12}^{(1)} & a_{33}^{(1)} - l_{31} a_{13}^{(1)} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} = \begin{bmatrix} b_1^{(1)} \\ b_2^{(1)} - l_{21} b_1^{(1)} \\ b_3^{(1)} - l_{31} b_1^{(1)} \end{bmatrix}$$

$$\begin{bmatrix} a_{11}^{(1)} & a_{12}^{(1)} & a_{13}^{(1)} \\ 0 & a_{22}^{(2)} & a_{23}^{(2)} \\ 0 & a_{32}^{(2)} & a_{33}^{(2)} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} = \begin{bmatrix} b_1^{(1)} \\ b_2^{(2)} \\ b_3^{(2)} \end{bmatrix}$$

$$\begin{bmatrix} a_{11}^{(1)} & a_{12}^{(1)} & a_{13}^{(1)} \\ & a_{22}^{(2)} & a_{23}^{(2)} \\ & 0 & a_{33}^{(3)} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} = \begin{bmatrix} b_1^{(1)} \\ b_2^{(2)} \\ b_3^{(3)} \end{bmatrix}$$

$$a_{33}^{(3)} = a_{33}^{(2)} - a_{22}^{(2)} \cdot l_{32}$$

$$b_3^{(3)} = b_3^{(2)} - b_2^{(2)} \cdot l_{32}$$

Upper triangular system

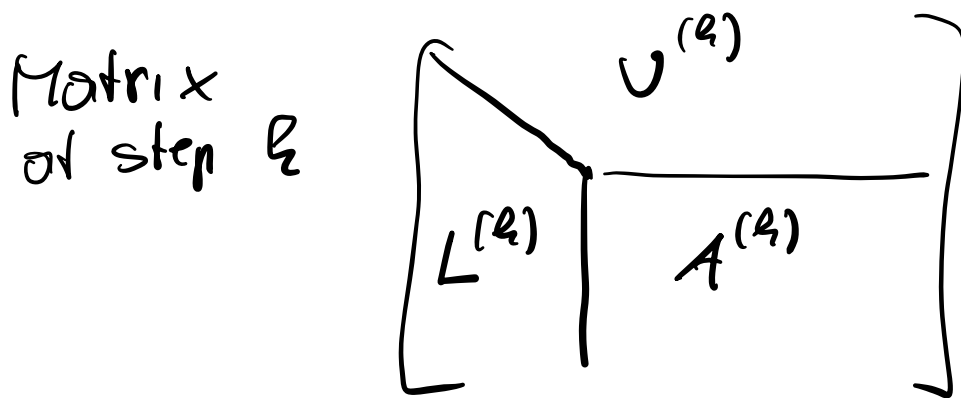
$$x_3 = b_3^{(3)} / a_{33}^{(3)} \Rightarrow \text{eliminate } x_3$$

$$\begin{bmatrix} a_{11}^{(1)} & a_{12}^{(1)} \\ 0 & a_{22}^{(2)} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} b_1^{(1)} - a_{13}^{(1)} x_3 \\ b_2^{(2)} - a_{23}^{(2)} x_3 \end{bmatrix}$$

Solve $x_2 = \frac{b_2^{(2)} - a_{23}^{(2)} x_3}{a_{22}^{(2)}}$

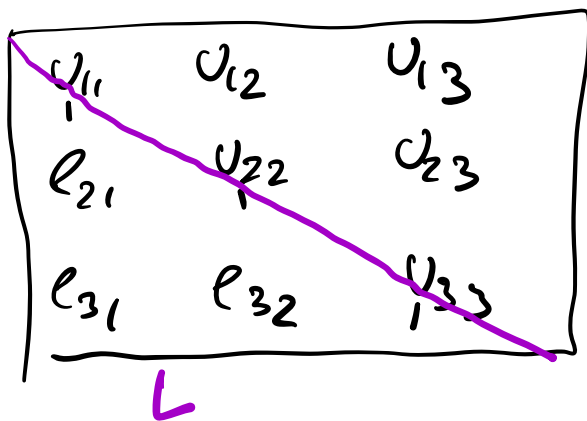
$\rightarrow x_1 = \dots$

Idea: store the multipliers in the lower triangle of A



$k=2$

$$\begin{bmatrix} u_{11} & u_{12} & u_{13} \\ l_{21} & \begin{array}{|c|} \hline a_{22}^{(2)} & a_{23}^{(2)} \\ \hline \end{array} \\ l_{31} & \begin{array}{|c|} \hline a_{32}^{(2)} & a_{33}^{(2)} \\ \hline \end{array} \end{bmatrix}$$



$$A = LU$$

$$\begin{matrix} & A & & x & & b \\ \begin{bmatrix} 1 & 1 & 3 \\ 2 & 2 & 2 \\ 3 & 6 & 4 \end{bmatrix} & & \begin{bmatrix} \\ \\ \end{bmatrix} & = & \begin{bmatrix} 5 \\ 6 \\ 13 \end{bmatrix} & & e_{21} = \frac{2}{1} \\ & & & & & & e_{31} = \frac{3}{1} \end{matrix}$$

$$\begin{bmatrix} 1 & 1 & 3 \\ 2 & 0 & 4 \\ 3 & 3 & -5 \end{bmatrix}$$

↕ swap rows

$$\begin{bmatrix} 1 & 1 & 3 \\ 3 & 3 & -5 \\ 2 & 0 & -4 \end{bmatrix}$$

Gauss elimination in Matlab

```
1: function A = mylu(A) % In-place LU factorization
2: % need square matrix
3: [n, m] = size(A);
4: assert(n == m);
5: for k=1:(n-1) % for variable x(k)
6:     % Assumed A(k, k) non-zero and then
7:     % calculate multipliers in column k
8:     A((k + 1):n, k) = A((k + 1):n, k)/A(k, k);
9:     for j = (k + 1):n
10:        % eliminate variable x(k)
11:        A((k + 1):n, j) = A((k + 1):n, j) - A((k + 1):
            n, k)*A(k, j);
12:     end
13: end
14: end
```

- ▶ Gaussian elimination is a general method for dense matrices and is commonly used
- ▶ Implementing Gaussian elimination efficiently is difficult and we will not discuss it
↪ course on HPC
- ▶ The LAPACK public-domain library is the main repository for excellent implementations of dense linear solvers
- ▶ Matlab (and numpy) use highly optimized variants of GEM by default, mostly based on LAPACK
- ▶ Matlab (and numpy) have specialized solvers for special cases of matrices, so always check help pages!

Problem?

Problem?

```
1: >> A = [1 1 3; 2 2 2; 3 6 4]
2:
3: A =
4:
5:      1      1      3
6:      2      2      2
7:      3      6      4
8:
9: >> mylu(A)
10:
11: ans =
12:
13:      1      1      3
14:      2      0     -4
15:      3     Inf     Inf
```

LU with pivoting

Zero diagonal entries (pivots) pose a problem \rightsquigarrow pivoting by swapping rows

\rightsquigarrow board

$$\begin{bmatrix} 1 & 1 & 3 \\ 1 & 3 & -5 \\ 3 & 1 & -4 \\ 2 & 0 & 1 \end{bmatrix}$$

$$LU = PA$$

permutation \nearrow

$$\begin{bmatrix} 1 & 0 & 0 \\ 0 & 0 & 1 \\ 0 & 1 & 0 \end{bmatrix}$$

```
1: >> A = [1 1 3; 2 2 2; 3 6 4];
2: >> P = [1 0 0; 0 0 1; 0 1 0] % swap row 2 and 3
3:
4: P =
5:
6:     1     0     0
7:     0     0     1
8:     0     1     0
9:
10: >> mylu(P*A)
11:
12: ans =
13:
14:     1     1     3
15:     3     3    -5
16:     2     0    -4
```

- ▶ For any square (regular or singular) matrix \mathbf{A} , partial (row) pivoting ensures exists of

$$\mathbf{PA} = \mathbf{LU}$$

where \mathbf{P} is a permutation matrix

- ▶ **Q:** What else could pivoting be useful for?
↪ let's see what Matlab does

```
1: >> [L, U, P] = lu(A) % built-in lu
2: L =
3:     1.0000         0         0
4:     0.6667     1.0000         0
5:     0.3333     0.5000     1.0000
6: U =
7:     3.0000     6.0000     4.0000
8:         0    -2.0000    -0.6667
9:         0         0     2.0000
10: P =
11:     0     0     1
12:     0     1     0
13:     1     0     0
14: >> norm(L*U - P*A)
15: ans = 0
```

```
1: >> [L, U, P] = lu(A) % built-in lu
2: L =
3:     1.0000         0         0
4:     0.6667     1.0000         0
5:     0.3333     0.5000     1.0000
6: U =
7:     3.0000     6.0000     4.0000
8:         0    -2.0000    -0.6667
9:         0         0     2.0000
10: P =
11:     0     0     1
12:     0     1     0
13:     1     0     0
14: >> norm(L*U - P*A)
15: ans = 0
```

Reverses order of rows rather than just swapping 2 and 3.
Leads to entries of L with magnitude ≤ 1

↪ board

$$A = \begin{bmatrix} 10^{-20} & 1 \\ 1 & 1 \end{bmatrix}$$

$$e_{21} = \frac{a_{21}^{(1)}}{a_{11}^{(1)}} = \frac{1}{10^{-20}} = 10^{20}$$

$$\left[\begin{array}{c|c} 10^{-20} & 1 \\ \hline 10^{20} & \nearrow 1 \end{array} \right]$$

$$= \begin{bmatrix} 10^{-20} & 1 \\ 10^{20} & 1 - 10^{20} \end{bmatrix}$$

$$\begin{aligned} & 1 - 1 \cdot 10^{20} \\ & = 1 - 10^{20} \end{aligned}$$

$$L = \begin{bmatrix} 1 & 0 \\ 10^{20} & 1 \end{bmatrix}, \quad U = \begin{bmatrix} 10^{-20} & 1 \\ 0 & 1 - 10^{20} \end{bmatrix}$$

$$LU = A$$

floating point: $1 - 10^{20} \approx -10^{20}$ ↙

↪ board

```
1: >> A = [1e-20 1; 1 1]
2:
3: A =
4:
5:      1.0000e-20      1.0000e+00
6:      1.0000e+00      1.0000e+00
7:
8: >> LUmat = mylu(A);
9: L = [1 0; LUmat(2, 1) 1];
10: U = LUmat; U(2, 1) = 0;
11: L*U
12:
13: ans =
14:
15:      1.0000e-20      1.0000e+00
16:      1.0000e+00           0
```

```
1: >> [L, U, P] = lu(A)
2: L =
3:     1.0000e+00           0
4:     1.0000e-20     1.0000e+00
5: U =
6:     1     1
7:     0     1
8: P =
9:     0     1
10:    1     0
11: >> P' * L * U
12: ans =
13:     1.0000e-20     1.0000e+00
14:     1.0000e+00     1.0000e+00
15: >> A
16: A =
17:     1.0000e-20     1.0000e+00
18:     1.0000e+00     1.0000e+00
```


Instability of LU decomposition *without* pivoting

If \mathbf{A} has an \mathbf{LU} factorization, then the computed $\tilde{\mathbf{L}}$ and $\tilde{\mathbf{U}}$ obtained in floating-point arithmetic with Gaussian elimination satisfy $\tilde{\mathbf{L}}\tilde{\mathbf{U}} = \mathbf{A} + \delta\mathbf{A}$ with the bound

$$\frac{\|\delta\mathbf{A}\|}{\|\mathbf{L}\|\|\mathbf{U}\|} \in \mathcal{O}(u),$$

where u is the roundoff unit.

- ▶ Notice that we would have liked to bound $\|\delta\mathbf{A}\|/\|\mathbf{A}\|$ but we got $\|\delta\mathbf{A}\|/(\|\mathbf{L}\|\|\mathbf{U}\|)$
- ▶ Thus, for matrices with $\|\mathbf{L}\|\|\mathbf{U}\| \approx \|\mathbf{A}\|$, the algorithm will show stable behavior
- ▶ However, if $\|\mathbf{L}\|\|\mathbf{U}\| \not\approx \|\mathbf{A}\|$, then we can get an unstable result

↪ Gaussian elimination is *not* stable in general

Example ↪ board

$$\mathbf{A} = \begin{bmatrix} 10^{-20} & 1 \\ 1 & 1 \end{bmatrix}, \mathbf{L} = \begin{bmatrix} 1 & 0 \\ 10^{20} & 1 \end{bmatrix}, \mathbf{U} = \begin{bmatrix} 10^{-20} & 1 \\ 0 & 1 - 10^{20} \end{bmatrix}$$

Compare $\|\mathbf{L}\|_{\infty}\|\mathbf{U}\|_{\infty}$ and $\|\mathbf{A}\|_{\infty}$

$$A = \begin{bmatrix} 10^{-20} & 1 \\ 1 & 1 \end{bmatrix}, \quad L = \begin{bmatrix} 1 & 0 \\ 10^{20} & 1 \end{bmatrix}, \quad U = \begin{bmatrix} 10^{-20} & 1 \\ 0 & 1 - 10^{20} \end{bmatrix}$$

Recall: operator norm induced by

$$\|x\|_\infty = \max_i |x_i|$$

$$\|A\|_\infty = \max_i \sum_{j=1}^n |a_{ij}|$$

$$\|A\|_\infty = 2$$

$$\|L\|_\infty = 1 + 10^{20}$$

$$\|U\|_\infty$$

$$\|L\|_\infty \|U\|_\infty \gg \|A\|_\infty$$

$$2 \cdot 10^{20}$$

LU with row pivoting to maximum element

```
1: function [A, P] = myplu(A) % In-place LU factorization
2: [n, m] = size(A); P = eye(n);
3: for k=1:(n-1) % for variable x(k)
4:     [~, selI] = max(abs(A(k:n, k))); % select pivot
5:     c = A(k, k:end); d = P(k, :);
6:     A(k, k:end) = A(selI + (k-1), k:end); P(k, :) = P(
        selI + (k-1), :);
7:     A(selI + (k-1), k:end) = c; P(selI + k-1, :) = d;
8:     % calculate multipliers in column k
9:     A((k + 1):n, k) = A((k + 1):n, k)/A(k, k);
10:    for j = (k + 1):n % eliminate variable x(k)
11:        A((k + 1):n, j) = A((k + 1):n, j) - A((k + 1):
            n, k)*A(k, j);
12:    end
13: end
14: end
```

```
1: >> A = [1e-20 1; 1 1]
2: [LUmat, P] = myplu(A);
3: L = [1 0; LUmat(2, 1) 1];
4: U = LUmat; U(2, 1) = 0;
5: P'*L*U
6:
7: A =
8:
9:      1.0000e-20      1.0000e+00
10:     1.0000e+00      1.0000e+00
11:
12:
13: ans =
14:
15:      1.0000e-20      1.0000e+00
16:     1.0000e+00      1.0000e+00
```

Stability of Gaussian elimination with pivoting

For $\mathbf{A} = \mathbf{LU}$, introduce the growth factor

$$\rho = \frac{\max_{i,j} |u_{ij}|}{\max_{i,j} |a_{ij}|}$$

where u_{ij} is the i, j -th element of \mathbf{U}

Consider the factorization $\mathbf{PA} = \mathbf{LU}$ with partial row pivoting w.r.t. taking the maximum element for a matrix \mathbf{A} of dimension $n \times n$. Gaussian elimination gives $\tilde{\mathbf{P}}, \tilde{\mathbf{L}}, \tilde{\mathbf{U}}$ that satisfy

$$\tilde{\mathbf{L}}\tilde{\mathbf{U}} = \tilde{\mathbf{P}}\mathbf{A} + \delta\mathbf{A}, \quad \frac{\|\delta\mathbf{A}\|}{\|\mathbf{A}\|} \in \mathcal{O}(\rho u),$$

where u is the roundoff unit. If all off-diagonal entries of \mathbf{L} are < 1 , implying that there are no ties in the selection of pivots in exact arithmetic, then $\mathbf{P} = \tilde{\mathbf{P}}$ for sufficiently small u .

This means that Gaussian elimination with partial pivoting is backward stable if ρ holds uniformly for matrices with $n \times n$.

- ▶ For any square (regular or singular) matrix \mathbf{A} , partial (row) pivoting ensures exists of

$$\mathbf{PA} = \mathbf{LU}$$

where \mathbf{P} is a permutation matrix

- ▶ Furthermore, pivoting (w.r.t. $\max |a_{ij}|$) leads to a backward stable algorithm. **However**, the growth factor ρ can be huge and grow with the dimension of \mathbf{A} ! Fortunately, large factors ρ “never seem to appear in real applications.” (Trefethen & Bau, Chapter 22)
- ▶ There also is full pivoting (rows + columns)

$$\mathbf{PAQ} = \mathbf{LU}$$

to further increases stability but it usually is not worth it in practice (higher costs to search for pivoting element over rows and columns but little improvement in terms of stability)

Solving linear systems

- ▶ Once an LU factorization is available, solving a linear system is cheap:

$$\mathbf{Ax} = \mathbf{LUx} = \mathbf{L(Ux)} = \mathbf{Ly} = \mathbf{b}$$

- ▶ Solve for \mathbf{y} using *forward substitution*
- ▶ Solve for \mathbf{x} by using *backward substitution* $\mathbf{Ux} = \mathbf{y}$

What is *forward/backward substitution*? \rightsquigarrow board

$$\begin{bmatrix} l_{11} & 0 & \cdots & \cdots & 0 \\ l_{21} & l_{22} & 0 & \cdots & 0 \\ \vdots & & & & \vdots \\ l_{n1} & \cdots & \cdots & \cdots & l_{nn} \end{bmatrix} \begin{bmatrix} x_1 \\ \vdots \\ \vdots \\ x_n \end{bmatrix} = \begin{bmatrix} b_1 \\ \vdots \\ \vdots \\ b_n \end{bmatrix}$$

Forward substitution

$$\begin{bmatrix} l_{11} & 0 & \dots & 0 \\ l_{21} & l_{22} & 0 & \dots & 0 \\ \vdots & & & & \\ l_{n1} & l_{n2} & \dots & l_{nn} \end{bmatrix} \begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix} = \begin{bmatrix} b_1 \\ b_2 \\ \vdots \\ b_n \end{bmatrix}$$

$$x_1 = b_1 / l_{11}$$

$$x_2 = (b_2 - l_{21} x_1) / l_{22}$$

⋮

$$x_n = (b_n - l_{n1} x_1 - \dots - l_{n,n-1} x_{n-1}) / l_{nn}$$

- ▶ If row pivoting is used, the same process works by also permuting the right-hand side \mathbf{b}

$$\mathbf{PAx} = \mathbf{LUx} = \mathbf{Ly} = \mathbf{Pb}$$

or formally (**never implement inverse for solving linear systems of equations**)

$$\mathbf{x} = (\mathbf{LU})^{-1}\mathbf{Pb} = \mathbf{U}^{-1}\mathbf{L}^{-1}\mathbf{Pb}$$

- ▶ Because \mathbf{P} is orthonormal, we have $\mathbf{P}^{-1} = \mathbf{P}^T$ and thus

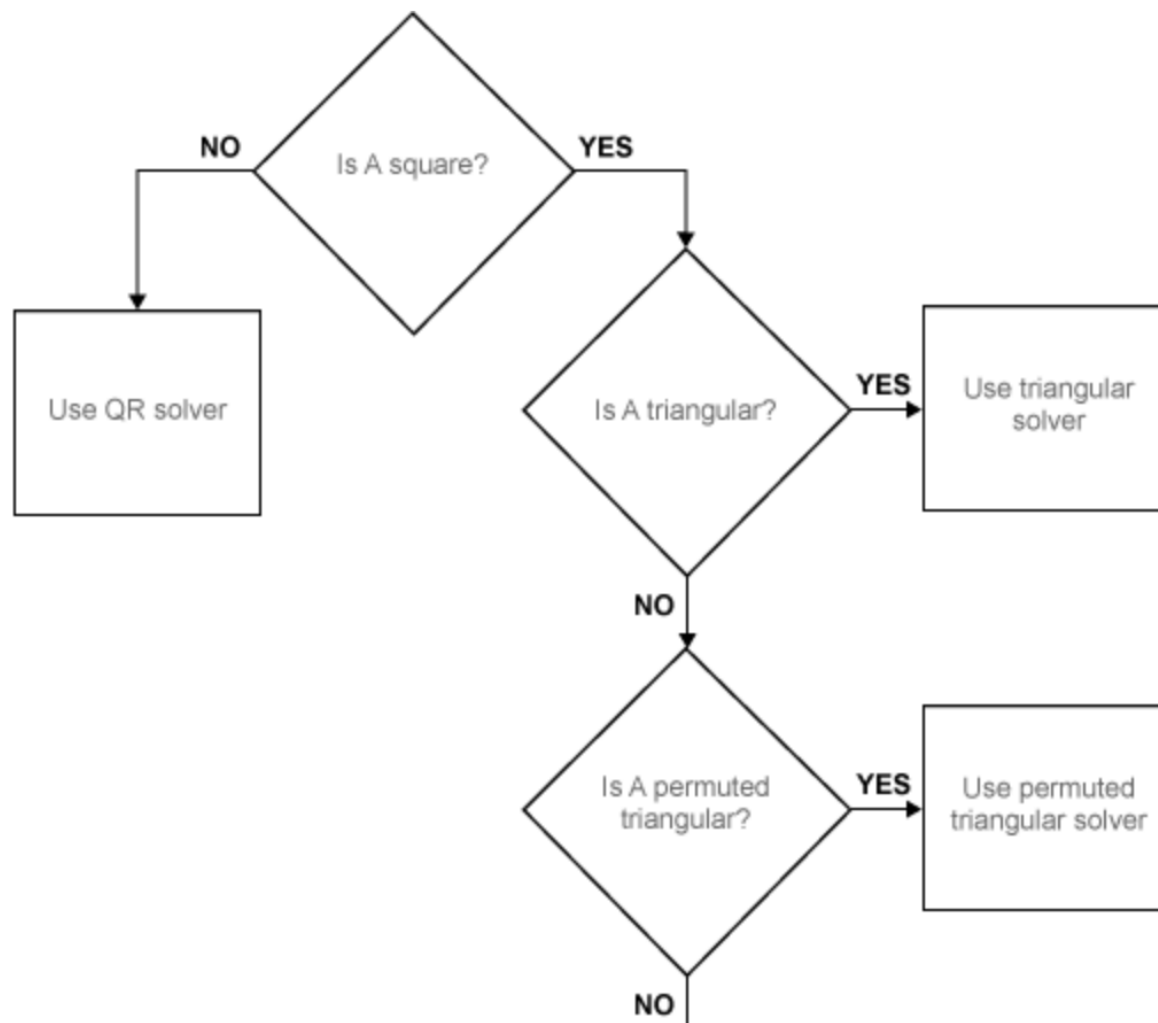
$$\mathbf{A} = \mathbf{P}^{-1}\mathbf{LU} = (\mathbf{P}^T\mathbf{L})\mathbf{U} = \tilde{\mathbf{L}}\mathbf{U},$$

with $\tilde{\mathbf{L}}$ a row permutation of a unit lower triangular matrix

In Matlab, the backslash operator solves linear systems (see help `mldivide`)

```
1: A = [1 2 3; 4 5 6; 7 8 0];
2: b = [2; 1; -1];
3: x = A\b; x'
4: [L, U] = lu(A)
5: y = L\b; x = U\y; x'
6: ans =
7:    -2.5556e+00    2.1111e+00    1.1111e-01
8: L =
9:    1.4286e-01    1.0000e+00    0
10:    5.7143e-01    5.0000e-01    1.0000e+00
11:    1.0000e+00    0    0
12: U =
13:    7.0000e+00    8.0000e+00    0
14:    0    8.5714e-01    3.0000e+00
15:    0    0    4.5000e+00
16: ans =
17:    -2.5556e+00    2.1111e+00    1.1111e-01
```

Is the permuted triangular matrix \tilde{L} (which we get from $[L, U] = \text{lu}(A)$) detected as such? Yes!



Numerical Methods I

MATH-GA 2010.001/CSCI-GA 2420.001

Benjamin Peherstorfer
Courant Institute, NYU

Based on slides by G. Stadler and A. Donev

Today

Last time

- ▶ Solving linear systems
- ▶ LU decomposition
- ▶ Pivoting

Today

- ▶ Cost analysis of LU decomposition
- ▶ Solving linear systems with sparse matrices
- ▶ Least-squares problems

Announcements

- ▶ Homework 2 has been posted; is due next week Mon, Oct 7 *before class*

Recap: Condition of solving system of linear equations (cont'd)

Now consider the general perturbations of the data

$$(\mathbf{A} + \delta\mathbf{A})(\mathbf{x} + \delta\mathbf{x}) = \mathbf{b} + \delta\mathbf{b}$$

One obtains the condition (proof in Quarteroni et al., Sec. 3.1)

$$\frac{\|\delta\mathbf{x}\|}{\|\mathbf{x}\|} \leq \frac{\kappa(\mathbf{A})}{1 - \kappa(\mathbf{A}) \frac{\|\delta\mathbf{A}\|}{\|\mathbf{A}\|}} \left(\frac{\|\delta\mathbf{b}\|}{\|\mathbf{b}\|} + \frac{\|\delta\mathbf{A}\|}{\|\mathbf{A}\|} \right)$$

Important practical estimate: Roundoff error in the data, with rounding unit u (recall $\approx 10^{-16}$ for double precision), produces a relative error

$$\frac{\|\delta\mathbf{x}\|_{\infty}}{\|\mathbf{x}\|_{\infty}} \leq \frac{2}{1 - \kappa(\mathbf{A})u} u\kappa(\mathbf{A})$$

\implies makes *no* sense to try to numerically solve systems with $\kappa(\mathbf{A}) > 10^{16}$ in double precision

Recap: LU decomposition

Gauss elimination and LU factorization

index of step \rightarrow

$$\begin{bmatrix} a_{11}^{(1)} & a_{12}^{(1)} & a_{13}^{(1)} \\ a_{21}^{(1)} & a_{22}^{(1)} & a_{23}^{(1)} \\ a_{31}^{(1)} & a_{32}^{(1)} & a_{33}^{(1)} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} = \begin{bmatrix} b_1^{(1)} \\ b_2^{(1)} \\ b_3^{(1)} \end{bmatrix}$$

multiply first row by l_{21} and subtract from 2nd row

$$\begin{bmatrix} a_{11}^{(1)} & a_{12}^{(1)} & a_{13}^{(1)} \\ 0 & a_{22}^{(2)} & a_{23}^{(2)} \\ 0 & a_{32}^{(2)} & a_{33}^{(2)} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} = \begin{bmatrix} b_1^{(1)} \\ b_2^{(2)} \\ b_3^{(2)} \end{bmatrix}$$

$l_{21} = \frac{a_{21}^{(1)}}{a_{11}^{(1)}}$

$$\begin{bmatrix} a_{11}^{(1)} & a_{12}^{(1)} & a_{13}^{(1)} \\ 0 & a_{22}^{(2)} & a_{23}^{(2)} \\ 0 & 0 & a_{33}^{(3)} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} = \begin{bmatrix} b_1^{(1)} \\ b_2^{(2)} \\ b_3^{(3)} \end{bmatrix}$$

Idea: Store multipliers l_{21}, l_{31}, l_{32} in matrix:

$$\begin{bmatrix} L & U \end{bmatrix}$$

$$A = LU, \quad L = \begin{bmatrix} 1 & & \\ l_{21} & 1 & \\ l_{31} & l_{32} & 1 \end{bmatrix}, \quad U = \begin{bmatrix} u_{11} & \dots & u_{13} \\ 0 & \dots & 0 \\ 0 & \dots & \vdots \end{bmatrix}$$

$$\begin{bmatrix} u_{11} & u_{12} & u_{13} \\ l_{21} & u_{22} & u_{23} \\ l_{31} & l_{32} & u_{33} \end{bmatrix}$$

Recap: Pivoting

LU with pivoting

Zero diagonal entries (pivots) pose a problem \rightsquigarrow pivoting by swapping rows

\rightsquigarrow board $\begin{bmatrix} 1 & 1 & 3 \\ 2 & 2 & 2 \\ 3 & 6 & 4 \end{bmatrix}$ $\xrightarrow{\text{Gauss elimination}}$ $\begin{bmatrix} 1 & 1 & 3 \\ l_{21}=2 & 0 & -4 \\ l_{31}=3 & 3 & -5 \end{bmatrix}$ Problem! Cannot divide by zero, pivot element is zero.

Swap rows $\begin{bmatrix} 1 & 1 & 3 \\ 3 & 3 & -5 \\ 2 & 0 & -4 \end{bmatrix}$ done - in general, use 3 as pivot.

$LU = PA$

P ... permutation matrix.

$P_2 = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 0 & 1 \\ 0 & 1 & 0 \end{bmatrix}$

Theorem: This always exists with appropriate perm. matrix P

Recap: Stability of Gaussian elimination with pivoting

For $\mathbf{A} = \mathbf{LU}$, introduce the growth factor

$$\rho = \frac{\max_{i,j} |u_{ij}|}{\max_{i,j} |a_{ij}|}$$

where u_{ij} is the i, j -th element of \mathbf{U}

Consider the factorization $\mathbf{PA} = \mathbf{LU}$ with partial row pivoting w.r.t. taking the maximum element for a matrix \mathbf{A} of dimension $n \times n$. Gaussian elimination gives $\tilde{\mathbf{P}}, \tilde{\mathbf{L}}, \tilde{\mathbf{U}}$ that satisfy

$$\tilde{\mathbf{L}}\tilde{\mathbf{U}} = \tilde{\mathbf{P}}\mathbf{A} + \delta\mathbf{A}, \quad \frac{\|\delta\mathbf{A}\|}{\|\mathbf{A}\|} \in \mathcal{O}(\rho u),$$

where u is the roundoff unit. If all off-diagonal entries of \mathbf{L} are < 1 , implying that there are no ties in the selection of pivots in exact arithmetic, then $\mathbf{P} = \tilde{\mathbf{P}}$ for sufficiently small u .

This means that Gaussian elimination with partial pivoting is backward stable if ρ holds uniformly for matrices with $n \times n$.

Recap: LU with row pivoting to maximum element

```
1: function [A, P] = myplu(A) % In-place LU factorization
2: [n, m] = size(A); P = eye(n);
3: for k=1:(n-1) % for variable x(k)
4:     [~, selI] = max(abs(A(k:n, k))); % select pivot
5:     c = A(k, k:end); d = P(k, :);
6:     A(k, k:end) = A(selI + (k-1), k:end); P(k, :) = P(
        selI + (k-1), :);
7:     A(selI + (k-1), k:end) = c; P(selI + k-1, :) = d;
8:     % calculate multipliers in column k
9:     A((k + 1):n, k) = A((k + 1):n, k)/A(k, k);
10:    for j = (k + 1):n % eliminate variable x(k)
11:        A((k + 1):n, j) = A((k + 1):n, j) - A((k + 1):
            n, k)*A(k, j);
12:    end
13: end
14: end
```

Costs: Forward/backward substitution

↪ board

Forward substitution

$$\begin{bmatrix} l_{11} & 0 & \dots & 0 \\ l_{21} & l_{22} & 0 & \dots & 0 \\ \vdots & & & & \\ l_{n1} & l_{n2} & \dots & l_{nn} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} = \begin{bmatrix} b_1 \\ b_2 \\ \vdots \\ b_n \end{bmatrix}$$

$$x_1 = b_1 / l_{11}$$

$$x_2 = (b_2 - l_{21}x_1) / l_{22}$$

\vdots

[\approx 1 division]

[\approx 1 div, 1 mul, 1 add]

\vdots

$$x_n = (b_n - l_{n1}x_1 - \dots - l_{n,n-1}x_{n-1}) / l_{nn}$$

[\approx 1 div, $(n-1)$ mul, $(n-1)$ add]

add and multiplications

$$\sum_{i=2}^n (i-1) = \sum_{i=1}^{n-1} i = \frac{n(n-1)}{2}$$

divisions: n

$$\text{total FLOPs: } n + 2 \frac{n(n-1)}{2} \in \mathcal{O}(n^2)$$

Costs: Forward/backward substitution

↪ board

Forward substitution requires

$$\frac{n(n-1)}{2} \text{ multiplications/additions,}$$

n divisions.

Overall: $\sim n^2$ **floating point operations** (flops) ↪ costs scale as $\mathcal{O}(n^2)$

Similarly, backward substitution has costs that scale as $\mathcal{O}(n^2)$

We count flops to estimate the computational time/effort. **What else matters?**

Costs: Forward/backward substitution

↪ board

Forward substitution requires

$$\frac{n(n-1)}{2} \text{ multiplications/additions,}$$

n divisions.

Overall: $\sim n^2$ **floating point operations** (flops) ↪ costs scale as $\mathcal{O}(n^2)$

Similarly, backward substitution has costs that scale as $\mathcal{O}(n^2)$

We count flops to estimate the computational time/effort. **What else matters?** Besides floating point operations, **computer memory access** has a significant influence on the efficiency of numerical methods.

Costs

For forward [backward] substitution at step k there are $\approx k [(n - k)]$ multiplications and subtractions plus a few divisions. The total over all n steps is

$$\sum_{k=1}^n k \in \mathcal{O}(n^2)$$

\rightsquigarrow the number of floating-point operations (FLOPs) scales as $\mathcal{O}(n^2)$

For Gaussian elimination, at step k , there are $\approx (n - k)^2$ operations. Thus, the total scales as

$$\sum_{k=1}^n (n - k)^2 \in \mathcal{O}(n^3)$$

Directly applying Gaussian elimination scales as

Directly applying Gaussian elimination scales as

$$\mathcal{O}(n^3)$$

Computing LU decomposition scales as

Directly applying Gaussian elimination scales as

$$\mathcal{O}(n^3)$$

Computing LU decomposition scales as

$$\mathcal{O}(n^3)$$

Forward/backward substitution scales as

Directly applying Gaussian elimination scales as

$$\mathcal{O}(n^3)$$

Computing LU decomposition scales as

$$\mathcal{O}(n^3)$$

Forward/backward substitution scales as

$$\mathcal{O}(n^2)$$

LU + forward/backward scales as

Directly applying Gaussian elimination scales as

$$\mathcal{O}(n^3)$$

Computing LU decomposition scales as

$$\mathcal{O}(n^3)$$

Forward/backward substitution scales as

$$\mathcal{O}(n^2)$$

LU + forward/backward scales as

$$\mathcal{O}(n^3)$$

↪ why useful?

Directly applying Gaussian elimination scales as

$$\mathcal{O}(n^3)$$

Computing LU decomposition scales as

$$\mathcal{O}(n^3)$$

Forward/backward substitution scales as

$$\mathcal{O}(n^2)$$

LU + forward/backward scales as

$$\mathcal{O}(n^3)$$

↪ **why useful?** can *reuse* *LU* for other *b*

Choleski factorization

A matrix is **symmetric positive definite (spd)**, if $A = A^T$ and for all $\mathbf{x} \in \mathbb{R}^n$, $\mathbf{x} \neq 0$, the inner product $\langle A\mathbf{x}, \mathbf{x} \rangle > 0$.

For spd matrices, we can compute the factorization:

$$A = LDL^T,$$

where L is a lower triangular matrix with 1's on the diagonal, and D is a positive diagonal matrix.

The Choleski factorization is obtained by multiplying the square root of D (which exists!) with L :

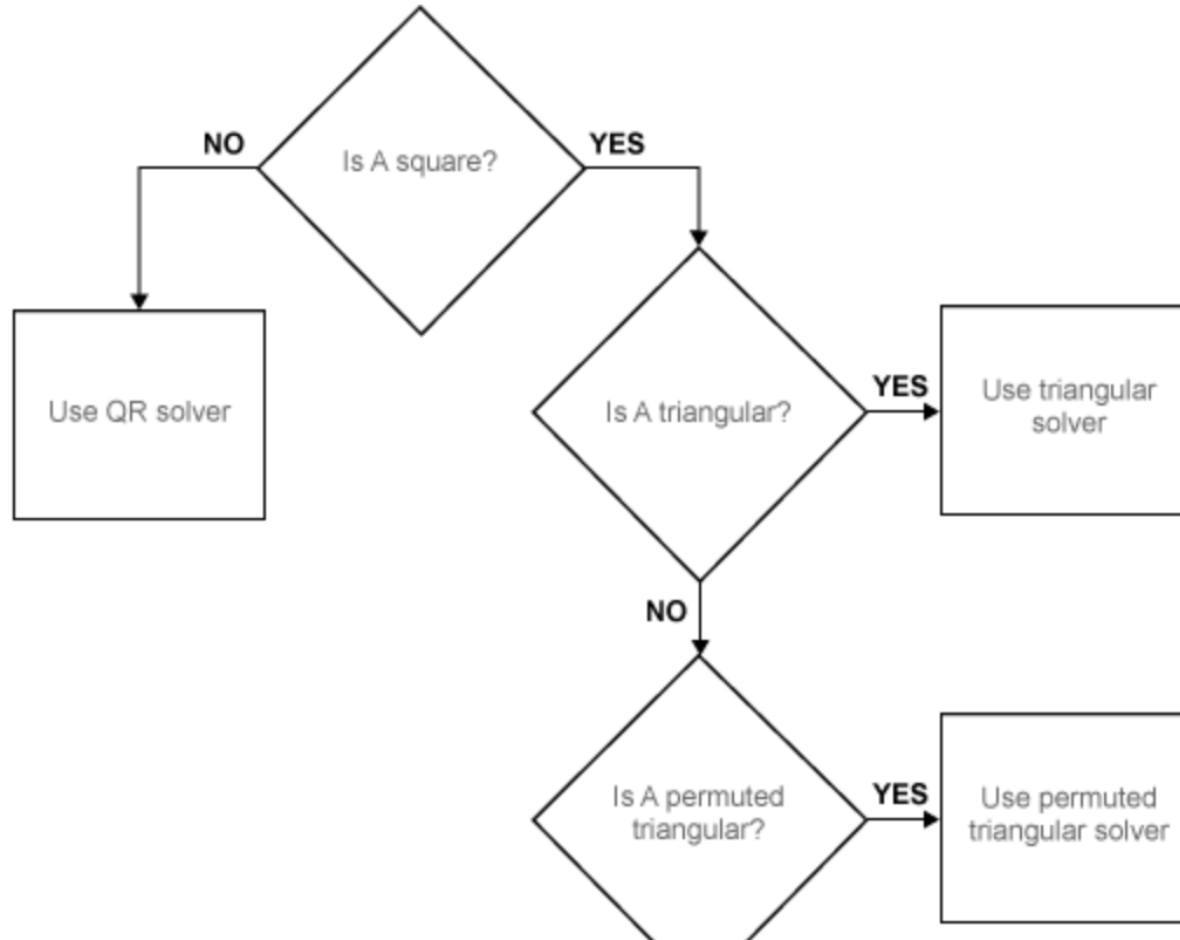
$$A = \bar{L}\bar{L}^T.$$

Algorithms for Choleski factorization are about twice as fast as Gaussian elimination but also scale as $O(n^3)$.

```
1: >> A = randn(1000, 1000)*diag(linspace(1, 10, 1000))*  
    randn(1000, 1000); A = A'*A;  
2: >> tic; chol(A); toc  
3: Elapsed time is 0.004863 seconds.  
4: >> tic; lu(A); toc  
5: Elapsed time is 0.010114 seconds.
```

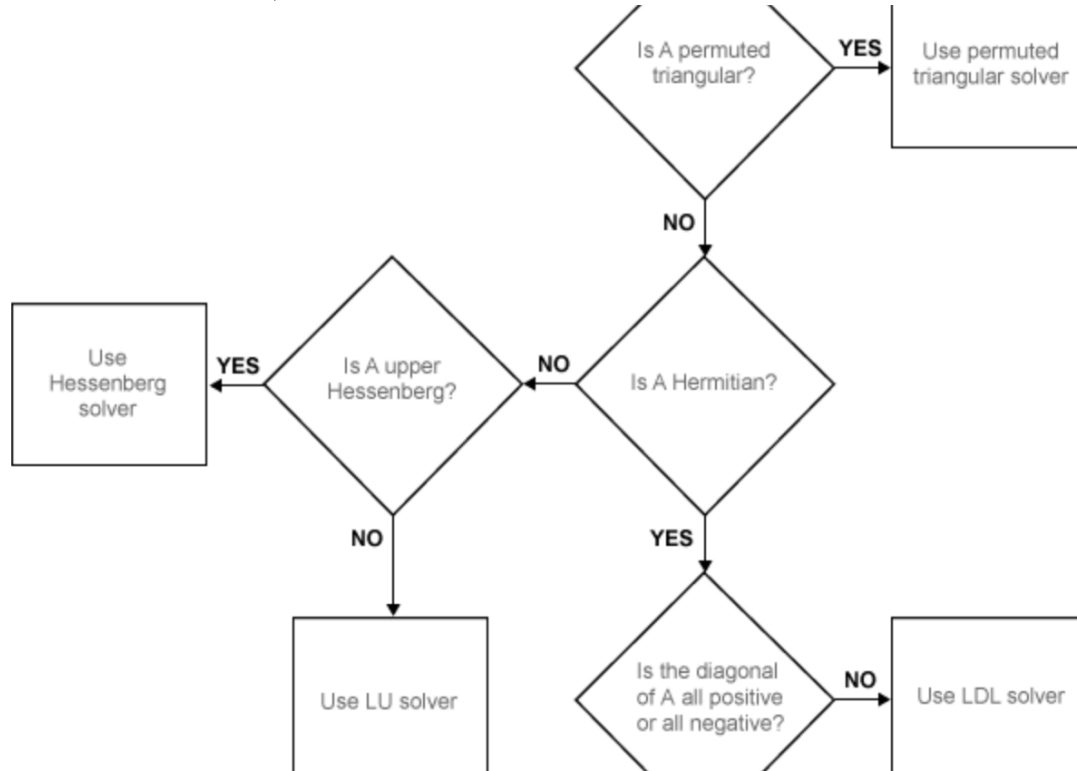
Special matrices in Matlab

Matlab tests for special matrices automatically and chooses a good decomposition/solver



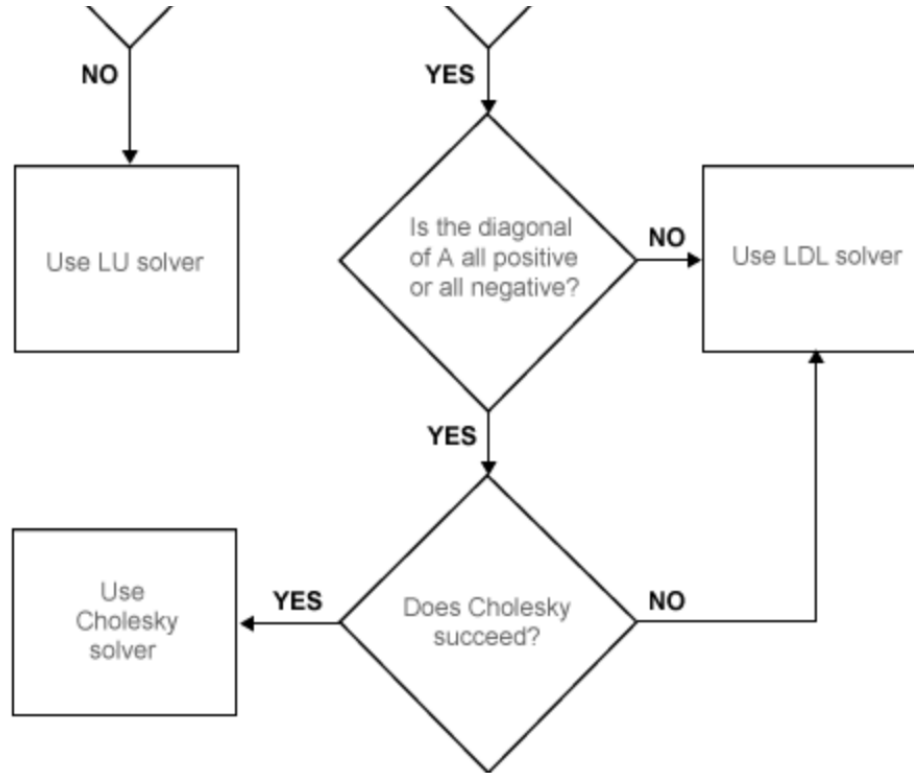
Special matrices in Matlab

Matlab tests for special matrices automatically and chooses a good decomposition/solver



Special matrices in Matlab

Matlab tests for special matrices automatically and chooses a good decomposition/solver



Conclusions/summary

- ▶ The condition of solving a linear system $\mathbf{Ax} = \mathbf{b}$ is determined by the condition number of the matrix \mathbf{A}

$$\kappa(\mathbf{A}) = \|\mathbf{A}\| \|\mathbf{A}^{-1}\| \geq 1$$

- ▶ Gaussian elimination can be used to solve general square linear systems and produces a factorization, if it exists

$$\mathbf{A} = \mathbf{LU}$$

- ▶ Partial pivoting is sufficient for existence and stability of the LU decomposition

$$\mathbf{PA} = \mathbf{LU}, \quad \mathbf{A} = \tilde{\mathbf{L}}\mathbf{U}$$

- ▶ The Cholesky factorization $\mathbf{A} = \mathbf{LL}^T$ exists if \mathbf{A} is spd and then it is the better choice (cheaper) than LU
- ▶ Rely on the highly optimized routines in Matlab (LAPACK) and other software packages than implementing these algorithms yourself \rightsquigarrow take the course on HPC next spring to learn more about the efficient *implementation* of these algorithms

Sparse matrices

Sparse matrix

- ▶ A matrix where a substantial fraction of the entries are zero is called a sparse matrix. Typically, only $\mathcal{O}(N)$ non-zero entries in an $N \times N$ matrix are allowed for sparse algorithms to show benefit over dense linear algebra routines.
- ▶ If we have only $\mathcal{O}(N)$ non-zero entries, then store only those; in contrast to dense matrices. Exploiting sparsity is important (life saving) for large matrices
- ▶ The structure of a sparse matrix refers to the set of indices i, j such that $|a_{ij}| > 0$ and is visualized in Matlab with `spy`
- ▶ The structure of sparse matrices comes from the nature of the problem, e.g., in an inter-city road transportation problem it corresponds to the pairs of cities connected by a road.
- ▶ In fact, just counting the number of non-zero elements is not enough: the sparsity structure is the most important property that determines whether an efficient method exists

Banded matrices

Banded matrices are a very special but common type of sparse matrices, e.g., tridiagonal matrices

$$\begin{bmatrix} a_1 & c_1 & & 0 \\ b_2 & a_2 & \ddots & \\ & \ddots & \ddots & c_{n-1} \\ 0 & & b_n & a_n \end{bmatrix}$$

For example, think of the Laplace problem $u''(x) = f(x)$ on the unit interval and a finite-difference discretization

$$u''(x) \approx \frac{u(x+h) - 2u(x) + u(x-h)}{h^2}$$

on an equidistant grid. This leads to a system of equations with a tridiagonal matrix
↪ Numerical Methods II (Spring semester)

There exist special techniques for banded matrices that are much faster than the general case, e.g., only $8n$ FLOPs

Decomposing sparse matrices

There also are general methods for dealing with sparse matrices, such as the sparse LU factorization.

How well they work depends on the structure of the matrix. **What could go wrong?**

Decomposing sparse matrices

There also are general methods for dealing with sparse matrices, such as the sparse LU factorization.

How well they work depends on the structure of the matrix. **What could go wrong?**

When factorizing sparse matrices, the factors, e.g., L and U , can be much less sparse than $A \rightsquigarrow$ fill-in

For many sparse matrices, there is a large fill-in

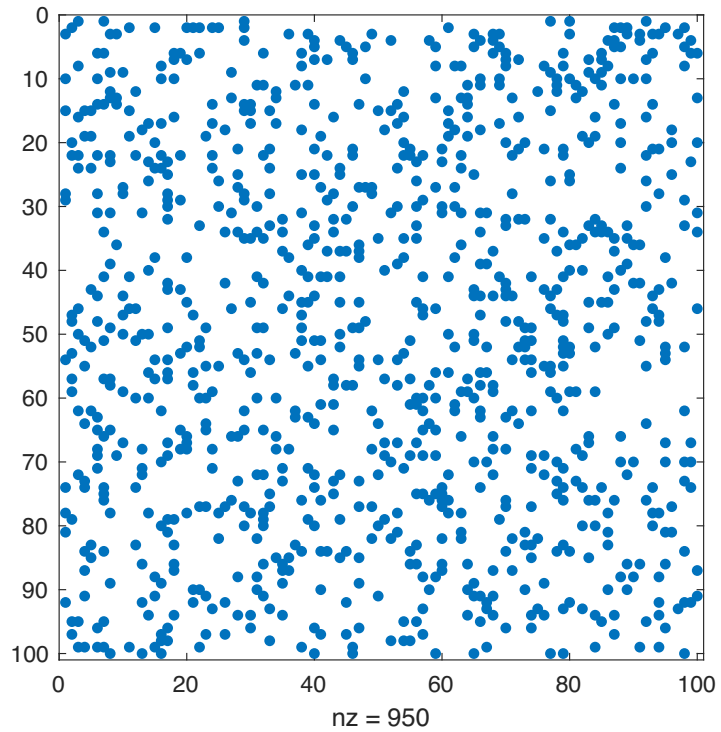
- ▶ Pivoting can help to reduce fill-in
- ▶ However, often “good” pivoting for sparsity leads to less stable behavior and vice versa

Sparse matrices in Matlab

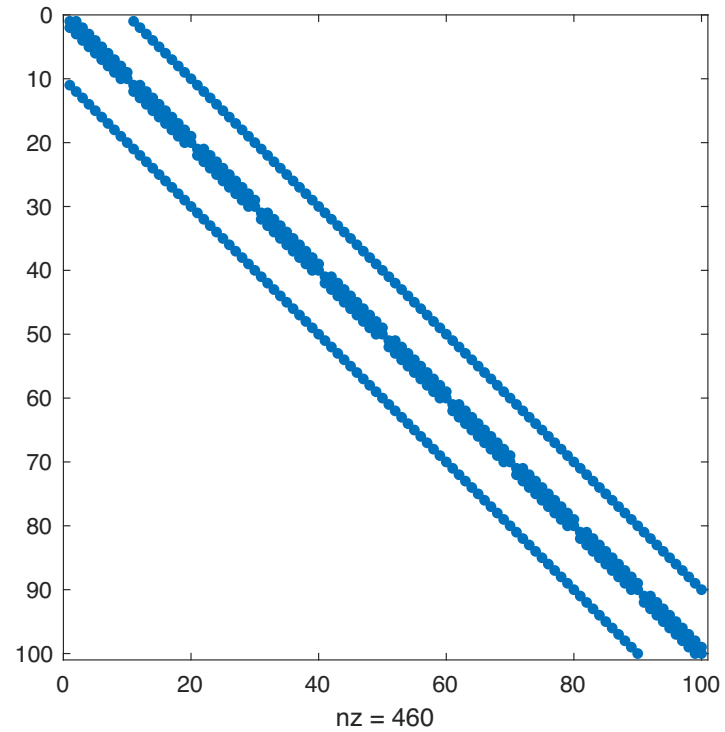
```
1: % S = sparse(i,j,v) generates a sparse matrix S
2: % from the triplets i, j, v with S(i(k),j(k)) = v(k).
3: >> A = sparse([1 2 2 4 4], [3 1 4 2 3], 1:5)
4: A =
5:      (2,1)          2
6:      (4,2)          4
7:      (1,3)          1
8:      (4,3)          5
9:      (2,4)          3
10: >> whos A
11:      Name          Size          Bytes    Class
12:      A             4x4             120     double    sparse
13: >> nnz(A)
14: ans =
15:      5
```

```
1: % S = sparse(i,j,v,m,n,nz) allocates space for nz
   nonzero elements.
2: % Use this syntax to allocate extra space for nonzero
   values to be filled in after construction.
3: >> A = sparse([], [], [], 4, 4, 5);
4: >> A(2, 1) = 2; A(4, 2) = 4; A(1, 3) = 1; A(4, 3) = 5;
   A(2, 4) = 3;
5: >> full(A)
6:
7: ans =
8:
9:      0      0      1      0
10:     2      0      0      3
11:     0      0      0      0
12:     0      4      5      0
```

```
1: % generate a random sparse matrix with density 10% and
    size 100x100
2: >> B = sprand(100, 100, 0.1);
3: % the sparse block tridiagonal matrix of order n^2
    resulting from discretizing Poisson's equation with
    the 5-point operator on an n-by-n mesh.
4: >> X = gallery('poisson', 10);
5: >> spy(B);
6: >> spy(X);
```

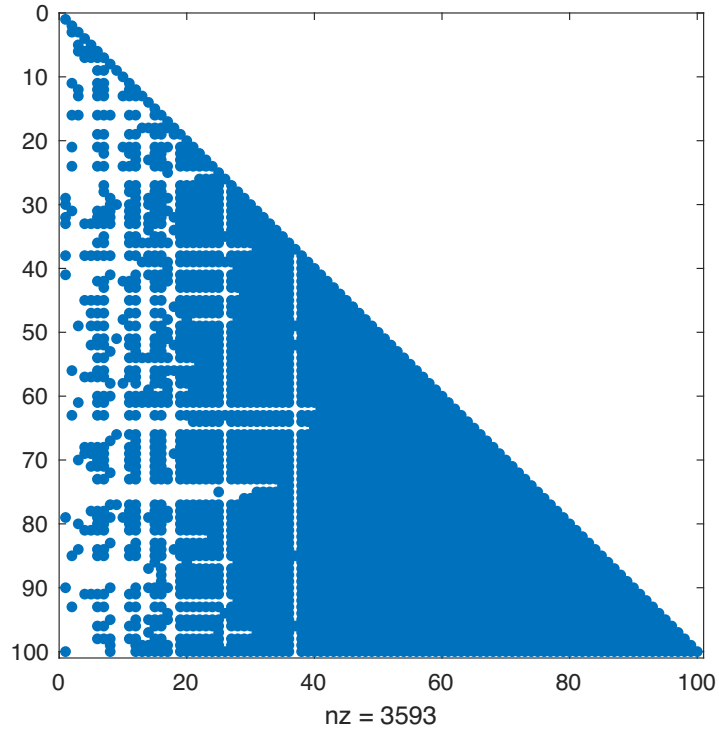


(a) matrix B

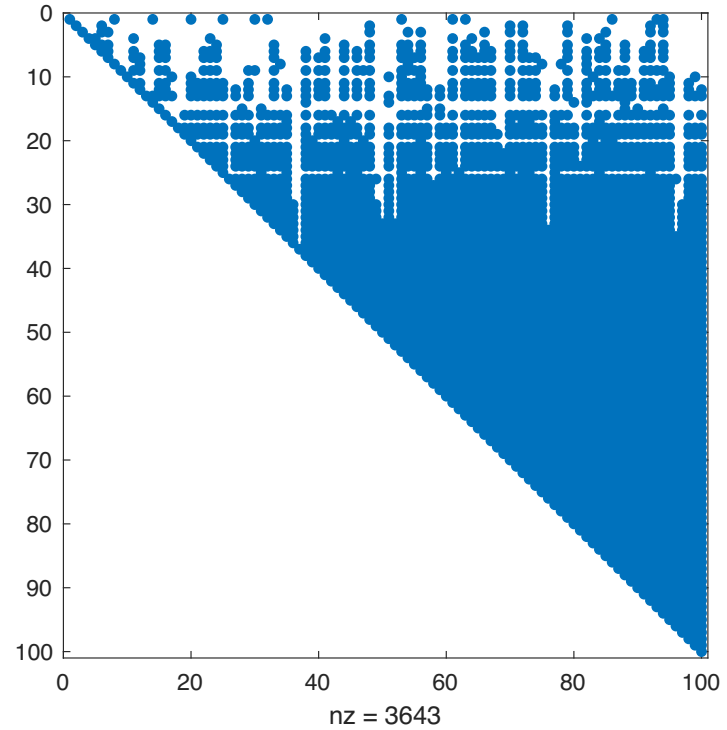


(b) matrix X

```
1: >> [L, U, P] = lu(B);
2: >> spy(L);
3: >> spy(U);
4:
5: >> [L, U, P] = lu(X);
6: >> spy(L);
7: >> spy(U);
```

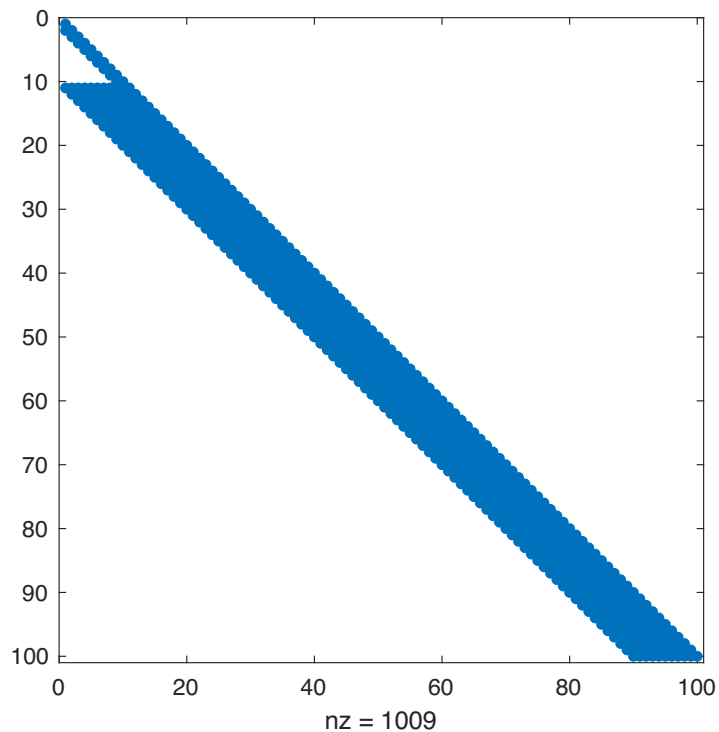


(a) factor L of B

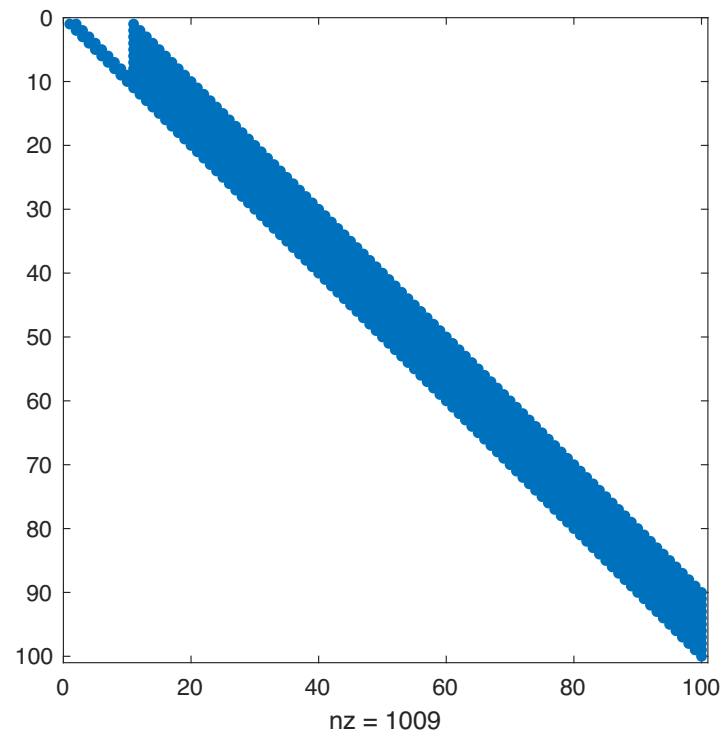


(b) factor U of B

A lot of fill-in! Factors L and U are not sparse, even though matrix B is sparse



(a) factor L of X



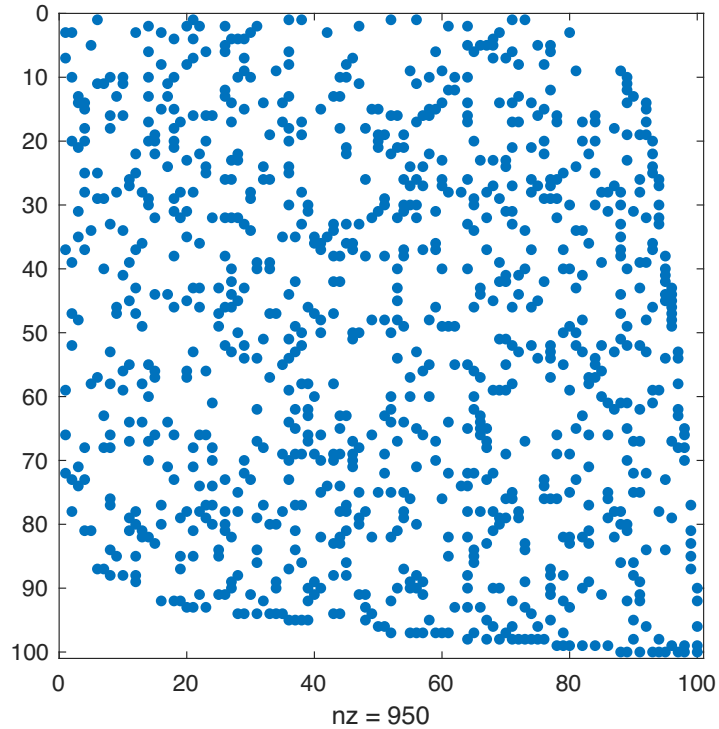
(b) factor U of X

Though better than for matrix B with random sparsity structure, there still are many more non-zero entries in the factors of X than in X itself.

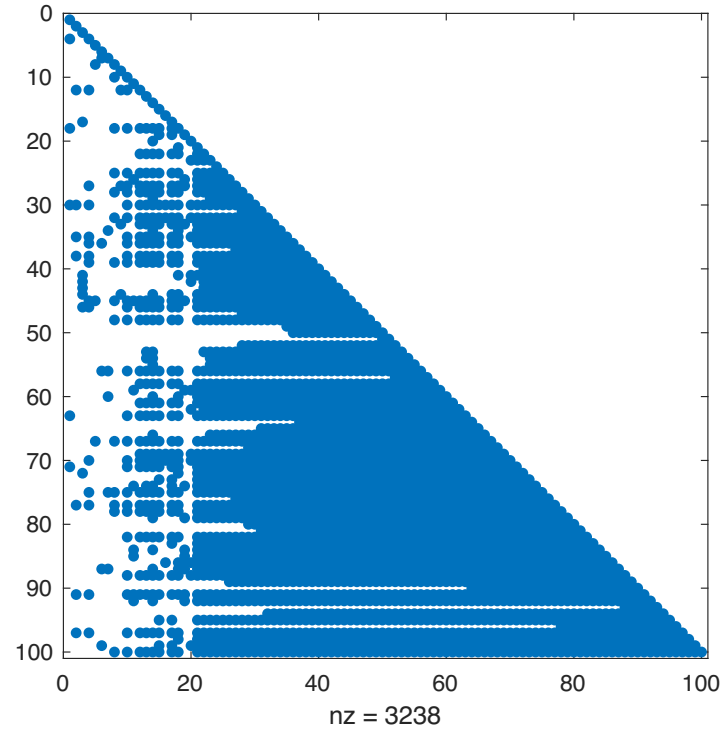
Changing the sparsity structure via re-ordering the matrix can help to reduce fill-in. For example, in Matlab the sparse reverse Cuthill-McKee ordering is implemented.

The re-ordered matrix tends to have its nonzero elements closer to the diagonal. This is a good reordering for LU or Cholesky factorization of matrices.

```
1: >> p = symrcm(B);
2: >> spy(B(p, p));
3: >> [L, U, P] = lu(B(p, p));
4: >> spy(L);
5: >> spy(U);
6:
7: >> p = symrcm(X)
8: >> spy(X(p, p));
9: >> [L, U, P] = lu(X(p, p));
10: >> spy(L);
11: >> spy(U);
```

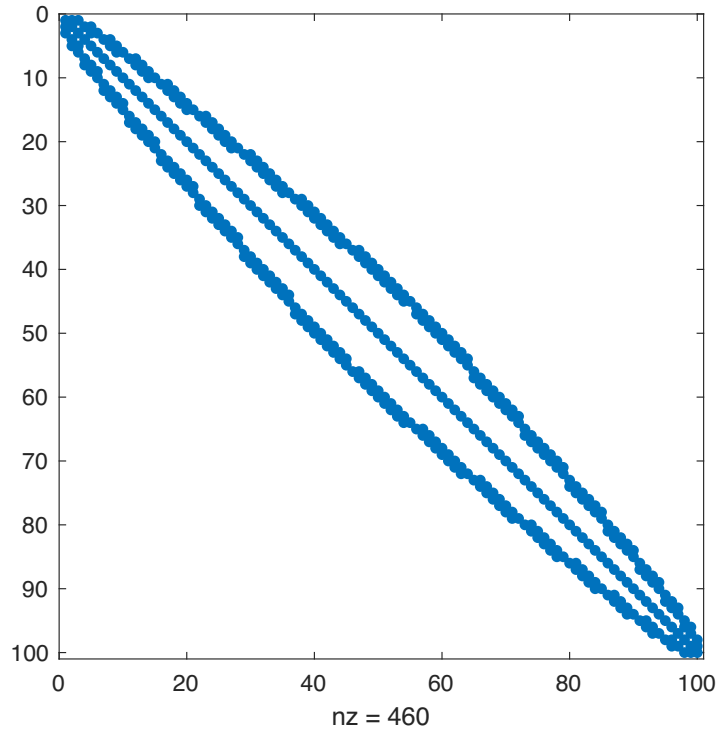


(a) re-ordered B matrix

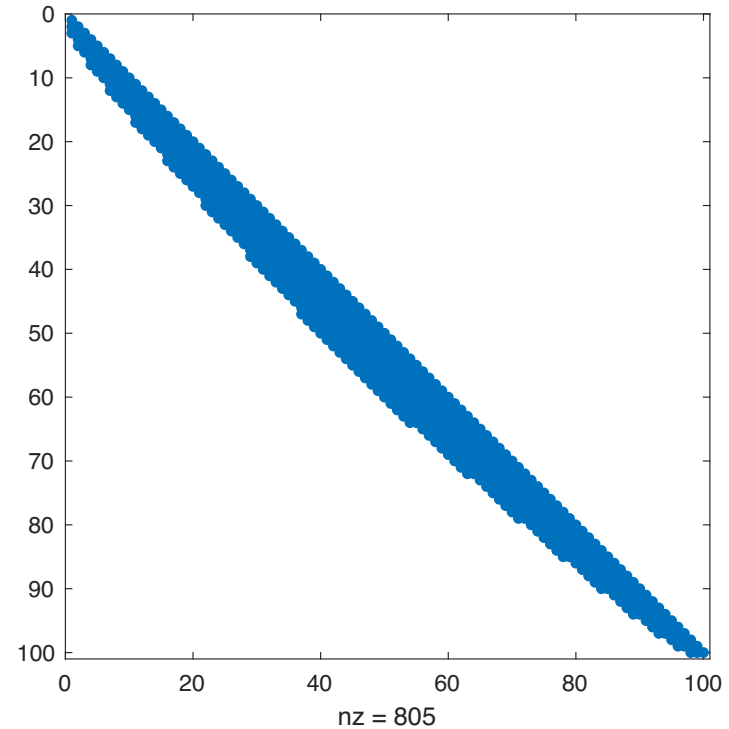


(b) factor L of re-ordered B

- ▶ Notice how the non-zero elements tend to be closer to the diagonal of the re-ordered B compared to the original B
- ▶ The fill-in is reduced from ~ 3500 to ~ 3200 non-zero entries; reduction of $< 10\%$
- ▶ It is hard to find a good ordering for matrix with a random sparsity structure



(a) re-ordered X matrix



(b) factor L of re-ordered X

- ▶ Non-zero entries of re-ordered X are closer to the diagonal than for the original X
- ▶ Fill-in is reduced by a roughly 20% from ~ 1000 to ~ 800 non-zero entries

Conclusions/summary

While there are general techniques for dealing with sparse matrices that help greatly, it all depends on the structure of the matrix

Pivoting has a dual, sometimes conflicting goal:

1. Reduce fill-in, i.e., improve memory use: Still active subject of research!
2. Reduce roundoff errors, i.e., improve stability. Typically some threshold pivoting is used only when needed

For many sparse matrices *iterative methods* (later) are required when large fill-in.

Feedback cards

Feedback cards

- ▶ Upload slides before class

Feedback cards

- ▶ Upload slides before class
- ▶ Neater handwriting and going slower when writing on board

Feedback cards

- ▶ Upload slides before class
- ▶ Neater handwriting and going slower when writing on board
- ▶ Define terms clearly instead of assuming everyone knows

Feedback cards

- ▶ Upload slides before class
- ▶ Neater handwriting and going slower when writing on board
- ▶ Define terms clearly instead of assuming everyone knows
- ▶ More intuitive explanations and motivations

Feedback cards

- ▶ Upload slides before class
- ▶ Neater handwriting and going slower when writing on board
- ▶ Define terms clearly instead of assuming everyone knows
- ▶ More intuitive explanations and motivations
- ▶ Show examples in Matlab rather than just on slides

Feedback cards

- ▶ Upload slides before class
- ▶ Neater handwriting and going slower when writing on board
- ▶ Define terms clearly instead of assuming everyone knows
- ▶ More intuitive explanations and motivations
- ▶ Show examples in Matlab rather than just on slides
- ▶ TA sessions

Feedback cards

- ▶ Upload slides before class
- ▶ Neater handwriting and going slower when writing on board
- ▶ Define terms clearly instead of assuming everyone knows
- ▶ More intuitive explanations and motivations
- ▶ Show examples in Matlab rather than just on slides
- ▶ TA sessions
- ▶ Code examples are interesting

Least-squares problems

Least-squares problems

Given **data points/measurements**

$$(t_i, b_i), \quad i = 1, \dots, m$$

and a **model function** ϕ that relates t and b :

$$b = \phi(t; x_1, \dots, x_n),$$

where x_1, \dots, x_n are model function parameters. If the model is supposed to describe the data, the deviations/errors

$$\Delta_i = b_i - \phi(t_i, x_1, \dots, x_n)$$

should be small. Thus, to fit the model to the measurements, one must choose x_1, \dots, x_n appropriately.

Least-squares problems

↪ visualization on board

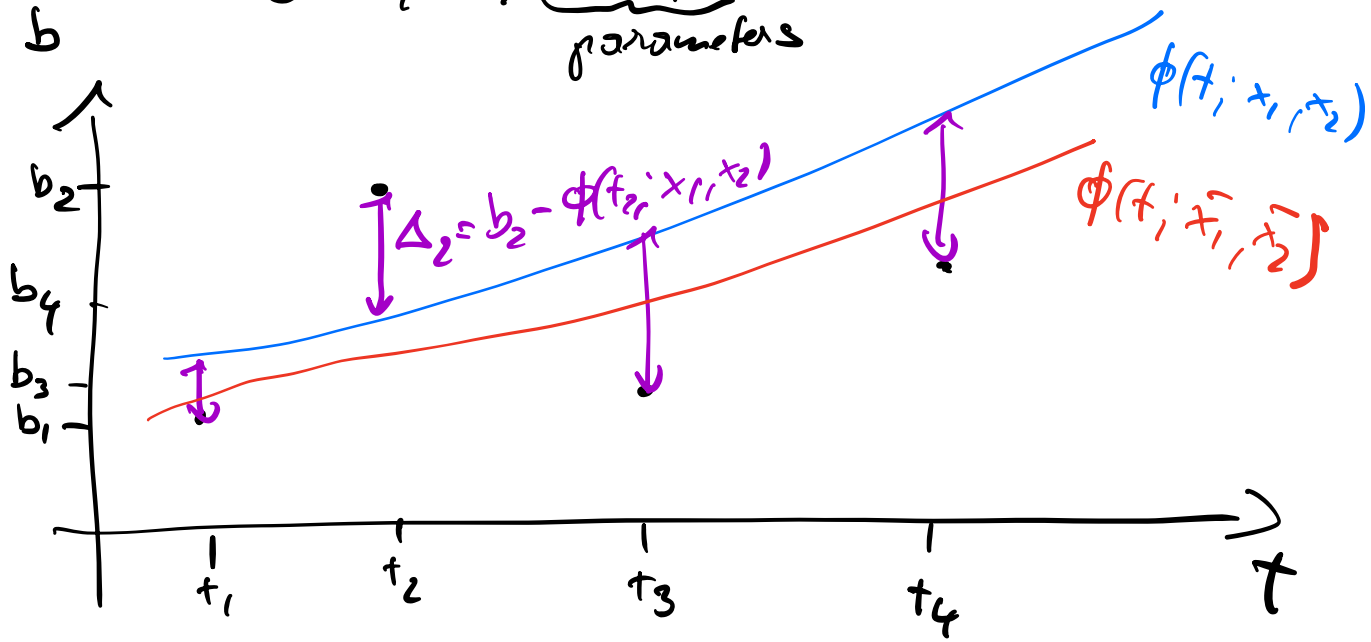
Least squares problem

data $\{(t_i, b_i)\}_{i=1}^m$

error $\Delta_i = b_i - \phi(t_i, x_1, \dots, x_n)$

parametrization

$$b = \phi(t; \underbrace{x_1, \dots, x_n}_{\text{parameters}})$$



Example:

- 1) $\phi(t; x_1, x_2) = x_1 t^2 + x_2 \exp(t)$
- 2) $\hat{\phi}(t; x_1, x_2) = x_1^2 t^2 + \log(x_2) t$

Least-squares problems

Least squares: Find x_1, \dots, x_n such that

$$\frac{1}{2} \sum_{i=1}^m \Delta_i^2 \rightarrow \min$$

Weighted least squares: Find x_1, \dots, x_n such that

$$\frac{1}{2} \sum_{i=1}^m \left(\frac{\Delta_i}{\delta b_i} \right)^2 \rightarrow \min,$$

where $\delta b_i > 0$ contain information about how much we trust the i th data point.

Least-squares problems (cont'd)

Alternatives to using squares:

L^1 error: Find x_1, \dots, x_n such that

$$\sum_{i=1}^m |\Delta_i| \rightarrow \min$$

Result can be very different, other statistical interpretation, more stable with respect to outliers.

L^∞ error: Find x_1, \dots, x_n such that

$$\max_{1 \leq i \leq m} |\Delta_i| \rightarrow \min$$

Keeps the worst-case error small (risk averse)

Numerical Methods I

MATH-GA 2010.001/CSCI-GA 2420.001

Benjamin Peherstorfer
Courant Institute, NYU

Based on slides by G. Stadler and A. Donev

Today

Last time

- ▶ Cost analysis of LU decomposition
- ▶ Solving linear systems with sparse matrices
- ▶ Least-squares problems

Today

- ▶ Least-squares problems

Announcements

- ▶ Homework 2 has been posted; is due next week Mon, Oct 7 *before class*

Recap: Least-squares problems

Given **data points/measurements**

$$(t_i, b_i), \quad i = 1, \dots, m$$

and a **model function** ϕ that relates t and b :

$$b = \phi(t; x_1, \dots, x_n),$$

where x_1, \dots, x_n are model function parameters. If the model is supposed to describe the data, the deviations/errors

$$\Delta_i = b_i - \phi(t_i, x_1, \dots, x_n)$$

should be small. Thus, to fit the model to the measurements, one must choose x_1, \dots, x_n appropriately.

Recap: Least-squares problems

Least squares: Find x_1, \dots, x_n such that

$$\frac{1}{2} \sum_{i=1}^m \Delta_i^2 \rightarrow \min$$

Weighted least squares: Find x_1, \dots, x_n such that

$$\frac{1}{2} \sum_{i=1}^m \left(\frac{\Delta_i}{\delta b_i} \right)^2 \rightarrow \min,$$

where $\delta b_i > 0$ contain information about how much we trust the i th data point.

Linear least-squares

We assume (for now) that the **model depends linearly on x_1, \dots, x_n** , e.g.:

$$\phi(t; x_1, \dots, x_n) = a_1(t)x_1 + \dots + a_n(t)x_n$$

↪ board

$$\phi(t; x_1, x_2) = x_1 t^2 + x_2 \exp(t) \quad ; \text{ data points } (t_i, b_i), i=1 \dots 3$$

$$\Delta_1 = b_1 - (x_1 t_1^2 + x_2 \exp(t_1))$$

$$\Delta_2 = b_2 - (x_1 t_2^2 + x_2 \exp(t_2))$$

$$\Delta_3 = b_3 - (x_1 t_3^2 + x_2 \exp(t_3))$$

$$\min_{x_1, x_2} \sum_{i=1}^3 \Delta_i^2$$

2 unknowns
3 equations

Linear least-squares

Choosing the least square error, this results in

$$\min_{\mathbf{x}} \|\mathbf{Ax} - \mathbf{b}\|^2,$$

where $\mathbf{x} = (x_1, \dots, x_n)^T$, $\mathbf{b} = (b_1, \dots, b_m)^T$, and $a_{ij} = a_j(t_i)$.

In the following, we study the **overdetermined case**, i.e., $m \geq n \rightsquigarrow$ board

$$A = \begin{bmatrix} t_1^2 & \exp(t_1) \\ t_2^2 & \exp(t_2) \\ t_3^2 & \exp(t_3) \end{bmatrix}, \quad \mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}, \quad \mathbf{b} = \begin{bmatrix} b_1 \\ b_2 \\ b_3 \end{bmatrix}$$

$\in \mathbb{R}^{3 \times 2} \qquad \qquad \in \mathbb{R}^2 \qquad \qquad \in \mathbb{R}^3$

$m = 3, \quad n = 2$

Linear least-squares

Different perspective:

Consider non-square matrices $A \in \mathbb{R}^{m \times n}$ with $m \geq n$ and $\text{rank}(A) = n$. Then the system

$$A\mathbf{x} = \mathbf{b}$$

does not necessarily have a solution (more equations than unknowns). We thus instead solve a minimization problem

$$\min_{\mathbf{x}} \|A\mathbf{x} - \mathbf{b}\|^2 = \min_{\mathbf{x}} \Phi(\mathbf{x})$$

How can we solve this optimization problem?

Because we consider the Euclidean norm $\|\cdot\|_2$, we obtain

$$\Phi(\mathbf{x}) = (\mathbf{Ax} - \mathbf{b})^T (\mathbf{Ax} - \mathbf{b}) = \mathbf{x}^T \mathbf{A}^T \mathbf{Ax} - \dots$$

which is quadratic in \mathbf{x} if \mathbf{A} has full rank \rightsquigarrow convex in \mathbf{x}

Therefore, the critical point is the global optimum

$$\nabla\Phi(\mathbf{x}) = \mathbf{A}^T (2(\mathbf{Ax} - \mathbf{b})) = \mathbf{0}$$

which satisfies the *normal equations*

$$\mathbf{A}^T \mathbf{Ax} = \mathbf{A}^T \mathbf{b}.$$

If \mathbf{A} is full rank, $\text{rank}(\mathbf{A}) = n$, then $\mathbf{A}^T \mathbf{A}$ is positive definite and the normal equations can be solved with the Cholesky factorization

Because we consider the Euclidean norm $\|\cdot\|_2$, we obtain

$$\Phi(\mathbf{x}) = (\mathbf{Ax} - \mathbf{b})^T (\mathbf{Ax} - \mathbf{b}) = \mathbf{x}^T \mathbf{A}^T \mathbf{Ax} - \dots$$

which is quadratic in \mathbf{x} if \mathbf{A} has full rank \rightsquigarrow convex in \mathbf{x}

Therefore, the critical point is the global optimum

$$\nabla\Phi(\mathbf{x}) = \mathbf{A}^T (2(\mathbf{Ax} - \mathbf{b})) = \mathbf{0}$$

which satisfies the *normal equations*

$$\mathbf{A}^T \mathbf{Ax} = \mathbf{A}^T \mathbf{b}.$$

If \mathbf{A} is full rank, $\text{rank}(\mathbf{A}) = n$, then $\mathbf{A}^T \mathbf{A}$ is positive definite and the normal equations can be solved with the Cholesky factorization (warning: we shouldn't do this, can you already see why?)

A geometry perspective on the normal equations

$$\min_{x \in \mathbb{R}^n} \|Ax - b\|_2^2, \quad A \in \mathbb{R}^{m \times n}, \quad m > n, \quad \text{rank}(A) = n$$
$$b \in \mathbb{R}^m$$

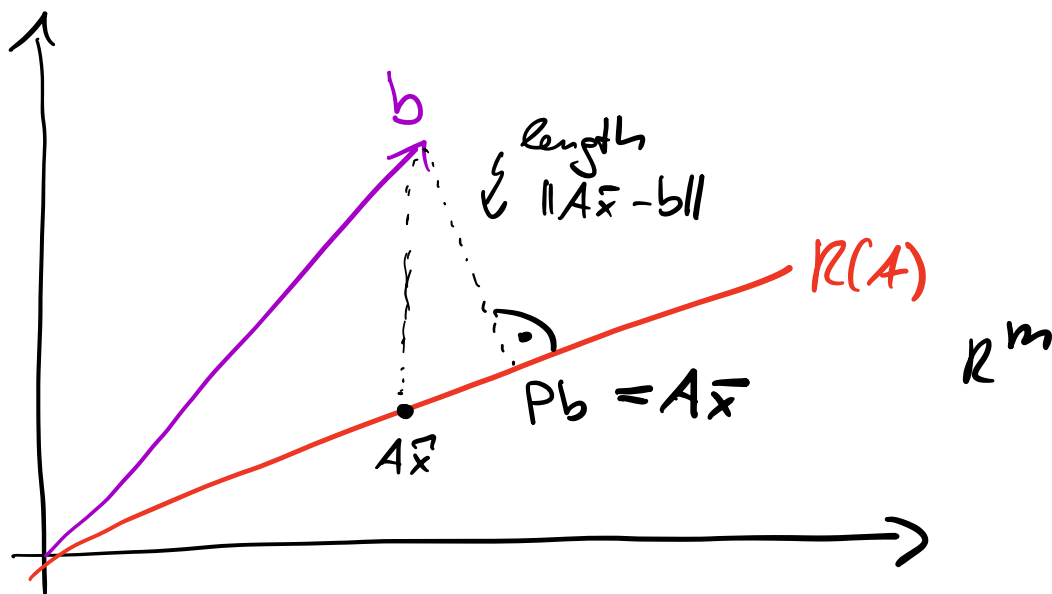
Geometric sketch

$$\bar{x} = \arg \min_{x \in \mathbb{R}^n} \|Ax - b\|_2^2$$



$$A\bar{x} - b \perp \mathcal{R}(A)$$

$$\| \{Ax \mid x \in \mathbb{R}^n\}$$



Orthogonal projection

Consider finite-dim space V with $\langle \cdot, \cdot \rangle$ and subspace $U \subset V$. Let U^\perp be the orthogonal complement w.r.t. V :

$$U^\perp = \{v \in V \mid \langle v, u \rangle = 0 \quad \forall u \in U\}$$

Then

$$\arg \min_{u' \in U} \|u' - v\| = u \iff v - u \in U^\perp$$

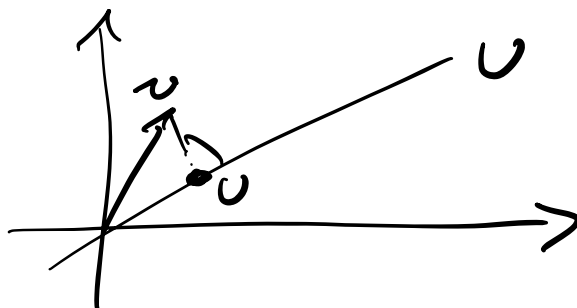
Proof:

" \Leftarrow ": Let $u \in U$ such that $v - u \in U^\perp$, then for all $u' \in U$

$$\begin{aligned} \|v - u'\|^2 &= \|v - u + u - u'\|^2 = \\ &= \|v - u\|^2 + \underbrace{2\langle \quad \rangle}_{=0} + \|u - u'\|^2 \\ &= \|v - u\|^2 + \|u - u'\|^2 \geq \|v - u\|^2 \end{aligned} \quad \Bigg| \begin{array}{l} \Rightarrow \\ \end{array}$$

This unique $u \in U$ is called the orthogonal projection of v onto U

$$u = P(v)$$



P is linear. Thus (finite-dim), there exists matrix P such that

$$u = \underline{P}v$$

Now apply geometric interpretation to find normal equations

$$\bar{x} = \arg \min_{x \in \mathbb{R}^n} \|Ax - b\|_2^2$$

$$V = \mathbb{R}^m, \quad U = R(A) \subset V$$

$$\|b - Ax\|^2 = \min \iff \langle b - Ax, Ax' \rangle = 0 \quad \forall x' \in \mathbb{R}^n$$

$$\iff \langle A^T(b - Ax), x' \rangle = 0 \quad \forall x' \in \mathbb{R}^n$$

$$\iff A^T Ax = A^T b$$

$$\iff x = (A^T A)^{-1} A^T b$$

$A(A^T A)^{-1} A^T \dots$ projector onto column space of A

$$\text{solve: } Ax = Pb = A(A^T A)^{-1} A^T b \iff x = (A^T A)^{-1} A^T b$$

Linear least-squares problems

Solving the **normal equations**

$$A^T A \bar{\mathbf{x}} = A^T \mathbf{b}$$

requires:

- ▶ computing $A^T A$ (which is $O(mn^2)$)
- ▶ **condition number of $A^T A$** ? \rightsquigarrow board

Condition number of $A^T A$ in $\|\cdot\|_2$ norm is

$$\kappa_2(A^T A) = \kappa_2(A)^2$$

Proof:

$$\begin{aligned} \kappa_2(A)^2 &= \|A^{-1}\|_2^2 \|A\|_2^2 = \frac{\max_{\|x\|_2=1} \|Ax\|_2^2}{\min_{\|x\|_2=1} \|Ax\|_2^2} \\ &= \frac{\max_{\|x\|_2=1} \langle Ax, Ax \rangle_2}{\min_{\|x\|_2=1} \langle Ax, Ax \rangle_2} \\ &= \frac{\max_{\|x\|_2=1} \langle A^T A x, x \rangle_2}{\min_{\|x\|_2=1} \langle A^T A x, x \rangle_2} \\ &= \frac{\lambda_{\max}(A^T A)}{\lambda_{\min}(A^T A)} \end{aligned}$$

Because $A^T A$ is spd, have

$$\begin{aligned} \lambda_{\max}(A^T A) &= \sqrt{\lambda_{\max}(A^T A)^T} = \sqrt{\lambda_{\max}(A^T A)^2} \\ &= \sqrt{\lambda_{\max}((A^T A)^T (A^T A))} \\ &= \|A^T A\|_2 \end{aligned}$$

analogous: $\frac{1}{\lambda_{\min}(A^T A)} = \|(A^T A)^{-1}\|_2$

$$\kappa_2(A)^2 = \|A^T A\|_2 \|(A^T A)^{-1}\|_2 = \kappa_2(A^T A)$$

Because $\kappa_2(A^T A) = \kappa_2(A)^2$ we should not compute with $A^T A$ even

Condition of orthogonal projection

Lemma:

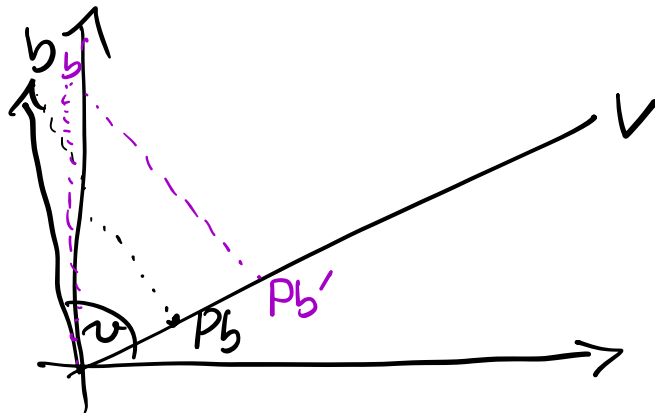
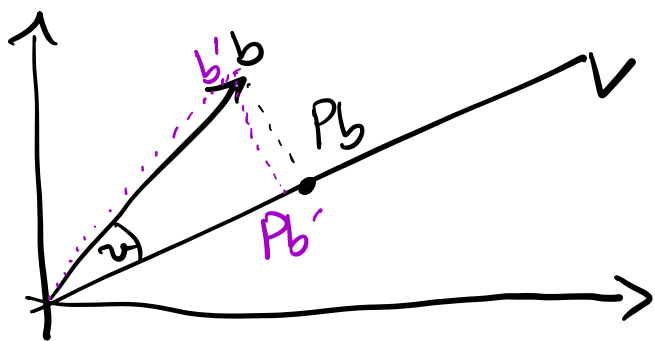
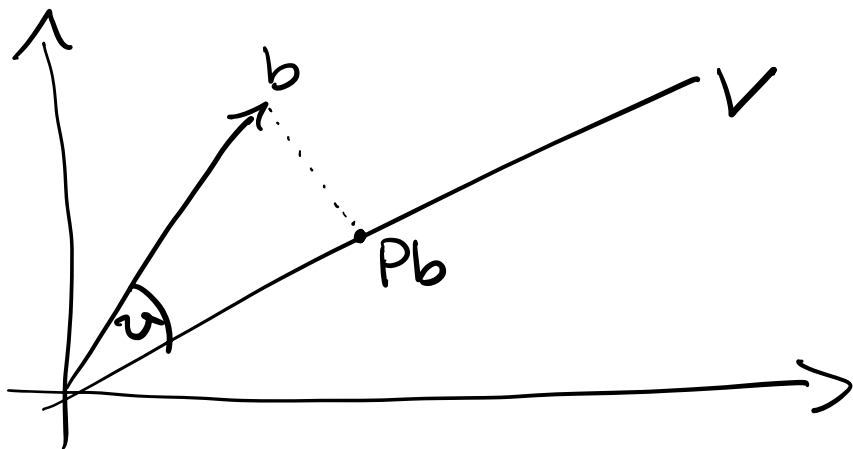
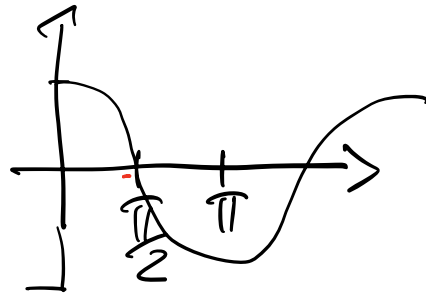
$P: \mathbb{R}^m \rightarrow V$ orthogonal projection onto $V \subseteq \mathbb{R}^m$

For $b \in \mathbb{R}^m$, denote by ν angle between b and V defined by

$$\sin(\nu) = \frac{\|b - Pb\|_2}{\|b\|_2}$$

Then the relative condition number of (P, b) w.r.t. 2-norm is

$$\kappa_{\text{rel}} = \frac{1}{\cos(\nu)} \|P\|_2$$



Linear least-squares problems

Solving the **normal equations**

$$A^T A \bar{\mathbf{x}} = A^T \mathbf{b}$$

requires:

- ▶ computing $A^T A$ (which is $O(mn^2)$)
- ▶ **condition number of $A^T A$** ? \rightsquigarrow board is square of condition number of A ;
(problematic for the Choleski factorization)

Numerical Methods I

MATH-GA 2010.001/CSCI-GA 2420.001

Benjamin Peherstorfer
Courant Institute, NYU

Based on slides by G. Stadler and A. Donev

Today

Last time

- ▶ Least-squares problems

Today

- ▶ Least-squares problems

Announcements

- ▶ Homework 2 has been posted; is due next week Mon, Oct 7 *before class*

Recap: Linear least-squares

Consider non-square matrices $A \in \mathbb{R}^{m \times n}$ with $m \geq n$ and $\text{rank}(A) = n$. Then the system

$$A\mathbf{x} = \mathbf{b}$$

does not necessarily have a solution (more equations than unknowns). We thus instead solve a minimization problem

$$\min_{\mathbf{x}} \|A\mathbf{x} - \mathbf{b}\|^2 = \min_{\mathbf{x}} \Phi(\mathbf{x})$$

Recap: Linear least-squares

Consider non-square matrices $A \in \mathbb{R}^{m \times n}$ with $m \geq n$ and $\text{rank}(A) = n$. Then the system

$$A\mathbf{x} = \mathbf{b}$$

does not necessarily have a solution (more equations than unknowns). We thus instead solve a minimization problem

$$\min_{\mathbf{x}} \|A\mathbf{x} - \mathbf{b}\|^2 = \min_{\mathbf{x}} \Phi(\mathbf{x})$$

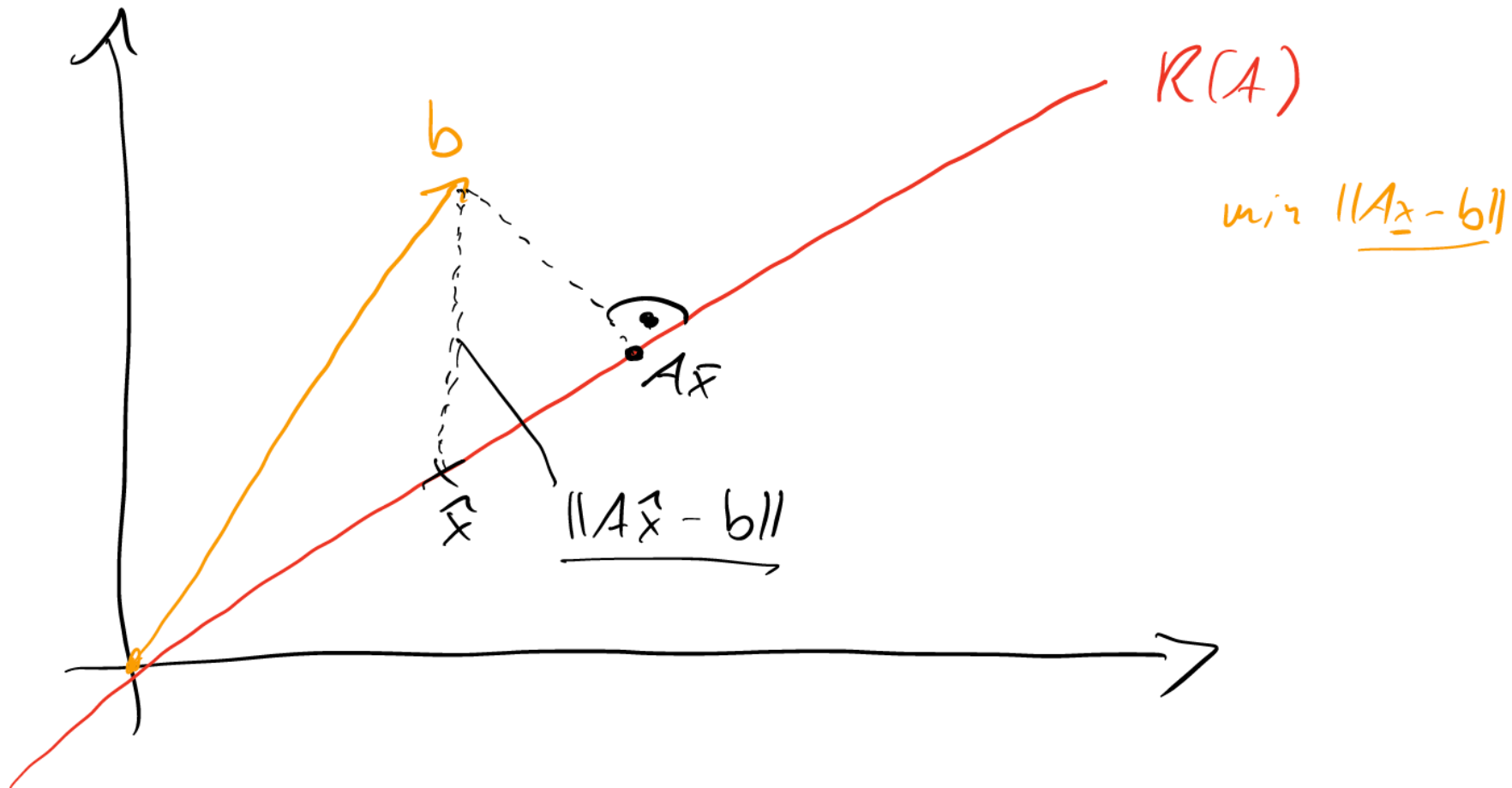
Solving the **normal equations**

$$A^T A \bar{\mathbf{x}} = A^T \mathbf{b}$$

requires:

- ▶ computing $A^T A$ (which is $O(mn^2)$)
- ▶ **condition number of $A^T A$** ? \rightsquigarrow is square of condition number of A ; (problematic for the Choleski factorization)

Recap: Least-squares problems



Recap: Linear least-squares problems

Conditioning

Solving the normal equation is equivalent to computing $P\mathbf{b}$, the **orthogonal projection of \mathbf{b} onto the subspace V spanned by columns of A .**

Let $P : \mathbb{R}^m \rightarrow V$ be an orthogonal projection onto $V \subseteq \mathbb{R}^n$. For $\mathbf{b} \in \mathbb{R}^m$, denote by θ the angle between \mathbf{b} and V defined by

$$\sin(\theta) = \frac{\|\mathbf{b} - P\mathbf{b}\|_2}{\|\mathbf{b}\|_2}.$$

The relative condition number of projecting \mathbf{b} onto V with P with respect to the 2-norm (\mathbf{b} is input) is

$$\kappa_{\text{rel}}(\mathbf{b}) = \frac{1}{\cos(\theta(\mathbf{b}))} \|P\|_2.$$

↪ board

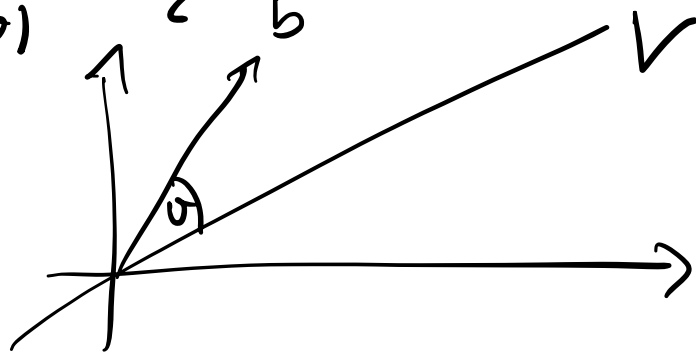
$P: \mathbb{R}^m \rightarrow V$ orthogonal projector $V \subseteq \mathbb{R}^m$

For $b \in \mathbb{R}^m$, denote angle ν between b and V that is defined as

$$\sin(\nu) = \frac{\|b - Pb\|_2}{\|b\|_2}$$

Then the relative condition number of (P, b) w.r.f. 2-norm is

$$\kappa_{\text{rel}} = \frac{1}{\cos(\nu)} \|P\|_2$$



Def: $f: X \rightarrow Y$ differentiable

$$\kappa_{\text{rel}} = \frac{\|x\|}{\|f(x)\|} \|f'(x)\|$$

Projection: P is linear \Rightarrow differentiable

$$\kappa_{\text{rel}}(b) = \frac{\|b\|}{\|Pb\|} \|P'(b)\|$$

(1) $P'(b) = P$

(2) $\frac{\|Pb\|^2}{\|b\|^2}$

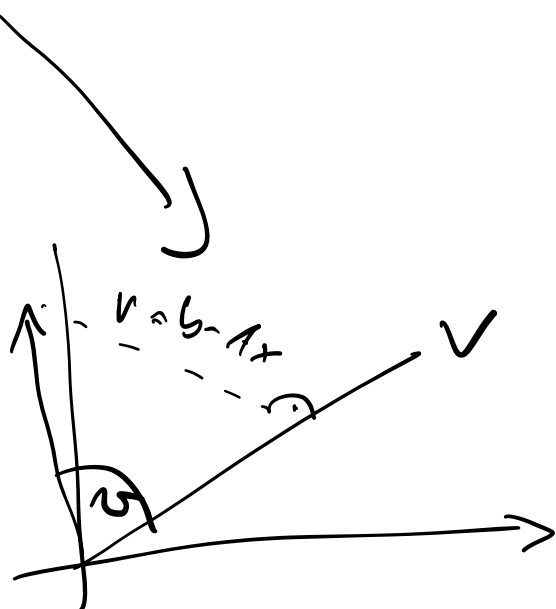
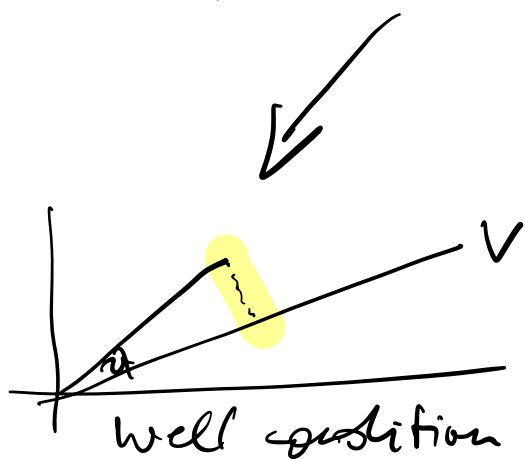
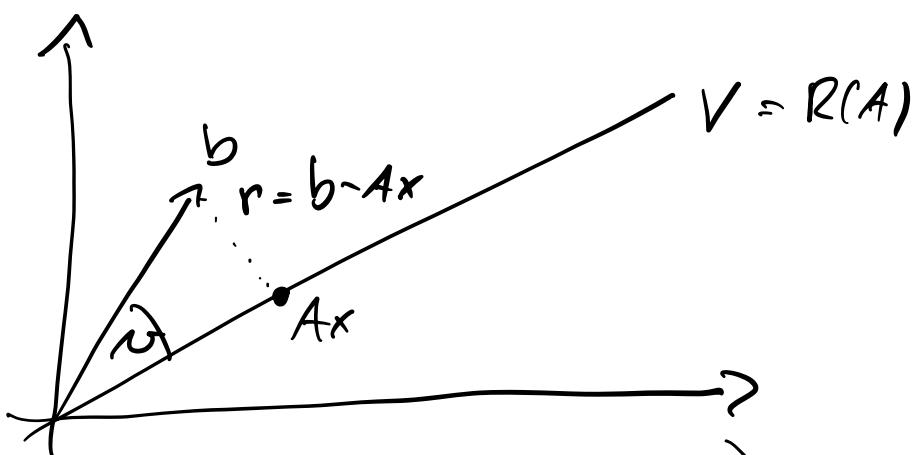
$$b - Pb \perp Pb$$

$$\|Pb\|^2 = \|b\|^2 - \|b - Pb\|^2$$

$$\frac{\|Pb\|^2}{\|b\|^2} = \frac{\|b\|^2 - \|b - Pb\|^2}{\|b\|^2}$$

$$= (1 - \sin^2(\nu)) = \cos^2(\nu)$$

$$P_{\text{re}}(b) = \frac{1}{\cos(\nu)} \|P\|_2$$



Linear least-squares problems

Now for the least-squares problem $\|\mathbf{Ax} - \mathbf{b}\|_2$. The relative condition number κ in the Euclidean norm is bounded by

- ▶ With respect to perturbations in \mathbf{b} :

$$\kappa \leq \frac{\kappa_2(A)}{\cos(\theta)}$$

- ▶ With respect to perturbations in \mathbf{A} :

$$\kappa \leq \kappa_2(A) + \kappa_2(A)^2 \tan(\theta)$$

Proof \rightsquigarrow next week

What are these bounds telling us?

Linear least-squares problems

Now for the least-squares problem $\|\mathbf{Ax} - \mathbf{b}\|_2$. The relative condition number κ in the Euclidean norm is bounded by

- ▶ With respect to perturbations in \mathbf{b} :

$$\kappa \leq \frac{\kappa_2(A)}{\cos(\theta)}$$

- ▶ With respect to perturbations in \mathbf{A} :

$$\kappa \leq \kappa_2(A) + \kappa_2(A)^2 \tan(\theta)$$

Proof \rightsquigarrow next week

What are these bounds telling us?

Small residual problems, small angle θ $\cos(\theta) \approx 1$, $\tan(\theta) \approx 0$: behavior similar to linear system.

Large residual problems, large angle θ $\cos(\theta) \ll 1$, $\tan(\theta) \approx 1$: behavior very different from linear system because $\kappa_2(A)^2$ shows up

How should we solve least-squares problems numerically?

We know from the previous slide that if the residual is large, then the condition κ is much larger than $\kappa_2(A)$ (closer to $\kappa_2(A)^2$)

- ▶ This is a poorly condition problem; however, do we care?

How should we solve least-squares problems numerically?

We know from the previous slide that if the residual is large, then the condition κ is much larger than $\kappa_2(A)$ (closer to $\kappa_2(A)^2$)

- ▶ This is a poorly condition problem; however, do we care?
- ▶ If the residual is large, then our \mathbf{Ax} won't explain well the right-hand side \mathbf{b}
- ▶ This means that “our curve doesn't fit well the data” and we probably should try to find another space in which to search for a solution (another \mathbf{A} with a range that better approximates the projected right-hand side \mathbf{b})

More relevant is the situation with a small residual and then $\kappa \approx \kappa_2(A)$

- ▶ Here we have a well condition problem; so what could go wrong?

More relevant is the situation with a small residual and then $\kappa \approx \kappa_2(\mathbf{A})$

- ▶ Here we have a well condition problem; so what could go wrong?
- ▶ If we choose a numerical method that solves the normal equations

$$\mathbf{A}^T \mathbf{A} \mathbf{x} = \mathbf{A}^T \mathbf{b}$$

then our problem becomes the problem of solving the linear system with matrix $\mathbf{A}^T \mathbf{A}$, which has condition number

$$\kappa_2(\mathbf{A}^T \mathbf{A}) = \kappa_2(\mathbf{A})^2$$

More relevant is the situation with a small residual and then $\kappa \approx \kappa_2(\mathbf{A})$

- ▶ Here we have a well condition problem; so what could go wrong?
- ▶ If we choose a numerical method that solves the normal equations

$$\mathbf{A}^T \mathbf{A} \mathbf{x} = \mathbf{A}^T \mathbf{b}$$

then our problem becomes the problem of solving the linear system with matrix $\mathbf{A}^T \mathbf{A}$, which has condition number

$$\kappa_2(\mathbf{A}^T \mathbf{A}) = \kappa_2(\mathbf{A})^2$$

\rightsquigarrow we are back in the situation of a poorly condition problem (“solving a linear system with $\mathbf{A}^T \mathbf{A}$ ”) even though our original problem (least-squares problem) is well condition

- ▶ Can we do better and solve the least-squares problem (the problem we are actually interested in) without having to solve a problem with condition that grows with $\kappa_2(\mathbf{A})^2$ on the way?

The QR decomposition

Recall that projecting \mathbf{b} onto the column span (range) of \mathbf{A} was the key step \rightsquigarrow let's try to find a numerical method that computes an orthonormal basis $\mathbf{q}_1, \dots, \mathbf{q}_n$ of the rank- n column span of \mathbf{A}

$$\mathbf{A} = \begin{bmatrix} | & & | \\ \mathbf{a}_1 & \dots & \mathbf{a}_n \\ | & & | \end{bmatrix} \in \mathbb{R}^{m \times n}, \quad m \geq n$$

$$\Downarrow$$

$$\underbrace{\begin{bmatrix} | & & | \\ \mathbf{a}_1 & \dots & \mathbf{a}_n \\ | & & | \end{bmatrix}}_{\mathbf{A}} = \underbrace{\begin{bmatrix} | & & | \\ \mathbf{q}_1 & \dots & \mathbf{q}_n \\ | & & | \end{bmatrix}}_{\mathbf{Q}} \underbrace{\begin{bmatrix} r_{11} & r_{12} & \dots & r_{1n} \\ & r_{22} & & \vdots \\ & & \ddots & \\ & & & r_{nn} \end{bmatrix}}_{\mathbf{R}}$$

with an invertible matrix \mathbf{R} so that

$$\text{span}(\mathbf{a}_1, \dots, \mathbf{a}_k) = \text{span}(\mathbf{q}_1, \dots, \mathbf{q}_k), \quad k = 1, \dots, n$$

$$\underbrace{\begin{bmatrix} | & & | \\ \mathbf{a}_1 & \dots & \mathbf{a}_n \\ | & & | \end{bmatrix}}_A = \underbrace{\begin{bmatrix} | & & | \\ \mathbf{q}_1 & \dots & \mathbf{q}_n \\ | & & | \end{bmatrix}}_Q \underbrace{\begin{bmatrix} r_{11} & r_{12} & \dots & r_{1n} \\ & r_{22} & & \dots \\ & & \ddots & \\ & & & r_{nn} \end{bmatrix}}_R$$

⇓ leads to system of equations ⇓

$$\mathbf{a}_1 = r_{11} \mathbf{q}_1$$

$$\mathbf{a}_2 = r_{12} \mathbf{q}_1 + r_{22} \mathbf{q}_2$$

$$\mathbf{a}_3 = r_{13} \mathbf{q}_1 + r_{23} \mathbf{q}_2 + r_{33} \mathbf{q}_3$$

⋮

$$\mathbf{a}_n = r_{1n} \mathbf{q}_1 + r_{2n} \mathbf{q}_2 + \dots + r_{nn} \mathbf{q}_n$$

What process does this motivate?

This motivates a process for computing the basis $\mathbf{q}_1, \dots, \mathbf{q}_n$

- ▶ At step j , we have $\mathbf{q}_1, \dots, \mathbf{q}_{j-1}$ that span $\text{span}(\mathbf{a}_1, \dots, \mathbf{a}_{j-1})$
- ▶ We want to find \mathbf{q}_j orthonormal to $\mathbf{q}_1, \dots, \mathbf{q}_{j-1}$ so that $\mathbf{q}_1, \dots, \mathbf{q}_j$ spans $\text{span}(\mathbf{a}_1, \dots, \mathbf{a}_j)$
- ▶ Thus, set

$$\mathbf{v}_j = \mathbf{a}_j - (\mathbf{q}_1^T \mathbf{a}_j) \mathbf{q}_1 - (\mathbf{q}_2^T \mathbf{a}_j) \mathbf{q}_2 - \dots - (\mathbf{q}_{j-1}^T \mathbf{a}_j) \mathbf{q}_{j-1}$$

and normalize

$$\mathbf{q}_j = \frac{\mathbf{v}_j}{\|\mathbf{v}_j\|_2}$$

Notice that at step j , the quantities $\mathbf{q}_1^T \mathbf{a}_j, \mathbf{q}_2^T \mathbf{a}_j, \dots, \mathbf{q}_{j-1}^T \mathbf{a}_j$ are the values $r_{j,1}, \dots, r_{j,j-1}$ and r_{jj} is responsible for the normalization and set to

$$r_{jj} = \left\| \mathbf{a}_j - \sum_{i=1}^{j-1} r_{ij} \mathbf{q}_i \right\|_2$$

This process is the *classical Gram-Schmidt* procedure to compute the QR factorization

This motivates a process for computing the basis $\mathbf{q}_1, \dots, \mathbf{q}_n$

- ▶ At step j , we have $\mathbf{q}_1, \dots, \mathbf{q}_{j-1}$ that span $\text{span}(\mathbf{a}_1, \dots, \mathbf{a}_{j-1})$
- ▶ We want to find \mathbf{q}_j orthonormal to $\mathbf{q}_1, \dots, \mathbf{q}_{j-1}$ so that $\mathbf{q}_1, \dots, \mathbf{q}_j$ spans $\text{span}(\mathbf{a}_1, \dots, \mathbf{a}_j)$
- ▶ Thus, set

$$\mathbf{v}_j = \mathbf{a}_j - (\mathbf{q}_1^T \mathbf{a}_j) \mathbf{q}_1 - (\mathbf{q}_2^T \mathbf{a}_j) \mathbf{q}_2 - \dots - (\mathbf{q}_{j-1}^T \mathbf{a}_j) \mathbf{q}_{j-1}$$

and normalize

$$\mathbf{q}_j = \frac{\mathbf{v}_j}{\|\mathbf{v}_j\|_2}$$

Notice that at step j , the quantities $\mathbf{q}_1^T \mathbf{a}_j, \mathbf{q}_2^T \mathbf{a}_j, \dots, \mathbf{q}_{j-1}^T \mathbf{a}_j$ are the values $r_{j,1}, \dots, r_{j,j-1}$ and r_{jj} is responsible for the normalization and set to

$$r_{jj} = \left\| \mathbf{a}_j - \sum_{i=1}^{j-1} r_{ij} \mathbf{q}_i \right\|_2$$

This process is the *classical Gram-Schmidt* procedure to compute the QR factorization
However, this process is numerically unstable!

Instead of directly computing

$$\mathbf{v}_j = \mathbf{a}_j - (\mathbf{q}_1^T \mathbf{a}_j) \mathbf{q}_1 - (\mathbf{q}_2^T \mathbf{a}_j) \mathbf{q}_2 - \cdots - (\mathbf{q}_{j-1}^T \mathbf{a}_j) \mathbf{q}_{j-1}$$

Instead of directly computing

$$\mathbf{v}_j = \mathbf{a}_j - (\mathbf{q}_1^T \mathbf{a}_j) \mathbf{q}_1 - (\mathbf{q}_2^T \mathbf{a}_j) \mathbf{q}_2 - \cdots - (\mathbf{q}_{j-1}^T \mathbf{a}_j) \mathbf{q}_{j-1}$$

based on \mathbf{a}_j , the *modified* Gram-Schmidt procedure computes \mathbf{v}_j iteratively

$$\mathbf{v}_j^{(1)} = \mathbf{a}_j,$$

$$\mathbf{v}_j^{(2)} = \mathbf{v}_j^{(1)} - \mathbf{q}_1 \mathbf{q}_1^T \mathbf{v}_j^{(1)}, \quad \text{"subtract from } \mathbf{v}_j^{(1)} \text{ what is already in } \mathbf{q}_1 \text{"}$$

$$\mathbf{v}_j^{(3)} = \mathbf{v}_j^{(2)} - \mathbf{q}_2 \mathbf{q}_2^T \mathbf{v}_j^{(2)}, \quad \text{"subtract from } \mathbf{v}_j^{(2)} \text{ what is already in } \mathbf{q}_2 \text{"}$$

\vdots

$$\mathbf{v}_j = \mathbf{v}_j^{(j)} = \mathbf{v}_j^{(j-1)} - \mathbf{q}_{j-1} \mathbf{q}_{j-1}^T \mathbf{v}_j^{(j-1)}$$

Computing a *QR* factorization with the modified Gram-Schmidt procedure is stabler than with the classical Gram-Schmidt procedure. However, even the modified Gram-Schmidt procedure can lead to vectors $\mathbf{q}_1, \dots, \mathbf{q}_n$ that are far from orthogonal if the condition number of \mathbf{A} is large (see, Golub et al., *Matrix Computations*, Section 5.2.9)

Let's recall what the Gram-Schmidt procedure is doing: It is applying a succession of triangular matrices R_k on the right of A so that the resulting matrix

$$A \underbrace{R_1 R_2 \dots R_n}_{R^{-1}} = Q$$

has orthonormal columns and R is upper-triangular.

Let's recall what the Gram-Schmidt procedure is doing: It is applying a succession of triangular matrices R_k on the right of A so that the resulting matrix

$$\mathbf{A} \underbrace{\mathbf{R}_1 \mathbf{R}_2 \dots \mathbf{R}_n}_{\mathbf{R}^{-1}} = \mathbf{Q}$$

has orthonormal columns and \mathbf{R} is upper-triangular.

Instead, we could try to find orthonormal matrices ($\mathbf{X}^T \mathbf{X} = \mathbf{X} \mathbf{X}^T = \mathbf{I}$) so that

$$\underbrace{\mathbf{Q}_n \dots \mathbf{Q}_2 \mathbf{Q}_1}_{\mathbf{Q}} \mathbf{A} = \mathbf{R}$$

is upper-triangular. The product $\mathbf{Q}_n \dots \mathbf{Q}_2 \mathbf{Q}_1 = \mathbf{Q}^T$ is orthonormal too and thus $\mathbf{A} = \mathbf{Q} \mathbf{R}$ a QR factorization of \mathbf{A} .

The Householder method judiciously finds the matrices $\mathbf{Q}_1, \mathbf{Q}_2, \dots, \mathbf{Q}_n$ via so-called Householder reflectors \rightsquigarrow **board**. The Householder method is backward stable.

The QR factorization

All these three algorithms (classical Gram-Schmidt, modified Gram-Schmidt, Householder triangularization) have roughly the FLOPs of $2mn^2$ for an $m \times n$ matrix

Why would we ever want to use (modified) Gram-Schmidt instead of Householder triangularization?

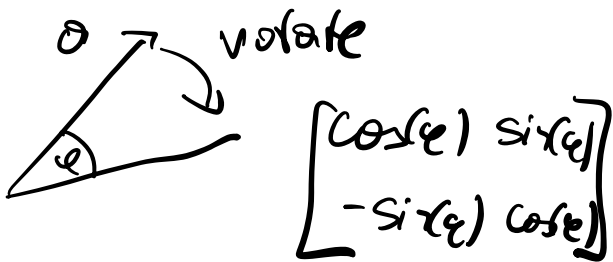
Transform A by multiplications with orthogonal matrices

$$A \sim Q_1 A \sim Q_2 Q_1 A \sim \dots$$

$$\|Q\|_2 = \|Q^{-1}\|_2 = 1 \quad | \quad \det(Q)^2 = 1$$

What are basic orthogonal transformations in \mathbb{R}^2 ?

rotations ($\det = 1$)

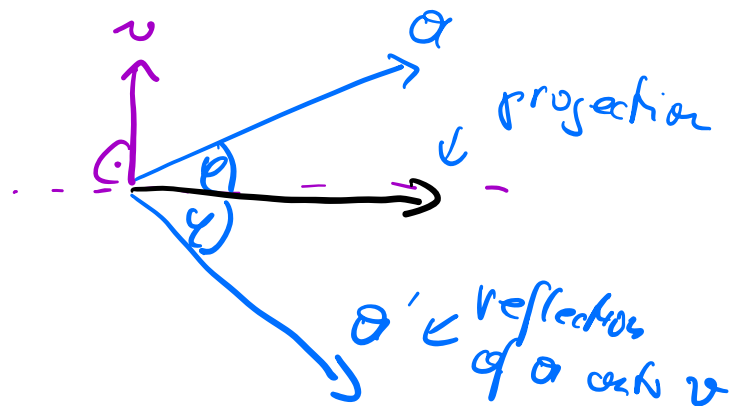


\Rightarrow Given rotations to construct Q with $A = QR$

\Rightarrow Jacobi method

reflections ($\det = -1$)

reflect at a hyperplane normal to v



projection

$$a \mapsto \left(I - \frac{vv^T}{v^T v} \right) a$$

reflection

$$a \mapsto \left(I - 2 \frac{vv^T}{v^T v} \right) a$$

\Rightarrow Householder reflections

Householder transformations

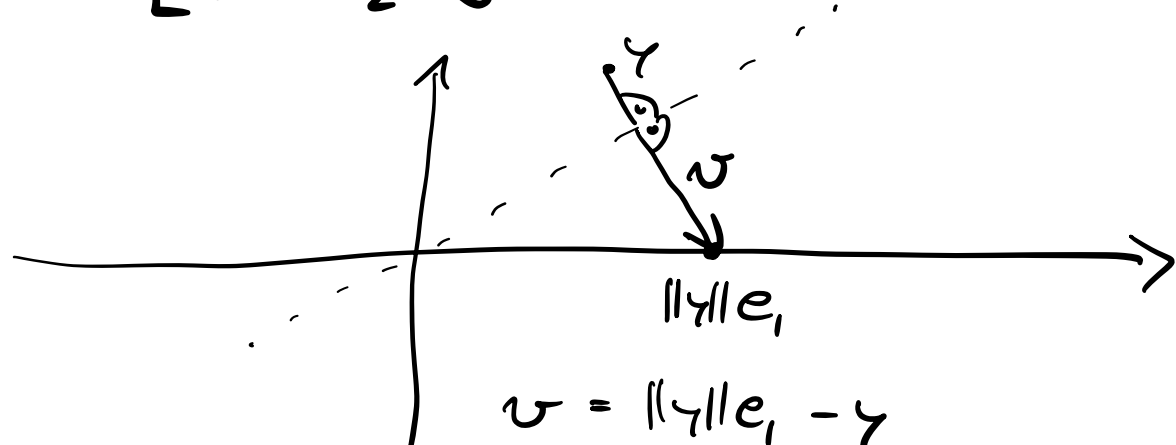
$$\begin{array}{c}
 \begin{bmatrix} x & x & x \\ x & x & x \\ x & x & x \\ x & x & x \\ x & x & x \end{bmatrix} \xrightarrow{Q_1} \begin{bmatrix} x & x & x \\ 0 & x & x \\ 0 & x & x \\ 0 & x & x \\ 0 & x & x \end{bmatrix} \xrightarrow{Q_2} \begin{bmatrix} x & x & x \\ x & x & x \\ 0 & x & x \\ 0 & x & x \\ 0 & x & x \end{bmatrix} \\
 A \qquad \qquad \qquad Q_1 A \qquad \qquad \qquad Q_2 Q_1 A
 \end{array}$$

of step $k=2$

$$y = \begin{bmatrix} x \\ x \\ x \\ x \end{bmatrix} \xrightarrow{F} \begin{bmatrix} * \\ 0 \\ 0 \\ 0 \end{bmatrix} = * e_1$$

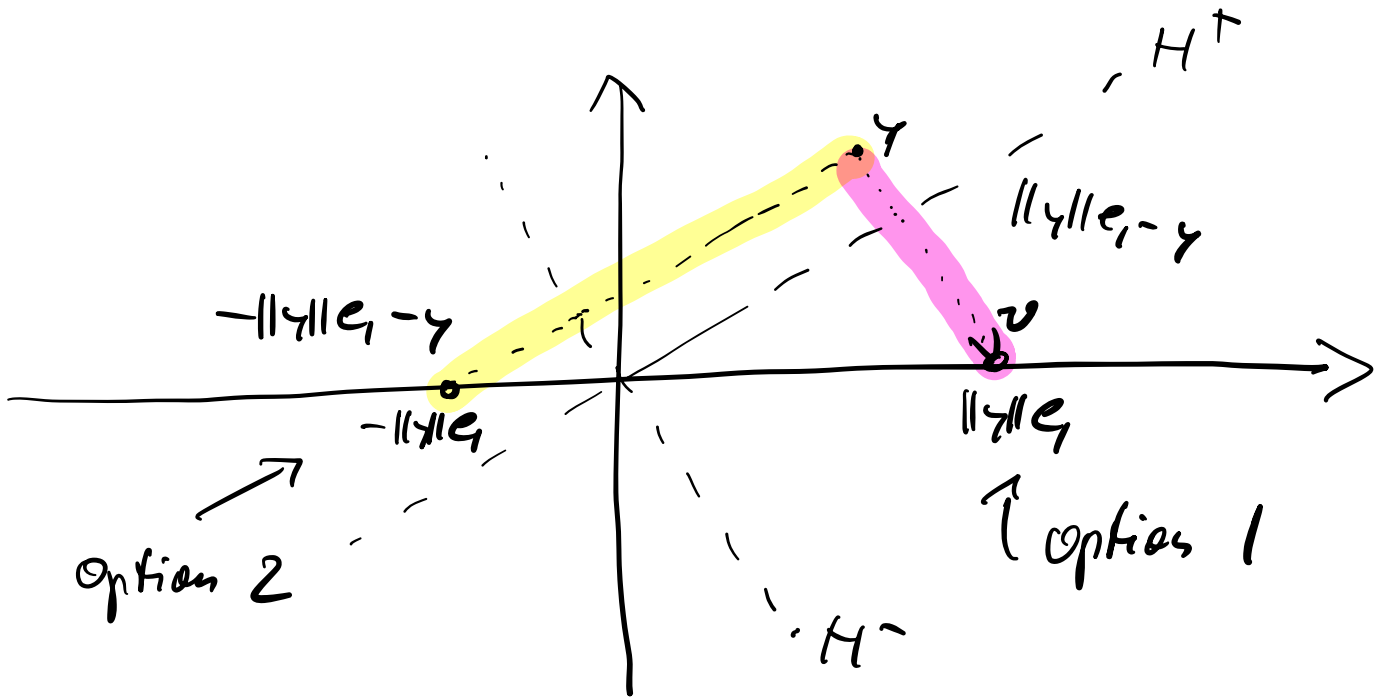
find

$$Q_2 = \begin{bmatrix} I_{k-1} & 0 \\ 0 & F_2 \end{bmatrix}$$



$$v = ||y||e_1 - y$$

$$F_2 = (I - 2 \frac{v v^T}{v^T v})$$



distance of reflection

→ taking the reflection corresponding to the longer distance is numerically more stable.

The QR factorization

All these three algorithms (classical Gram-Schmidt, modified Gram-Schmidt, Householder triangularization) have roughly the FLOPs of $2mn^2$ for an $m \times n$ matrix

Why would we ever want to use (modified) Gram-Schmidt instead of Householder triangularization? Gram-Schmidt can be easier to parallelize, for example (Recall that best algorithm depends also on what hardware we want to implement it on.)

Every matrix $A \in \mathbb{R}^{m \times n}$ with $m \geq n$ has a QR factorization. It is unique if we require the diagonal elements of R to be positive.

If $m > n$ and $\mathbf{Q} \in \mathbb{R}^{m \times n}$, then we speak of a reduced QR factorization. Otherwise, we have $\mathbf{Q} \in \mathbb{R}^{m \times m}$ and we speak of a full QR factorization.

```
1: >> A = randn(10, 10); [Q, R] = qr(A);
2: >> size(Q)
3: ans =
4:      10      10
5: >> size(R)
6: ans =
7:      10      10
```

```
1: >> A = randn(10, 4); [Q, R] = qr(A)
2: >> size(Q)
3: ans =
4:     10     10
5: >> size(R)
6: ans =
7:     10     4
8: >>
9: >> [Q, R] = qr(A, 0); % reduced QR
10: >> size(Q)
11: ans =
12:     10     4
13: >> size(R)
14: ans =
15:     4     4
```

Back to our least-squares problem

One would like to avoid the multiplication $A^T A$ and use a suitable factorization of A that avoids solving the normal equation directly:

$$A = QR = [Q_1, Q_2] \begin{bmatrix} R_1 \\ 0 \end{bmatrix} = Q_1 R_1,$$

where $Q \in \mathbb{R}^{m \times m}$ is an orthonormal matrix ($QQ^T = I$), and $R \in \mathbb{R}^{m \times n}$ consists of an upper triangular matrix and a block of zeros.

How can the QR factorization be used to solve the least-squares problem?

Back to our least-squares problem

One would like to avoid the multiplication $A^T A$ and use a suitable factorization of A that avoids solving the normal equation directly:

$$A = QR = [Q_1, Q_2] \begin{bmatrix} R_1 \\ 0 \end{bmatrix} = Q_1 R_1,$$

where $Q \in \mathbb{R}^{m \times m}$ is an orthonormal matrix ($QQ^T = I$), and $R \in \mathbb{R}^{m \times n}$ consists of an upper triangular matrix and a block of zeros.

How can the QR factorization be used to solve the least-squares problem?

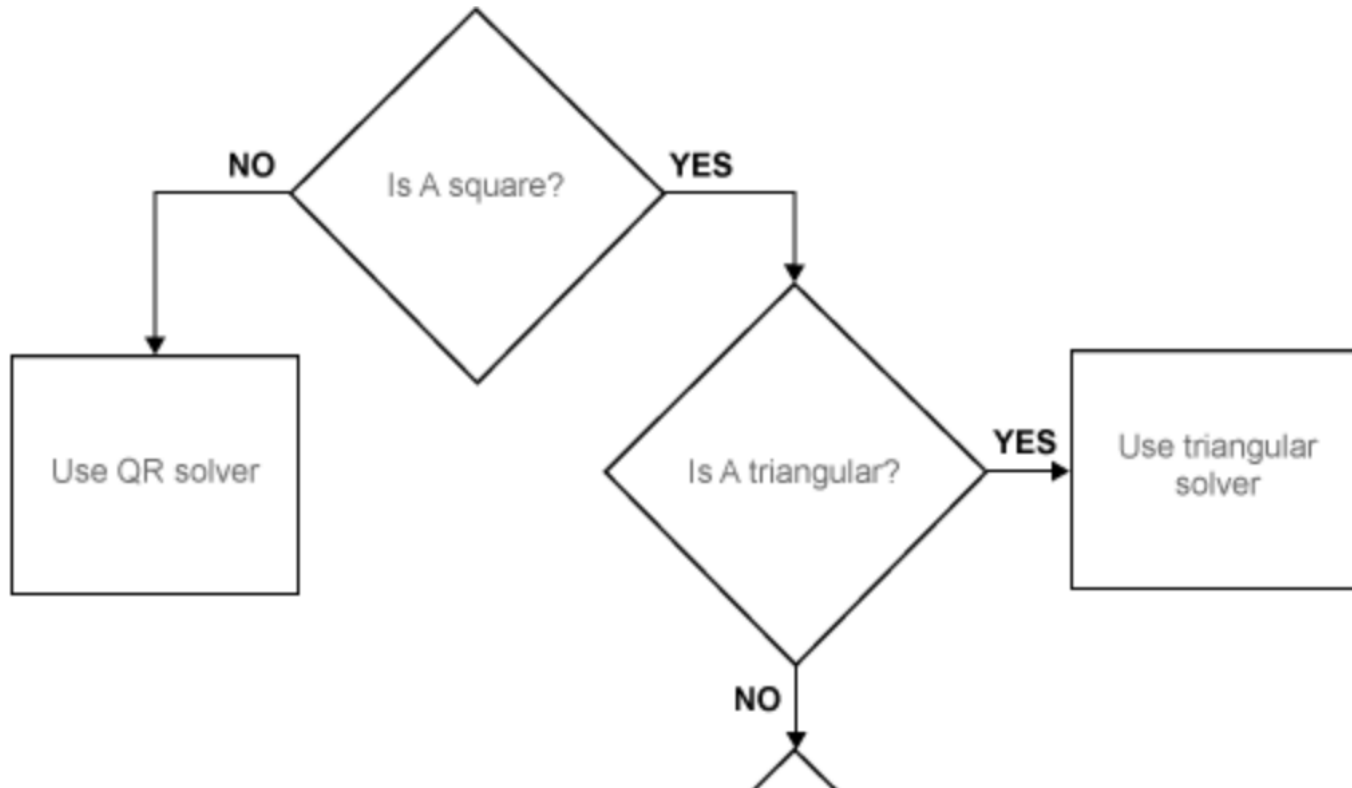
$$\begin{aligned} \min_{\mathbf{x}} \|A\mathbf{x} - \mathbf{b}\|^2 &= \min_{\mathbf{x}} \|Q^T(A\mathbf{x} - \mathbf{b})\|^2 &&= \min_{\mathbf{x}} \left\| \begin{bmatrix} \mathbf{b}_1 - R_1\mathbf{x} \\ \mathbf{b}_2 \end{bmatrix} \right\|^2, \\ &= \min_{\mathbf{x}} \|\mathbf{b}_1 - R_1\mathbf{x}\|^2 + \|\mathbf{b}_2\|^2 \end{aligned}$$

where $Q^T \mathbf{b} = \begin{bmatrix} \mathbf{b}_1 \\ \mathbf{b}_2 \end{bmatrix}$.

Thus, the least squares solution is $\mathbf{x} = R^{-1} \mathbf{b}_1$ and the residual is $\|\mathbf{b}_2\|$.

Stability of solving least-squares problem with Householder triangularization

Solving a least-squares problem with $\mathbf{A} \in \mathbb{R}^{m \times n}$, $m \geq n$ and $\text{rank}(\mathbf{A}) = n$ via QR factorization computed with Householder triangularization is backward stable.



Numerical Methods I

MATH-GA 2010.001/CSCI-GA 2420.001

Benjamin Peherstorfer
Courant Institute, NYU

Based on slides by G. Stadler and A. Donev

Eigen decomposition

Eigen decomposition

- ▶ For a square matrix $\mathbf{A} \in \mathbb{C}^{n \times n}$, there exists at least one λ such that

$$\mathbf{Ax} = \lambda\mathbf{x} \implies (\mathbf{A} - \lambda\mathbf{I})\mathbf{x} = 0$$

- ▶ Putting the eigenvectors \mathbf{x}_j as columns in a matrix \mathbf{X} , and the eigenvalues λ_j on the diagonal of a diagonal matrix $\mathbf{\Lambda}$, we get

$$\mathbf{AX} = \mathbf{X}\mathbf{\Lambda}$$

- ▶ A matrix is non-defective or diagonalizable if there exist n linearly independent eigenvectors, which means that \mathbf{X} is invertible

$$\mathbf{X}^{-1}\mathbf{AX} = \mathbf{\Lambda}$$

$$\mathbf{A} = \mathbf{X}\mathbf{\Lambda}\mathbf{X}^{-1}$$

- ▶ The transformation from \mathbf{A} to $\mathbf{\Lambda} = \mathbf{X}^{-1}\mathbf{AX}$ is called a similarity transformation and it preserves the eigenvalues.

- ▶ A matrix is unitarily diagonalizable if there exist n linearly independent orthogonal eigenvectors, i.e., if the matrix \mathbf{X} can be chosen to be unitary (orthonormal), $\mathbf{X} = \mathbf{U}$, where $\mathbf{U}^{-1} = \mathbf{U}^H$

$$\mathbf{A} = \mathbf{U}\mathbf{\Lambda}\mathbf{U}^H$$

Note that unitary matrices generalize orthogonal matrices to the complex domain, so we use adjoints (conjugate transpose) instead of transpose throughout

- ▶ Theorem: A matrix is unitarily diagonalizable iff it is normal, i.e., it commutes with its adjoint:

$$\mathbf{A}^H \mathbf{A} = \mathbf{A} \mathbf{A}^H$$

- ▶ Theorem: Hermitian (symmetric) matrices, $\mathbf{A}^H = \mathbf{A}$, are unitarily diagonalizable and have *real* eigenvalues.

- ▶ The usual eigenvectors are more precisely called *right* eigenvectors. There are also *left* eigenvectors corresponding to a given eigenvalue λ

$$\mathbf{y}^H \mathbf{A} = \lambda \mathbf{y}^H \implies \mathbf{A}^H \mathbf{y} = \bar{\lambda} \mathbf{y},$$

$$\mathbf{Y}^H \mathbf{A} = \mathbf{\Lambda} \mathbf{Y}^H$$

with conjugate $\bar{\lambda}$ of λ

- ▶ For a matrix that is diagonalizable, observe that

$$\mathbf{Y}^H = \mathbf{X}^{-1}$$

and so the left eigenvectors provide no new information

- ▶ For unitarily diagonalizable matrices, $\mathbf{Y} = (\mathbf{X}^{-1})^H = (\mathbf{X}^H)^H = \mathbf{X} = \mathbf{U}$, so that the left and right eigenvectors coincide.

Numerically finding eigenvalues

For a matrix $A \in \mathbb{C}^{n \times n}$ (potentially real), we want to find $\lambda \in \mathbb{C}$ and $\mathbf{x} \neq 0$ such that

$$A\mathbf{x} = \lambda\mathbf{x}.$$

Most relevant problems:

- ▶ A symmetric (and large)
- ▶ A spd (and large)
- ▶ A stochastic matrix, i.e., all entries $0 \leq a_{ij} \leq 1$ are probabilities, and thus $\sum_j a_{ij} = 1$.

How hard are they to find numerically?

How hard are they to find numerically?

- ▶ This is a **nonlinear** problem.

How hard are they to find numerically?

- ▶ This is a **nonlinear** problem.
- ▶ How **difficult** is this?

How hard are they to find numerically?

- ▶ This is a **nonlinear** problem.
- ▶ How **difficult** is this? Eigenvalues are the roots of the characteristic polynomial. Also, any polynomial is the characteristic polynomial of a matrix \rightsquigarrow For matrices larger than 4×4 , eigenvalues cannot be computed in closed form (Abel's theorem).

How hard are they to find numerically?

- ▶ This is a **nonlinear** problem.
- ▶ How **difficult** is this? Eigenvalues are the roots of the characteristic polynomial. Also, any polynomial is the characteristic polynomial of a matrix \rightsquigarrow For matrices larger than 4×4 , eigenvalues cannot be computed in closed form (Abel's theorem).
- ▶ Must use an **iterative** algorithm

How hard are they to find numerically?

- ▶ This is a **nonlinear** problem.
- ▶ How **difficult** is this? Eigenvalues are the roots of the characteristic polynomial. Also, any polynomial is the characteristic polynomial of a matrix \rightsquigarrow For matrices larger than 4×4 , eigenvalues cannot be computed in closed form (Abel's theorem).
- ▶ Must use an **iterative** algorithm \rightsquigarrow this is fundamentally different from what we have seen previously when solving systems of *linear* equations! These algorithms (LU, QR) give the *exact* solution in *exact* arithmetic in finite number of steps. We *cannot* expect something similar for computing eigenvalues!

Condition of finding eigenvalues of a matrix

The absolute condition number of determining a simple eigenvalue λ_0 of a matrix $\mathbf{A} \in \mathbb{C}^{n \times n}$ with respect to the $\|\cdot\|_2$ is

$$\kappa_{\text{abs}} = \frac{1}{|\cos(\angle(\mathbf{x}, \mathbf{y}))|}, \quad \cos(\angle(\mathbf{x}, \mathbf{y})) = \frac{|\langle \mathbf{x}, \mathbf{y} \rangle|}{\|\mathbf{x}\| \|\mathbf{y}\|}$$

and the relative condition number is

$$\kappa_{\text{rel}} = \frac{\|\mathbf{A}\|}{|\lambda_0 \cos(\angle(\mathbf{x}, \mathbf{y}))|},$$

where \mathbf{x} is an eigenvector of \mathbf{A} for the eigenvalue λ_0 ($\mathbf{A}\mathbf{x} = \lambda_0\mathbf{x}$) and \mathbf{y} an adjoint eigenvector ($\mathbf{A}^H\mathbf{y} = \bar{\lambda}_0\mathbf{y}$).

Sketch of proof \rightsquigarrow board \rightsquigarrow next time

(see also Deufhard, Theorem 5.2)

Interpretation

Perturbations of order δ in entries of matrix \mathbf{A} induce changes of the order $\delta\lambda = \delta / \cos(\angle(\mathbf{x}_0, \mathbf{y}_0))$

In particular, for normal matrices* ($\mathbf{A}\mathbf{A}^H = \mathbf{A}^H\mathbf{A}$), we have $\mathbf{x}_0 = \mathbf{y}_0$ and thus $\angle(\mathbf{x}_0, \mathbf{y}_0) = 0$ and thus $\cos(\angle(\mathbf{x}_0, \mathbf{y}_0)) = 1$, which means $\kappa_{\text{abs}} = 1$, which can be considered well conditioned

Finding non-simple eigenvalues can have very high absolute condition number (but can still be done numerically). For a detailed treatment have a look at textbook by Golub et al. on Matrix Computations.

*Equivalent: Have orthonormal eigenbasis of \mathbb{C} ; diagonalizable by unitary matrix.

Bounding error in eigenvalue computation

Let $\mathbf{A} \in \mathbb{C}^{n \times n}$ be a Hermitian matrix and let $(\hat{\lambda}, \hat{\mathbf{x}})$ be a computed approximation of an eigenvalue/eigenvector pair (λ, \mathbf{x}) of \mathbf{A} . Defining the residual

$$\hat{\mathbf{r}} = \mathbf{A}\hat{\mathbf{x}} - \hat{\lambda}\hat{\mathbf{x}}, \quad \hat{\mathbf{x}} \neq \mathbf{0},$$

it then follows that

$$\min_{\lambda_i \in \sigma(\mathbf{A})} |\hat{\lambda} - \lambda_i| \leq \frac{\|\hat{\mathbf{r}}\|_2}{\|\hat{\mathbf{x}}\|_2},$$

where $\sigma(\mathbf{A}) = \{\lambda | \lambda \text{ is an eigenvalue of } \mathbf{A}\}$ is the spectrum of \mathbf{A} .

Proof \rightsquigarrow board

What is special about this bound?

Let $A \in \mathbb{C}^{n \times n}$ Hermitian, let $(\hat{\lambda}, \hat{x})$ approximation of eigen pair (λ, x) . Define residual

$$\hat{r} = A\hat{x} - \hat{\lambda}\hat{x}, \quad \hat{x} \neq 0$$

it follows

$$\min_{\lambda_i \in \sigma(A)} |\hat{\lambda} - \lambda_i| \leq \frac{\|\hat{r}\|}{\|\hat{x}\|}$$

A is Hermitian, exists an orthogonal eigenbasis of \mathbb{C}^n

$$\hat{x} = \sum_{i=1}^n \alpha_i u_i \quad \text{with} \quad \alpha_i = u_i^H \hat{x}$$

$$\hat{r} = A\hat{x} - \hat{\lambda}\hat{x} = A\left(\sum_{i=1}^n \alpha_i u_i\right) - \hat{\lambda} \sum_{i=1}^n \alpha_i u_i$$

$$= \sum_{i=1}^n \alpha_i \lambda_i u_i - \hat{\lambda} \sum_{i=1}^n \alpha_i u_i$$

$$= \sum_{i=1}^n \alpha_i (\lambda_i - \hat{\lambda}) u_i$$

$$\frac{\|\hat{r}\|^2}{\|\hat{x}\|^2} = \frac{\sum_{i=1}^n |\alpha_i|^2 (\lambda_i - \hat{\lambda})^2 \|u_i\|^2}{\sum_{j=1}^n |\alpha_j|^2 \|u_j\|^2}$$

$$= \sum_{i=1}^n \left(\frac{|\alpha_i|^2}{\sum_{j=1}^n |\alpha_j|^2} \right) (\lambda_i - \hat{\lambda})^2 = \sum_{i=1}^n \beta_i (\lambda_i - \hat{\lambda})^2$$

Now note that $\beta_i \geq 0$

$$\sum_{i=1}^n \beta_i = \sum_{i=1}^n \frac{|\alpha_i|^2}{\sum_{j=1}^n |\alpha_j|^2} = 1$$

Thus

$$\begin{aligned} \frac{\|\hat{r}\|^2}{\|\hat{x}\|^2} &= \sum_{i=1}^n \beta_i (\lambda_i - \hat{\lambda})^2 \geq \sum_{i=1}^n \beta_i \min_j (\lambda_j - \hat{\lambda})^2 \\ &= \min_j (\lambda_j - \hat{\lambda})^2 \underbrace{\sum_{i=1}^n \beta_i}_{=1} = \min_j (\lambda_j - \hat{\lambda})^2 \end{aligned}$$

$$\min_j |\lambda_j - \hat{\lambda}| \leq \frac{\|\hat{r}\|}{\|\hat{x}\|}$$

Bounding error in eigenvalue computation

Let $\mathbf{A} \in \mathbb{C}^{n \times n}$ be a Hermitian matrix and let $(\hat{\lambda}, \hat{\mathbf{x}})$ be a computed approximation of an eigenvalue/eigenvector pair (λ, \mathbf{x}) of \mathbf{A} . Defining the residual

$$\hat{\mathbf{r}} = \mathbf{A}\hat{\mathbf{x}} - \hat{\lambda}\hat{\mathbf{x}}, \quad \hat{\mathbf{x}} \neq \mathbf{0},$$

it then follows that

$$\min_{\lambda_i \in \sigma(\mathbf{A})} |\hat{\lambda} - \lambda_i| \leq \frac{\|\hat{\mathbf{r}}\|_2}{\|\hat{\mathbf{x}}\|_2},$$

where $\sigma(\mathbf{A}) = \{\lambda | \lambda \text{ is an eigenvalue of } \mathbf{A}\}$ is the spectrum of \mathbf{A} .

Proof \rightsquigarrow board

What is special about this bound?

- ▶ This is an *a posteriori* bound that bounds the error *after* we have computed the result
- ▶ We will see many more residual-based *a posteriori* bounds (broadly speaking: the residual is something we can compute, and if the problem is “well-behaved” then the norm of the residual is a reasonable bound of the norm of the error.)

Condition of computing eigenvectors

- ▶ The condition of computing eigenvector \mathbf{x}_i for an eigenvalue λ_i depends on the separation between the eigenvalues

$$\kappa = \frac{1}{\min_{i \neq j} |\lambda_i - \lambda_j|}$$

(Quarteroni et al., Section 5)

- ▶ Computing \mathbf{x}_i can be ill-conditioned if some eigenvalue λ_j is “very close” to the eigenvalue λ_i
- ▶ This indicates that multiple eigenvalues require care. Even for Hermitian matrices *eigenvectors* can be hard to compute

The Power Method

Let $\mathbf{A} \in \mathbb{C}^{n \times n}$ be diagonalizable matrix and λ_1 be a simple eigenvalue with

$$|\lambda_1| > |\lambda_2| \geq \dots \geq |\lambda_n|$$

Let \mathbf{x}_0 be an initial guess that is not orthogonal to the eigenspace of λ_1 , then \mathbf{x}_k obtained via the iterations

$$\mathbf{z}_{k+1} = \mathbf{A}\mathbf{x}_k \tag{1}$$

$$\mathbf{x}_{k+1} = \mathbf{z}_{k+1} / \|\mathbf{z}_{k+1}\|_2 \tag{2}$$

will converge to the normalized eigenvector of \mathbf{A} corresponding to λ_1 for $k \rightarrow \infty$.

This process is called the power method.

Proof \rightsquigarrow board

Power method

↳ first
$$x_k = \frac{A^k x_0}{\|A^k x_0\|_2}$$

↳ Because A is diagonalizable, have vectors

$$\{v_i\}_{i=1}^n$$

with

$$x_0 = \sum_{i=1}^n \alpha_i v_i, \quad \alpha_1 \neq 0$$

↳ Further

$$\begin{aligned} A^k x_0 &= V \Lambda^k V^{-1} x_0 = V \Lambda^k \begin{bmatrix} \alpha_1 \\ \vdots \\ \alpha_n \end{bmatrix} \\ &= \sum_{i=1}^n \alpha_i \lambda_i^k v_i \\ &= \alpha_1 \lambda_1^k \left(v_1 + \underbrace{\sum_{i=2}^n \frac{\alpha_i}{\alpha_1} \left(\frac{\lambda_i}{\lambda_1}\right)^k v_i}_{\gamma(k)} \right) \end{aligned}$$

Because $\left| \frac{\lambda_i}{\lambda_1} \right| < 1$ per definition of λ_1 ,
the vector $A^k x_0$ assumes an increasingly
significant component in direction v_1

$$x_k = \frac{\alpha_1 \lambda_1^k (v_1 + \gamma^{(k)})}{\|\alpha_1 \lambda_1^k (v_1 + \gamma^{(k)})\|_2} = \gamma_k \frac{v_1 + \gamma^{(k)}}{\|v_1 + \gamma^{(k)}\|_2}$$

where γ_k is the sign of $\alpha_1 \lambda_1^k$.

With $\|\gamma^{(k)}\| \Rightarrow 0$ for $k \Rightarrow \infty$

$$x_k \rightarrow \frac{v_1}{\|v_1\|} \curvearrowright$$

Convergence speed

Component in direction v_2 converges slowest with $O\left(\left|\frac{\lambda_2}{\lambda_1}\right|^k\right)$

\rightarrow the bigger the gap $\lambda_1 \leftrightarrow \lambda_2$, the faster

\rightarrow again, well separated EV are 'easier' to find

Compute eigenvalue $\lambda_1^{(k)}$ via

$$\frac{\gamma^T A \gamma}{\gamma^T \gamma} = \frac{\gamma^T \lambda \gamma}{\gamma^T \gamma} = \lambda \quad \text{if } \gamma \text{ is EV with EV } \lambda \text{ of } A$$

$$\lambda_1^{(k)} = \frac{x_k^T A x_k}{x_k^T x_k}$$

Power method (cont'd)

Start with initial guess \mathbf{x}_0 and then iterate

- ▶ Compute matrix-vector product and normalize it

$$\mathbf{x}_k = \frac{\mathbf{A}\mathbf{x}_{k-1}}{\|\mathbf{A}\mathbf{x}_{k-1}\|}$$

- ▶ Obtain eigenvalue estimate (note that $\|\mathbf{x}_k\| = 1$)

$$\lambda_1^{(k)} = \mathbf{x}_k^H \mathbf{A} \mathbf{x}_k$$

- ▶ Test for convergence? **How?**

Power method (cont'd)

Start with initial guess \mathbf{x}_0 and then iterate

- ▶ Compute matrix-vector product and normalize it

$$\mathbf{x}_k = \frac{\mathbf{A}\mathbf{x}_{k-1}}{\|\mathbf{A}\mathbf{x}_{k-1}\|}$$

- ▶ Obtain eigenvalue estimate (note that $\|\mathbf{x}_k\| = 1$)

$$\lambda_1^{(k)} = \mathbf{x}_k^H \mathbf{A} \mathbf{x}_k$$

- ▶ Test for convergence? **How?** Compute residual

$$\mathbf{r}_k = \mathbf{A}\mathbf{x}_k - \lambda_1^{(k)} \mathbf{x}_k$$

and terminate if the error estimate is small enough (bound if \mathbf{A} Hermitian, heuristic otherwise)

$$\min_i |\lambda_i - \lambda_1^{(k)}| \leq \frac{\|\mathbf{r}_k\|}{\|\mathbf{x}_k\|} < \epsilon$$

Let us recall the condition number

$$\|f(\tilde{x}) - f(x)\| \leq \kappa_{\text{abs}} \|\tilde{x} - x\|, \quad x \rightarrow \tilde{x}$$

We define " $\dot{\leq}$ " to mean

$$g(x) \dot{\leq} h(x), \quad x \rightarrow x_0 \quad \hat{=} \quad g(x) \leq h(x) + o(\|h(x)\|) \quad \hat{=} \quad x \rightarrow x_0$$

Sometime you see

$$g(x) \leq h(x) + o(\|h(x)\|^2)$$

Thus, the condition number tells us:

$$\|f(x + \delta x) - f(x)\| \leq \kappa_{\text{abs}} \|\delta x\| + o(\|\delta x\|), \quad \delta x \rightarrow 0$$

This is a perturbation bound.

With condition number, here is a receipt how to derive perturbation bounds

- ① formulate problem via function f
- ② establish f is differentiable
- ③ take derivative of f and compute the norm

Let us start with " $Ax = b$ "

Bound for perturbation in b

$$f: \mathbb{R}^n \rightarrow \mathbb{R}^n, \quad b \mapsto f(b) = A^{-1}b \quad \text{①}$$

Function f is linear \rightarrow differentiable (2)

Derivative $f'(b) = A^{-1}$ with norm $\|A^{-1}\|$ (3)

Obtain

$$\|x - \bar{x}\| \leq \|A^{-1}\| \|\delta b\| + o(\|\delta b\|), \|\delta b\| \rightarrow 0$$

Now bound for " $Ax = b$ " with perturbations in A

$$f: GL(n) \rightarrow \mathbb{R}^n, A \mapsto f(A) = A^{-1}b \quad (1)$$

Differentiable

\hookrightarrow Cramer's rule

\hookrightarrow Neumann series

} Derif(hord

(3)

Now show: The mapping $g: GL(n) \rightarrow GL(n)$ with $g(A) = A^{-1}$ is differentiable and

$$g'(A)C = -A^{-1}CA^{-1} \quad C \in \mathbb{R}^{n \times n}$$

Take derivative of $(A+tC)(A+tC)^{-1} = I$

$$0 = C(A+tC)^{-1} + (A+tC) \frac{d}{dt}(A+tC)^{-1}$$

$$\frac{d}{dt}(A+tC)^{-1} = \underline{-(A+tC)^{-1}C(A+tC)^{-1}}$$

$$\underline{g'(A)C} = \left. \frac{d}{dt}(A+tC)^{-1} \right|_{t=0} = -A^{-1}CA^{-1}$$

Now back to our condition number

$$f(A) = A^{-1}b$$

$$f'(A)C = -A^{-1}CA^{-1}b = -A^{-1}Cx$$

$$\kappa_{\text{rel}} = \frac{\|A\|}{\|f(A)\|} \|f'(A)\|_{\text{op}}$$

$$= \frac{\|A\|}{\|x\|} \sup_{\|C\|=1} \|-A^{-1}Cx\| \leq$$

$$\leq \frac{\|A\|}{\|x\|} \sup_{\|C\|=1} \|A^{-1}\| \|C\| \|x\|$$

$$= \|A\| \|A^{-1}\| = \kappa(A)$$

Perturbation bound

$$\frac{\|x - \tilde{x}\|}{\|x\|} \leq \|A\| \|A^{-1}\| \frac{\|A - \tilde{A}\|}{\|A\|} + o(\|\delta A\|)$$

$\|\delta A\| \rightarrow 0$

Now least squares

Let $A \in \mathbb{R}^{m \times n}$, $m \geq n$, full rank(A) = n,

and x the unique solution

$$\|b - Ax\|_2 = \min$$

Assume $x \neq 0$ and define angle ν between b and $R(A)$ of A as

$$\sin(\nu) = \frac{\|b - Ax\|_2}{\|b\|_2} = \frac{\|r\|_2}{\|b\|_2}$$

Then

(a) rel condition number w.r.t. perturbations in b is

$$\kappa_{\text{rel}} \leq \frac{\kappa_2(A)}{\cos(\nu)}$$

(b) rel cond. number w.r.t. A

$$\kappa_{\text{rel}} \leq \kappa_2(A) + \kappa_2(A)^2 \tan(\nu)$$

Let's look at (a) first:

Solution

$$x = \phi(b) = (A^T A)^{-1} A^T b \quad (1)$$

ϕ is differentiable because is linear in b

$$\phi'(b) = (A^T A)^{-1} A^T \quad (2)$$

Assemble condition number

$$\kappa_{\text{rel}} = \frac{\|b\|}{\|x\|} \|\phi'(b)\|$$

$$= \frac{\|b\|}{\|x\|} \|(A^T A)^{-1} A^T\|_2 \frac{\|A\|_2}{\|A\|_2}$$

• For or full column rank matrix A

$$\kappa_2(A) = \|A\|_2 \|(A^T A)^{-1} A^T\|_2$$

• Recall where P is proj onto $R(A)$

$$\frac{\|b\|}{\|Pb\|} = \frac{1}{\cos(\nu)} \quad P = A(A^T A)^{-1} A^T$$

$$\begin{aligned}
 \rho_{\text{rel}} &= \frac{\|b\|}{\|A\| \|x\|} \kappa_2(A) \\
 &\leq \frac{\|b\|}{\|Ax\|} \kappa_2(A) = \frac{\|b\|}{\|A(A^T A)^{-1} A^T b\|} \kappa_2(A) \\
 &= \frac{\|b\|}{\|Pb\|} \kappa_2(A) = \frac{\kappa_2(A)}{\cos(\nu)}
 \end{aligned}$$

Perturbation bound

$$\frac{\|x - \tilde{x}\|}{\|x\|} \leq \frac{\kappa_2(A)}{\cos(\nu)} \frac{\|Sb\|}{\|b\|} + o(\|Sb\|), \quad \|Sb\| \rightarrow 0$$

Now (b): Our function

$$\phi(A) = (A^T A)^{-1} A^T b \quad \textcircled{1}$$

ϕ is differentiable, convinced us already earlier $\textcircled{2}$

Construct derivative by using equation that defines ϕ

$$(A+tC)^T (A+tC) \underbrace{\phi(A+tC)}_{=x} = (A+tC)^T b$$

Total derivative

$$\begin{aligned}
 &C^T (A+tC) \phi(A+tC) + \\
 &(A+tC)^T C \phi(A+tC) + \\
 &(A+tC)^T (A+tC) \frac{d}{dt} \phi(A+tC) = C^T b
 \end{aligned}$$

Now for $t=0$ obtain $\lim_{t \rightarrow 0} \phi'(A)C$

$$C^T A \underbrace{\phi(A)}_{=x} + A^T C \underbrace{\phi(A)}_{=x} + A^T A \phi'(A)C = C^T b$$

$$C^T A x + A^T C x + A^T A \phi(A) C = C^T b$$

$$\begin{aligned} \phi'(A) C &= (A^T A)^{-1} (C^T (b - A x) - A^T C x) \\ &= (A^T A)^{-1} C^T (b - A x) - (A^T A)^{-1} A^T C x \end{aligned}$$

Estimate norm

$$\begin{aligned} \|\phi'(A)\| &= \sup_{\|c\|=1} \|\phi'(A) c\| \\ &\leq \|(A^T A)^{-1}\| \|b - A x\| + \|(A^T A)^{-1} A^T\| \|x\| \end{aligned}$$

Condition number

$$\begin{aligned} \kappa_{\text{rel}} &= \frac{\|A\|}{\|\phi(A)\|} \|\phi'(A)\| \\ &\leq \frac{\|A\|}{\|x\|} \|(A^T A)^{-1} A^T\| \|x\| + \frac{\|A\|}{\|x\|} \|(A^T A)^{-1}\| \|v\| \frac{\|A\|}{\|A\|} \\ &= \underbrace{\|A\| \|(A^T A)^{-1} A^T\|}_{\kappa_2(A)} + \underbrace{\|A\|^2 \|(A^T A)^{-1}\|}_{\kappa_2(A^T A) = \kappa_2(A)^2} \left(\frac{\|v\|}{\|A\| \|x\|} \right) \end{aligned}$$

$$\frac{\|v\|}{\|A\| \|x\|} = \frac{\|v\|}{\|b\|} \frac{\|b\|}{\|A\| \|x\|} \leq \sin(v) \frac{1}{\cos(v)} = \tan(v)$$

Again this gives perturbation bound.

Now eigenvalues

Let $\lambda_0 \in \mathbb{C}$ be a simple eigenvalue of $A \in \mathbb{C}^{n \times n}$.

Then there is a cont. diff. mapping

$$\lambda: V \subset \mathbb{C}^{n \times n} \rightarrow \mathbb{C}, \quad B \mapsto \lambda(B)$$

from a neighborhood V of A in $\mathbb{C}^{n \times n}$ such that

$$\lambda(A) = \lambda_0$$

and $\lambda(B)$ is a simple eigenvalue of B for all $B \in V$.

If x_0 is an eigenvector of A for λ_0 and y_0 is a (left) eigenvector of $A^* = A^H$ for the eigenvalue

$$\overline{\lambda_0}$$

so that

$$Ax_0 = \lambda_0 x_0 \quad \text{and} \quad A^* y_0 = \overline{\lambda_0} y_0,$$

then the derivative of λ at A satisfies

$$\lambda'(A)C = \frac{\langle Cx_0, y_0 \rangle}{\langle x_0, x_0 \rangle} \quad \text{for all } C \in \mathbb{C}^{n \times n}.$$

Similar to before, we consider the matrix

$$A + tC, \quad C \in \mathbb{C}^{n \times n}$$

and want to understand the solution map $\lambda(t)$.

✓ $A \in \mathbb{C}^{n \times n}$ has n eigenvalues of which $d \leq n$ are distinct, then the char. polynomial is

$$\begin{aligned}\chi_A(\lambda) &= \det(A - \lambda I) = \prod_{i=1}^n (\lambda_i - \lambda) \\ &= \prod_{i=1}^d (\lambda_i' - \lambda)^{\mu(\lambda_i')}\end{aligned}$$

with $\mu(\lambda_i')$ being the algebraic multiplicity of the eigenvalue λ_i' .

Simple means

$$\mu(\lambda_0) = 1$$

Because we have a simple eigenvalue

$$\chi_A'(\lambda_0) \neq 0$$

Define function

$$F(t, \lambda) = \chi_{A+tC}(\lambda)$$

Apply implicit function theorem to show

$\lambda: (-\epsilon, \epsilon) \rightarrow \mathbb{C}$ is cont. diff. with $\lambda(0) = \lambda_0$

$\lambda(t)$ simple EV of $A+tC$

↳ Consider (t_0, λ_0) , $t_0 = 0$

$$\text{↳ } F(t_0, \lambda_0) = \chi_A(\lambda_0) = 0$$

$$\hookrightarrow \frac{\partial}{\partial \lambda} F(t, \lambda) \Big|_{t_0, \lambda_0} = \chi_{A+tC}'(\lambda) \Big|_{t_0, \lambda_0} = \chi_A'(\lambda) \neq 0$$

\Rightarrow there exists $\lambda: (-\epsilon, \epsilon) \rightarrow \mathbb{C}$ as stated above.

Similarly, use that λ_0 is simple to determine cont. diff

$$x: (-\epsilon, \epsilon) \rightarrow \mathbb{C}^n, \quad t \mapsto x(t)$$

such that $x(0) = x_0$ and $x(t)$ is eigenvector of $A+tC$ for eigenvalue $\lambda(t)$.

Now we can differentiate w.r.t. t the equation

$$(A+tC)x(t) = \lambda(t)x(t)$$

Numerical Methods I

MATH-GA 2010.001/CSCI-GA 2420.001

Benjamin Peherstorfer
Courant Institute, NYU

Based on slides by G. Stadler and A. Donev

Today

Last time

- ▶ Computing eigenvalues
- ▶ Perturbation bounds

Today

- ▶ More on computing the eigenvalues

Announcements

- ▶ Homework 3 posted, is due Mon, Oct 21 before class

Numerically finding eigenvalues

For a matrix $A \in \mathbb{C}^{n \times n}$ (potentially real), we want to find $\lambda \in \mathbb{C}$ and $\mathbf{x} \neq 0$ such that

$$A\mathbf{x} = \lambda\mathbf{x}.$$

Most relevant problems:

- ▶ A symmetric (and large)
- ▶ A spd (and large)
- ▶ A stochastic matrix, i.e., all entries $0 \leq a_{ij} \leq 1$ are probabilities, and thus $\sum_j a_{ij} = 1$.

How hard are they to find numerically?

How hard are they to find numerically?

- ▶ This is a **nonlinear** problem.

How hard are they to find numerically?

- ▶ This is a **nonlinear** problem.
- ▶ How **difficult** is this?

How hard are they to find numerically?

- ▶ This is a **nonlinear** problem.
- ▶ How **difficult** is this? Eigenvalues are the roots of the characteristic polynomial. Also, any polynomial is the characteristic polynomial of a matrix \rightsquigarrow For matrices larger than 4×4 , eigenvalues cannot be computed in closed form (Abel's theorem).

How hard are they to find numerically?

- ▶ This is a **nonlinear** problem.
- ▶ How **difficult** is this? Eigenvalues are the roots of the characteristic polynomial. Also, any polynomial is the characteristic polynomial of a matrix \rightsquigarrow For matrices larger than 4×4 , eigenvalues cannot be computed in closed form (Abel's theorem).
- ▶ Must use an **iterative** algorithm

How hard are they to find numerically?

- ▶ This is a **nonlinear** problem.
- ▶ How **difficult** is this? Eigenvalues are the roots of the characteristic polynomial. Also, any polynomial is the characteristic polynomial of a matrix \rightsquigarrow For matrices larger than 4×4 , eigenvalues cannot be computed in closed form (Abel's theorem).
- ▶ Must use an **iterative** algorithm \rightsquigarrow this is fundamentally different from what we have seen previously when solving systems of *linear* equations! These algorithms (LU, QR) give the *exact* solution in *exact* arithmetic in finite number of steps. We *cannot* expect something similar for computing eigenvalues!

Condition of finding eigenvalues of a matrix

The absolute condition number of determining a simple eigenvalue λ_0 of a matrix $\mathbf{A} \in \mathbb{C}^{n \times n}$ with respect to the $\|\cdot\|_2$ is

$$\kappa_{\text{abs}} = \frac{1}{|\cos(\angle(\mathbf{x}, \mathbf{y}))|}, \quad \cos(\angle(\mathbf{x}, \mathbf{y})) = \frac{|\langle \mathbf{x}, \mathbf{y} \rangle|}{\|\mathbf{x}\| \|\mathbf{y}\|}$$

and the relative condition number is

$$\kappa_{\text{rel}} = \frac{\|\mathbf{A}\|}{|\lambda_0 \cos(\angle(\mathbf{x}, \mathbf{y}))|},$$

where \mathbf{x} is an eigenvector of \mathbf{A} for the eigenvalue λ_0 ($\mathbf{A}\mathbf{x} = \lambda_0\mathbf{x}$) and \mathbf{y} an adjoint eigenvector ($\mathbf{A}^H\mathbf{y} = \bar{\lambda}_0\mathbf{y}$).

Sketch of proof \rightsquigarrow board \rightsquigarrow finish proof sketch

(see also Deufhard, Theorem 5.2)

Last time

↳ looked $A+TC$, $C \in \mathbb{C}^{n \times n}$

↳ cont. diff

$$\lambda: (-\varepsilon, \varepsilon) \rightarrow \mathbb{C}$$

so that $\lambda(0) = \lambda_0$ and $\lambda(t)$ is a simple EV of $A+TC$

↳ cont. diff

$$x: (-\varepsilon, \varepsilon) \rightarrow \mathbb{C}^n$$

so that $x(0) = x_0$ and $x(t)$ is EV of $A+TC$

Now we can differentiate w.r.t. t the equation

$$(A+TC)x(t) = \lambda(t)x(t)$$

to obtain

$$Cx(t) + (A+TC)x'(t) = \lambda'(t)x(t) + \lambda(t)x'(t)$$

At $t=0$

$$Cx_0 + Ax'(0) = \lambda'(0)x_0 + \lambda_0 x'(0)$$

Multiply from right with y_0

$$\langle Cx_0, y_0 \rangle + \underline{\langle Ax'(0), y_0 \rangle} = \langle \lambda'(0)x_0, y_0 \rangle + \underline{\langle \lambda_0 x'(0), y_0 \rangle}$$

Now use that

$$\langle \lambda'(0)x_0, y_0 \rangle = \lambda'(0) \langle x_0, y_0 \rangle$$

and

$$\begin{aligned} \langle Ax'(0), y_0 \rangle &= \langle x'(0), A^T y_0 \rangle \\ &= \langle x'(0), \bar{\lambda}_0 y_0 \rangle \\ &= \langle \lambda_0 x'(0), y_0 \rangle \end{aligned}$$

Obtain derivative

$$\lambda'(0) = \frac{\langle Cx_0, y_0 \rangle}{\langle x_0, y_0 \rangle}$$

$$\lambda'(A)C = \lambda'(0) = \frac{\langle Cx_0, y_0 \rangle}{\langle x_0, y_0 \rangle}$$

The absolute condition number of determining a simple eigenvalue λ_0 of a matrix $A \in \mathbb{C}^{n \times n}$ with respect to the 2-norm is

$$\kappa_{\text{abs}} = \|\lambda'(A)\| = \frac{\|x\| \|y\|}{|\langle x, y \rangle|} = \frac{1}{|\cos(\angle(x, y))|}$$

and the rel. condition number

$$\kappa_{\text{rel}} = \frac{\|A\|}{\|\lambda_0\|} \|\lambda'(A)\| = \frac{\|A\|}{|\lambda_0 \cos(\angle(x, y))|}$$

where x is eigenvector of A for λ_0
 y is left eigenvector. --

For all $C \in \mathbb{C}^{n \times n}$, we have

$$|\langle Cx, y \rangle| \leq \|Cx\| \|y\| \leq \|C\| \|x\| \|y\|$$

Equality if $C = yx^H$

$$\| \mathcal{L}'(A)C \| = \sup_{C \neq 0} \frac{|\langle Cx, y \rangle / \langle x, y \rangle|}{\|C\|}$$

$$\underline{C = yx^H} \quad \frac{|\langle yx^H x, y \rangle / \langle x, y \rangle|}{\|yx^H\| [= \|x\| \|y\|]}$$

$$= \frac{| \|yx^H\|^2 / \langle x, y \rangle |}{\|y\| \|x\|}$$

$$= \frac{\|y\|^2 \|x\|^2}{\|y\| \|x\|} \frac{1}{|\langle x, y \rangle|}$$

$$= \frac{\|x\| \|y\|}{|\langle x, y \rangle|} = \frac{1}{\cos(\angle(x, y))}$$

The Power Method

Recap: Power method

Let $\mathbf{A} \in \mathbb{C}^{n \times n}$ be diagonalizable matrix and λ_1 be a simple eigenvalue with

$$|\lambda_1| > |\lambda_2| \geq \dots \geq |\lambda_n|$$

Let \mathbf{x}_0 be an initial guess that is not orthogonal to the eigenspace of λ_1 , then \mathbf{x}_k obtained via the iterations

$$\mathbf{z}_{k+1} = \mathbf{A}\mathbf{x}_k \tag{3}$$

$$\mathbf{x}_{k+1} = \mathbf{z}_{k+1} / \|\mathbf{z}_{k+1}\|_2 \tag{4}$$

will converge to the normalized eigenvector of \mathbf{A} corresponding to λ_1 for $k \rightarrow \infty$.

This process is called the power method.

Proof \rightsquigarrow last week

Recap Power method (cont'd)

Start with initial guess \mathbf{x}_0 and then iterate

- ▶ Compute matrix-vector product and normalize it

$$\mathbf{x}_k = \frac{\mathbf{A}\mathbf{x}_{k-1}}{\|\mathbf{A}\mathbf{x}_{k-1}\|}$$

- ▶ Obtain eigenvalue estimate (note that $\|\mathbf{x}_k\| = 1$)

$$\lambda_1^{(k)} = \mathbf{x}_k^H \mathbf{A} \mathbf{x}_k$$

- ▶ Test for convergence? **How?** Compute residual

$$\mathbf{r}_k = \mathbf{A}\mathbf{x}_k - \lambda_1^{(k)} \mathbf{x}_k$$

and terminate if the error estimate is small enough (bound if \mathbf{A} Hermitian, heuristic otherwise)

$$\min_i |\lambda_i - \lambda_1^{(k)}| \leq \frac{\|\mathbf{r}_k\|}{\|\mathbf{x}_k\|} < \epsilon$$

- ▶ The power method converges linearly

$$\|\mathbf{x}_k - (\pm \mathbf{v}_1)\| \in \mathcal{O}((|\lambda_2|/|\lambda_1|)^k)$$

- ▶ If \mathbf{A} is normal, then the eigenvalue estimate converges a bit faster but still linearly

$$|\lambda_1^{(k)} - \lambda_1| \in \mathcal{O}((|\lambda_2|/|\lambda_1|)^{2k})$$

- ▶ The power method is fast when the dominant eigenvalue is well separated from the rest
- ▶ This conclusion is rather general for all iterative methods, convergence is often good if eigenvalues are well separated and bad otherwise
- ▶ The power method is typically too slow to be used in practice

We have that eigenvector converges as

$$\|x_k - (\pm v_1)\| \in \mathcal{O}\left(\left|\frac{\lambda_2}{\lambda_1}\right|^k\right)$$

linear convergence: $\exists c > 0$

$$\lim_{k \rightarrow \infty} \frac{\|x - x_{k+1}\|}{\|x - x_k\|} < c, \quad c < 1$$

quadratic convergence

$$\lim_{k \rightarrow \infty} \frac{\|x - x_{k+1}\|}{\|x - x_k\|^2} < c$$

Try linear convergence

$$\begin{aligned} \lim_{k \rightarrow \infty} \frac{\|x_{k+1} - (\pm v_1)\|}{\|x_k - (\pm v_1)\|} &= \lim_{k \rightarrow \infty} \frac{\left|\frac{\lambda_2}{\lambda_1}\right|^{k+1}}{\left|\frac{\lambda_2}{\lambda_1}\right|^k} \\ &= \lim_{k \rightarrow \infty} \left|\frac{\lambda_2}{\lambda_1}\right| < c \leq 1 \end{aligned}$$

Try quadratic

$$\begin{aligned} \lim_{k \rightarrow \infty} \frac{\|x_{k+1} - (\pm v_1)\|}{\|x_k - (\pm v_1)\|^2} &= \lim_{k \rightarrow \infty} \frac{\left|\frac{\lambda_2}{\lambda_1}\right|^{k+1}}{\left|\frac{\lambda_2}{\lambda_1}\right|^{2k}} \\ &= \lim_{k \rightarrow \infty} \frac{\left|\frac{\lambda_2}{\lambda_1}\right|}{\left|\frac{\lambda_2}{\lambda_1}\right|^k} \text{ diverges} \end{aligned}$$

Now eigenvalue

$$\| \lambda_1^{(n)} - \lambda_1 \| \in O(|\lambda_2/\lambda_1|^{2n})$$

Quadratic

$$\lim_{h \rightarrow \infty} \frac{|\lambda_2/\lambda_1|^{2(n+1)}}{|\lambda_2/\lambda_1|^{2n+2}} = \lim_{h \rightarrow \infty} \frac{|\lambda_2/\lambda_1|^2}{|\lambda_2/\lambda_1|^{2n}} \text{ diverges}$$

linear

$$\lim_{h \rightarrow \infty} \frac{|\lambda_2/\lambda_1|^{2(n+1)}}{|\lambda_2/\lambda_1|^{2n}} = \lim_{h \rightarrow \infty} |\lambda_2/\lambda_1|^2 < 1$$

The inverse power method

For any μ not an eigenvalue of \mathbf{A} :

- ▶ The eigenvectors of $(\mathbf{A} - \mu\mathbf{I})^{-1}$ are the same as the eigenvectors of \mathbf{A}
- ▶ The eigenvalues of $(\mathbf{A} - \mu\mathbf{I})^{-1}$ are $\{(\lambda_j - \mu)^{-1}\} \rightsquigarrow$ why useful?

The inverse power method

For any μ not an eigenvalue of \mathbf{A} :

- ▶ The eigenvectors of $(\mathbf{A} - \mu\mathbf{I})^{-1}$ are the same as the eigenvectors of \mathbf{A}
- ▶ The eigenvalues of $(\mathbf{A} - \mu\mathbf{I})^{-1}$ are $\{(\lambda_j - \mu)^{-1}\} \rightsquigarrow$ **why useful?**

Thus, if we have a good estimate μ of an eigenvalue λ_J of matrix \mathbf{A} , then

$$\frac{|\lambda_j - \mu|^{-1}}{|\lambda_J - \mu|^{-1}} \ll 1, \quad j \neq J$$

and thus the power method applied to $(\mathbf{A} - \mu\mathbf{I})^{-1}$ converges rapidly to \mathbf{v}_J :

$$(\mathbf{A} - \mu\mathbf{I})\mathbf{y}_{k+1} = \mathbf{x}_k \tag{5}$$

$$\mathbf{x}_{k+1} = \mathbf{y}_{k+1} / \|\mathbf{y}_{k+1}\| \tag{6}$$

- ▶ What do we need to keep in mind?

The inverse power method

For any μ not an eigenvalue of \mathbf{A} :

- ▶ The eigenvectors of $(\mathbf{A} - \mu\mathbf{I})^{-1}$ are the same as the eigenvectors of \mathbf{A}
- ▶ The eigenvalues of $(\mathbf{A} - \mu\mathbf{I})^{-1}$ are $\{(\lambda_j - \mu)^{-1}\} \rightsquigarrow$ **why useful?**

Thus, if we have a good estimate μ of an eigenvalue λ_J of matrix \mathbf{A} , then

$$\frac{|\lambda_j - \mu|^{-1}}{|\lambda_J - \mu|^{-1}} \ll 1, \quad j \neq J$$

and thus the power method applied to $(\mathbf{A} - \mu\mathbf{I})^{-1}$ converges rapidly to \mathbf{v}_J :

$$(\mathbf{A} - \mu\mathbf{I})\mathbf{y}_{k+1} = \mathbf{x}_k \tag{5}$$

$$\mathbf{x}_{k+1} = \mathbf{y}_{k+1} / \|\mathbf{y}_{k+1}\| \tag{6}$$

- ▶ What do we need to keep in mind? Costs: Requires matrix solve in every iteration; same matrix, different right-hand sides (\rightsquigarrow LU and Cholesky decompositions)
- ▶ This algorithm is used in practice to find *eigenvectors* if the *eigenvalues* are already known

Rayleigh quotient iterations

- ▶ The convergence speed of the inverse power method increases with a better eigenvalue estimate \rightsquigarrow **what could we do?**

Rayleigh quotient iterations

- ▶ The convergence speed of the inverse power method increases with a better eigenvalue estimate \rightsquigarrow **what could we do?**
- ▶ Combine estimating eigenvalue and eigenvector \rightsquigarrow Rayleigh quotient iteration

Accelerated version of the inverse power method using changing shifts:

- ▶ Choose starting vector \mathbf{x}^0 with $\|\mathbf{x}^0\| = 1$. Compute $\lambda^{(0)} = (\mathbf{x}^0)^T A \mathbf{x}^0$.
- ▶ For $i = 0, 1, \dots$ do

$$(A - \lambda^{(k)} I) \mathbf{x}^{k+1} = \mathbf{x}^k, \quad \mathbf{y}^{k+1} = \mathbf{x}^{k+1} / \|\mathbf{x}^{k+1}\|.$$

- ▶ Compute $\lambda^{(k+1)} = (\mathbf{y}^{k+1})^T A \mathbf{y}^{k+1}$, and go back.

If it converges (depends on starting point), then it converges *cubically* \rightsquigarrow details in Trefethen & Bau

(This is the only method we will see that converges so quickly!)

The QR algorithm

The QR algorithm

The power method is not well suited for finding all eigenvalues of a matrix \mathbf{A}

Idea of the QR method: Build a matrix \mathbf{A}' that shares the eigenvalues of \mathbf{A} via similarity transformations

$$\mathbf{A}' = \mathbf{P}^{-1}\mathbf{A}\mathbf{P}, \quad \mathbf{A}, \mathbf{A}' \text{ have the same eigenvalues}$$

and for which we know the eigenvalues. \rightsquigarrow What matrix would we like \mathbf{A}' to be?

The QR algorithm

The power method is not well suited for finding all eigenvalues of a matrix \mathbf{A}

Idea of the QR method: Build a matrix \mathbf{A}' that shares the eigenvalues of \mathbf{A} via similarity transformations

$$\mathbf{A}' = \mathbf{P}^{-1}\mathbf{A}\mathbf{P}, \quad \mathbf{A}, \mathbf{A}' \text{ have the same eigenvalues}$$

and for which we know the eigenvalues. \rightsquigarrow **What matrix would we like \mathbf{A}' to be?**

The **QR algorithm** for finding eigenvalues is as follows ($A_0 := A$), and for $k = 0, 1, \dots$:

- ▶ Compute QR decomposition of A_k , i.e., $A_k = Q_k R_k$.
- ▶ $A_{k+1} := R_k Q_k$, $k := k + 1$ and go back.

All iterates A_1, A_2, \dots have the same eigenvalues because

$$Q_k A_{k+1} Q_k^T = Q_k R_k Q_k Q_k^T = Q_k R_k = A_k$$

and A_k converges to a *diagonal* matrix if A is Hermitian and eigenvalues well separated

Intuition why QR method converges

Think of it as the power method applied to many linearly independent vectors $z_1^{(0)}, \dots, z_n^{(0)}$ at once

Define

$$Z^{(0)} = \begin{bmatrix} | & & | \\ z_1^{(0)} & \dots & z_n^{(0)} \\ | & & | \end{bmatrix}$$

and define

$$Z^{(k)} = A^k Z^{(0)} = \begin{bmatrix} | & & | \\ z_1^{(k)} & \dots & z_n^{(k)} \\ | & & | \end{bmatrix}$$

Recall that in the power method we had to re-normalize after each step \rightsquigarrow now we have multiple vectors and therefore also need to orthogonalize \rightsquigarrow QR

With orthogonalization after each iteration, we obtain the algorithm

1. $Z^{(k)} = A\bar{Q}^{(k-1)}$
2. $\bar{Q}^{(k)}R^{(k)} = Z^{(k)}$
3. $A^{(k)} = (\bar{Q}^{(k)})^T A\bar{Q}^{(k)}$

↪ equivalent to the QR method

Summary: Let the QR algorithm be applied to a symmetric real matrix A with well separated eigenvalues

$$|\lambda_1| > |\lambda_2| > \dots > |\lambda_n|$$

and eigenvectors matrix V that has nonsingular leading principal submatrices (all upper-left 1×1 , 2×2 , ... submatrices). Then, for $k \rightarrow \infty$, the iterates $A^{(k)}$ converge linearly in

$$\mathcal{O} \left(\max_j \frac{|\lambda_{j+1}|^k}{|\lambda_j|^k} \right)$$

to $\text{diag}(\lambda_1, \dots, \lambda_n)$ and $\bar{Q}^{(k)}$ to V (up to \pm)

- ▶ The convergence of the QR algorithm is closely related to that of the power method: It is only fast if all eigenvalues are well separated
- ▶ For more general (e.g., non-symmetric) matrices in complex arithmetic, the algorithm converges to the Schur decomposition $A = UTU^H$ with triangular matrix T and unitary matrix $U \rightsquigarrow$ read eigenvalues from diagonal of T
- ▶ The work *per iteration* of the basic QR algorithm that we discussed is in $\mathcal{O}(n^3)$ because of the QR factorization in each step; the power method has cost $\mathcal{O}(n^2)$ (mat-vec) per iteration
- ▶ There are several key improvements to the basic QR algorithm that bring down the cost per iteration to $\mathcal{O}(n^2)$ (Hessenberg matrices)
- ▶ There also can be shifts (compare power method) to accelerate the convergence
- ▶ As always with linear algebra routines, the “best” are implemented in LAPACK and can be called via Matlab, numpy, etc

Eigenvalues in Matlab

- In MATLAB, sophisticated variants of the **QR algorithm** (LAPACK library) are implemented in the function *eig*:

$$\Lambda = \text{eig}(A)$$

$$[X, \Lambda] = \text{eig}(A)$$

- For large or sparse matrices, iterative methods based on the **Arnoldi iteration** (ARPACK library), can be used to obtain a few of the largest/smallest/closest-to- μ eigenvalues:

$$\Lambda = \text{eigs}(A, n_{\text{eigs}})$$

$$[X, \Lambda] = \text{eigs}(A, n_{\text{eigs}})$$

- The **Schur decomposition** is provided by $[U, T] = \text{schur}(A)$.

Conclusions/summary

- Eigenvalues are **well-conditioned** for **unitarily diagonalizable matrices** (includes Hermitian matrices), but ill-conditioned for nearly non-diagonalizable matrices.
- Eigenvectors are **well-conditioned** only when **eigenvalues are well-separated**.
- Eigenvalue algorithms are **always iterative**.
- The **power method** and its variants can be used to find the **largest or smallest eigenvalue**, and they converge fast if there is a **large separation** between the target eigenvalue and nearby ones.
- Estimating **all eigenvalues and/or eigenvectors** can be done by combining the power method with *QR* factorizations.
- MATLAB has high-quality implementations of sophisticated variants of these algorithms.

Numerical Methods I

MATH-GA 2010.001/CSCI-GA 2420.001

Benjamin Peherstorfer
Courant Institute, NYU

Based on slides by G. Stadler and A. Donev

Today

Last time

- ▶ Computing eigenvalues

Today

- ▶ Singular value decomposition

Announcements

- ▶ Homework 3 posted, is due Mon, Oct 21 before class

Recap: Power method for computing eigenvectors and eigenvalues

Let $\mathbf{A} \in \mathbb{C}^{n \times n}$ be diagonalizable matrix and λ_1 be a simple eigenvalue with

$$|\lambda_1| > |\lambda_2| \geq \dots \geq |\lambda_n|$$

Let \mathbf{x}_0 be an initial guess that is not orthogonal to the eigenspace of λ_1 , then \mathbf{x}_k obtained via the iterations

$$\mathbf{z}_{k+1} = \mathbf{A}\mathbf{x}_k \tag{7}$$

$$\mathbf{x}_{k+1} = \mathbf{z}_{k+1} / \|\mathbf{z}_{k+1}\|_2 \tag{8}$$

will converge to the normalized eigenvector of \mathbf{A} corresponding to λ_1 for $k \rightarrow \infty$.

This process is called the power method.

Recap: The QR algorithm

The power method is not well suited for finding all eigenvalues of a matrix \mathbf{A}

Idea of the QR method: Build a matrix \mathbf{A}' that shares the eigenvalues of \mathbf{A} via similarity transformations

$$\mathbf{A}' = \mathbf{P}^{-1}\mathbf{A}\mathbf{P}, \quad \mathbf{A}, \mathbf{A}' \text{ have the same eigenvalues}$$

and for which we know the eigenvalues. \rightsquigarrow What matrix would we like \mathbf{A}' to be?

Recap: The QR algorithm

The power method is not well suited for finding all eigenvalues of a matrix \mathbf{A}

Idea of the QR method: Build a matrix \mathbf{A}' that shares the eigenvalues of \mathbf{A} via similarity transformations

$$\mathbf{A}' = \mathbf{P}^{-1}\mathbf{A}\mathbf{P}, \quad \mathbf{A}, \mathbf{A}' \text{ have the same eigenvalues}$$

and for which we know the eigenvalues. \rightsquigarrow **What matrix would we like \mathbf{A}' to be?**

The **QR algorithm** for finding eigenvalues is as follows ($A_0 := A$), and for $k = 0, 1, \dots$:

- ▶ Compute QR decomposition of A_k , i.e., $A_k = Q_k R_k$.
- ▶ $A_{k+1} := R_k Q_k$, $k := k + 1$ and go back.

All iterates A_1, A_2, \dots have the same eigenvalues because

$$Q_k A_{k+1} Q_k^T = Q_k R_k Q_k Q_k^T = Q_k R_k = A_k$$

and A_k converges to a *diagonal* matrix if A is Hermitian and eigenvalues well separated

Singular Value Decomposition (SVD)

Let $A \in \mathbb{C}^{m \times n}$. A singular value decomposition of A is a factorization

$$A = U\Sigma V^H,$$

where

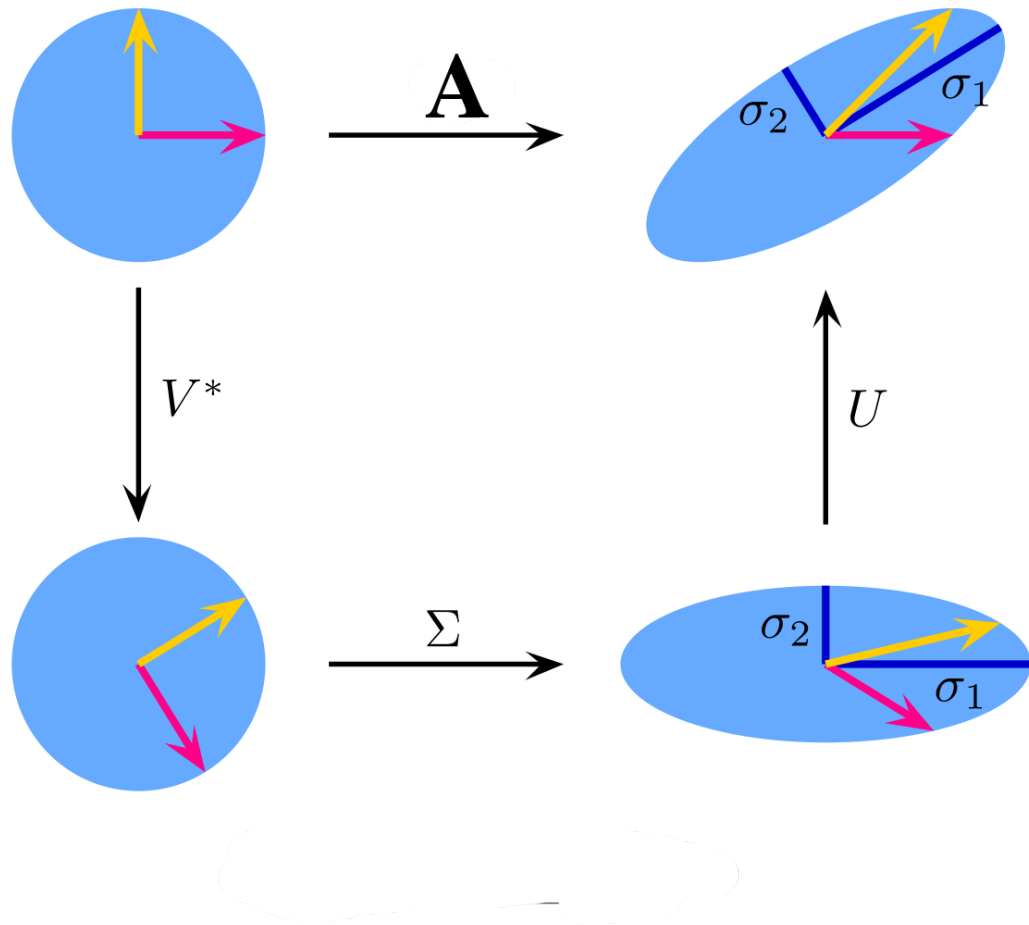
$$U \in \mathbb{C}^{m \times m} \text{ is unitary} \quad (9)$$

$$V \in \mathbb{C}^{n \times n} \text{ is unitary} \quad (10)$$

$$\Sigma \in \mathbb{R}^{m \times n} \text{ is diagonal.} \quad (11)$$

Additionally, the diagonal entries σ_j of Σ are non-negative and in non-decreasing order so that $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_p \geq 0$ where $p \in \min(m, n)$.

- ▶ The diagonal matrix Σ is real and has the same shape as A even when A is not square
- ▶ The matrices U and V are always square



The image of the unit sphere under a map A is a hyperellipse (in \mathbb{R}^m). Thus, with $A = U\Sigma V^H$, have

- ▶ The unitary map V^H preserves the sphere (rotating a sphere is a sphere)
- ▶ The diagonal matrix Σ stretches the sphere into a hyperellipse aligned with the canonical basis
- ▶ The unitary map U rotates or reflects the hyperellipse without changing shape

[Edited from figure by Georg-Johann, Wikipedia]

Existence and uniqueness of SVD

Theorem: Every matrix $A \in \mathbb{C}^{m \times n}$ has a singular value decomposition. Furthermore, the singular values $\{\sigma_j\}$ are uniquely determined, and, if A is square and the σ_j are distinct, the left and right singular vectors $\{u_j\}$ and $\{v_j\}$ are uniquely determined up to complex signs (i.e., complex scalar factors of absolute value 1.) Proof in Trefethen & Bau.

Reduced SVD

$$\underbrace{\begin{bmatrix} A \end{bmatrix}}_{m \times n} = \underbrace{\begin{bmatrix} \hat{U} \end{bmatrix}}_{m \times n} \underbrace{\begin{bmatrix} \hat{\Sigma} \end{bmatrix}}_{n \times n} \underbrace{\begin{bmatrix} \hat{V}^H \end{bmatrix}}_{n \times n}$$

SVD vs. eigenvalue decomposition

SVD expresses a matrix in proper bases for domain and range space to represent it as a diagonal matrix

$$A = U\Sigma V^H$$

We have seen something similar with eigenvectors: A non-defective square matrix A can be expressed as a diagonal matrix of eigenvalues Λ if the range and domain are presented in a basis of eigenvectors

$$A = X\Lambda X^{-1}$$

There are fundamental differences between the SVD and eigenvalue decomposition

- ▶ SVD uses two different bases (left and right singular vectors); eigenvalue decomposition uses just one (eigenvectors)
- ▶ SVD uses orthonormal bases, whereas eigenvalue basis generally is not orthogonal
- ▶ Not all matrices (even square ones) have an eigendecomposition, but all matrices (even rectangular ones) have a singular value decomposition

Typically, eigenvalues tell us something about the behavior of iterative processes that involve the matrix A such as A^k and e^{tA}

Singular values tend to tell us something about A itself

The SVD and matrix properties

In the following:

- ▶ The matrix A is of dimension $m \times n$
- ▶ $p = \min(m, n)$
- ▶ $r \leq p$ is the number of non-zero singular values of A

We now list how the SVD is related to fundamental properties of the matrix A

- ▶ The rank of A is r , the number of non-zero singular values.
- ▶ The range (column span) of A is $\text{span}(u_1, \dots, u_r)$, the kernel is $\text{span}(v_{r+1}, \dots, v_n)$
- ▶ $\|A\|_2 = \sigma_1$ and $\|A\|_F = \sqrt{\sigma_1^2 + \dots + \sigma_r^2}$
- ▶ For square A , $|\det(A)| = \prod_{i=1}^m \sigma_i$
- ▶ The non-zero singular values of A are the square roots of the non-zero eigenvalues of $A^H A$ and $AA^H \rightsquigarrow$ **proof**

The non-zero SV of A are the square roots of the non-zero EV of $A^H A$ and $A A^H$

$$\begin{aligned} A^H A &= (U \Sigma V^H)^H (U \Sigma V^H) \\ &= V \Sigma^H \underbrace{U^H U}_{I} \Sigma V^H \\ &= V \Sigma^H \Sigma V^H \end{aligned}$$

Need that $\Sigma^H \Sigma$ is diagonal, then we have eigen-decomposition of $A^H A$

$\Sigma^H \Sigma$ has size $n \times n$ and the diagonal $\sigma_1^2, \dots, \sigma_p^2$

$$\Lambda = \begin{bmatrix} \sigma_1^2 & & \\ & \ddots & \\ & & \sigma_p^2 \end{bmatrix} = \Sigma^H \Sigma$$

Full SVD

$$\underbrace{\mathbf{A}}_{m \times n} = \underbrace{\mathbf{U}}_{m \times m} \underbrace{\mathbf{\Sigma}}_{m \times n} \underbrace{\mathbf{V}^H}_{n \times n}$$

Reduced SVD

$$\underbrace{\mathbf{A}}_{m \times n} = \underbrace{\mathbf{U}}_{m \times n} \underbrace{\mathbf{\Sigma}}_{n \times n} \underbrace{\mathbf{V}^H}_{n \times n}$$

Rank-revealing SVD of a rank r matrix with dimension $m \times n$

$$\underbrace{\mathbf{A}}_{m \times n} = \underbrace{\mathbf{U}}_{m \times r} \underbrace{\mathbf{\Sigma}}_{r \times r} \underbrace{\mathbf{V}^H}_{r \times n}$$

Representing matrices via sums of rank-one matrices

Represent A as a sum of M rank-one matrices

$$A = \sum_{j=1}^M A^{(j)}$$

There are many possibilities of choosing $A^{(j)}$

- ▶ Let $A^{(j)}$ contain the j -th of the m rows of A
- ▶ Let $A^{(j)}$ contain the j -th of the n columns of A
- ▶ Let $A^{(j)}$ contain one of the mn entries of A
- ▶ ...

What is a property of a “sum representation” that we like to see in numerical analysis?

Representing matrices via sums of rank-one matrices

Represent A as a sum of M rank-one matrices

$$A = \sum_{j=1}^M A^{(j)}$$

There are many possibilities of choosing $A^{(j)}$

- ▶ Let $A^{(j)}$ contain the j -th of the m rows of A
- ▶ Let $A^{(j)}$ contain the j -th of the n columns of A
- ▶ Let $A^{(j)}$ contain one of the mn entries of A
- ▶ ...

What is a property of a “sum representation” that we like to see in numerical analysis?

↪ we can truncate it after a few terms and get a good approximation

$$A \approx \sum_{j=1}^{M'} A^{(j)}$$

with $M' \ll M$

SVD for low-rank approximation

Consider the SVD $A = U\Sigma V^H$, then

$$A = \sum_{j=1}^r \sigma_j u_j v_j^H$$

Let us truncate the sum after $1 \leq q \leq r$ terms and define

$$A_q = \sum_{j=1}^q \sigma_j u_j v_j^H.$$

Then, A_q is a best rank q approximation of A in the $\|\cdot\|_2$ norm

$$\|A - A_q\|_2 = \inf_{\substack{B \in \mathbb{C}^{m \times n} \\ \text{rank}(B) \leq q}} \|A - B\|_2$$

\rightsquigarrow proof

Best rank q approximation:

Consider $A = U \Sigma V^H$

$$A_r = \sum_{j=1}^r \sigma_j u_j v_j^H$$

want show

$$\begin{cases} (1) \|A - A_q\|_2 = \sigma_{q+1} \\ (2) \|A - A_q\|_2 = \inf_{\substack{B \in \mathbb{C}^{m \times n} \\ \text{rank}(B) \leq q}} \|A - B\|_2 \end{cases}$$

$$(1) A - A_q = \sum_{i=q+1}^r \sigma_i u_i v_i^H$$

$$A - A_q = \begin{bmatrix} | & & | \\ u_{q+1} & \dots & u_r \\ | & & | \end{bmatrix} \begin{bmatrix} \sigma_{q+1} & & \\ & \dots & \\ & & \sigma_r \end{bmatrix} \begin{bmatrix} - & v_{q+1}^H & - \\ & \vdots & \\ - & v_r^H & - \end{bmatrix}$$

We found the SVD of $A - A_q$ with SV

$$\sigma_{q+1} \geq \dots \geq \sigma_r$$

Because $\|Z\|_2 = \sigma_{\max}(Z)$

$$\|A - A_q\|_2 = \sigma_{q+1}$$

Suppose there is B with $\text{rank}(B) \leq q$

and

$$\|A - B\|_2 \leq \|A - A_q\|_2 = \sigma_{q+1} \quad (*)$$

Because $B \in \mathbb{C}^{n \times n}$ has $\text{rank} \leq q$, there is an $n - q$ dimensional null space $W \subseteq \mathbb{C}^n$ such that

$$w \in W \Rightarrow Bw = 0$$

for any $w \in W$ have

$$Aw = (A - B)w$$

and

$$\begin{aligned} \|Aw\|_2 &= \|(A - B)w\|_2 \leq \|A - B\|_2 \|w\|_2 \\ &\leq \underline{\underline{\sigma_{q+1}}} \|w\|_2 \end{aligned}$$

Now take v_i , then

$$\begin{aligned} Av_i &= U \Sigma V^H v_i = \\ &= U \Sigma \begin{bmatrix} 0 \\ \vdots \\ 0 \\ \vdots \\ 0 \end{bmatrix} = U \begin{bmatrix} 0 \\ \vdots \\ \sigma_i \\ \vdots \\ 0 \end{bmatrix} = \sigma_i u_i \end{aligned}$$

The vectors v_1, \dots, v_{q+1} span a $q+1$ dim space $Z \subseteq \mathbb{C}^n$

$$\forall z \in Z: \quad z = \sum_{i=1}^{q+1} \alpha_i v_i, \quad \|z\|_2 = \sqrt{\sum_{i=1}^{q+1} |\alpha_i|^2}$$

$$Az = \sum_{i=1}^{q+1} \sigma_i v_i \alpha_i$$

with

$$\|Az\|_2 = \sqrt{\sum_{i=1}^{q+1} \sigma_i^2 \alpha_i^2 \underbrace{\|v_i\|_2^2}_{=1}} \geq \sigma_{q+1} \sqrt{\sum_{i=1}^{q+1} \alpha_i^2} = \sigma_{q+1} \|z\|_2$$

However, the subspace Z has $\dim q+1$ and the subspace W has $\dim n-q$, thus there must be a non-zero intersection, which contradicts that $\forall w \in W: \|Aw\| < \sigma_{q+1} \|w\|_2$

SVD for low-rank approximation

Consider the SVD $A = U\Sigma V^H$, then

$$A = \sum_{j=1}^r \sigma_j u_j v_j^H$$

Let us truncate the sum after $1 \leq q \leq r$ terms and define

$$A_q = \sum_{j=1}^q \sigma_j u_j v_j^H.$$

Then, A_q is a best rank q approximation of A in the $\|\cdot\|_2$ norm

$$\|A - A_q\|_2 = \inf_{\substack{B \in \mathbb{C}^{m \times n} \\ \text{rank}(B) \leq q}} \|A - B\|_2 = \sigma_{q+1}$$

The error is the first left out singular value σ_{q+1} ! \rightsquigarrow **proof**

Furthermore, in the Frobenius norm $\|\cdot\|_F$, for any $0 \leq q \leq r$, the matrix A_q from the previous slide also satisfies

$$\|A - A_q\|_F = \inf_{\substack{B \in \mathbb{C}^{m \times n} \\ \text{rank}(B) \leq q}} \|A - B\|_F = \sqrt{\sigma_{q+1}^2 + \cdots + \sigma_r^2}$$

SVD and pseudo inverse

The (Moore-Penrose) pseudo inverse of an $m \times n$ matrix that is regular and square is $A^+ = A^{-1}$.

Otherwise, the pseudo inverse is given by

$$A^+ = V\Sigma^+U^H,$$

where

$$\Sigma^+ = \text{diag}(\sigma_1^{-1}, \sigma_2^{-1}, \dots, \sigma_r^{-1}, 0, \dots, 0)$$

Thus, we can solve least-squares problems $\min_x \|b - Ax\|_2$ by taking the SVD of A , computing A^+ , and setting $x = A^+b$ for the minimal norm solution w.r.t. $\|\cdot\|_2$

The approach we discussed via the QR decomposition is cheaper but the approach via the SVD is sometimes preferred because it allows to easily regularize the problem by truncating small singular values.

How to compute the SVD?

How to compute the SVD?

The SVD of an $m \times n$ ($m \geq n$) matrix A is related to the eigenvalue decomposition of $A^H A$

$$A^H A = V \Sigma^H \Sigma V^H,$$

How to compute the SVD?

The SVD of an $m \times n$ ($m \geq n$) matrix A is related to the eigenvalue decomposition of $A^H A$

$$A^H A = V \Sigma^H \Sigma V^H,$$

Since we know how to numerically compute the eigendecomposition, we could compute the SVD of A as follows

1. Form $A^H A$
2. Compute the eigendecomposition $A^H A = V \Lambda V^H$ (Notice that $Z = A^H A$ is normal because $Z^H Z = Z Z^H$)
3. Let Σ be the $m \times n$ non-negative diagonal square root of Λ
4. Solve the system $U \Sigma = A V$ for unitary U

Is this a good idea?

How **not** to compute the SVD?

The SVD of an $m \times n$ ($m \geq n$) matrix A is related to the eigenvalue decomposition of $A^H A$

$$A^H A = V \Sigma^H \Sigma V^H,$$

Since we know how to numerically compute the eigendecomposition, we could compute the SVD of A as follows

1. Form $A^H A$
2. Compute the eigendecomposition $A^H A = V \Lambda V^H$ (Notice that $Z = A^H A$ is normal because $Z^H Z = Z Z^H$)
3. Let Σ be the $m \times n$ non-negative diagonal square root of Λ
4. Solve the system $U \Sigma = A V$ for unitary U

Is this a good idea? \rightsquigarrow as we have seen before, it is typically dangerous[†] from a stability perspective to compute something of the matrix A via the matrix $A^H A$ (think of least-squares regression) \rightsquigarrow **board**

[†]<https://nhigham.com/2022/10/11/seven-sins-of-numerical-linear-algebra/>

Motivate why computing SV via EV
is not good idea

↳ perturbation of EV of Hermitic $A^H A$

$$|\lambda_{\ell_2}(A^H A + \delta B) - \lambda_{\ell_2}(A^H A)| \leq \|\delta B\|_2 \quad \leftarrow \text{condition number}$$

↳ similar perturbation bound holds for SV

$$|\sigma_{\ell_2}(A + \delta A) - \sigma_{\ell_2}(A)| \leq \|\delta A\|$$

Backward stable algo for SV computes $\tilde{\sigma}_{\ell_2}$
that satisfies

$$\tilde{\sigma}_{\ell_2} = \sigma_{\ell_2}(A + \delta A) \quad \frac{\|\delta A\|}{\|A\|} \in \mathcal{O}(\epsilon)$$

implies

$$|\tilde{\sigma}_{\ell_2} - \sigma_{\ell_2}| \leq \|\delta A\| \quad \text{with} \quad \|\delta A\| \in \mathcal{O}(\epsilon \|A\|)$$

Now let us compute $\lambda_{\ell_2}(A^H A)$. With backward
stable algo

$$\tilde{\lambda}_{\ell_2} = \lambda_{\ell_2}(A^H A + \delta B), \quad \frac{\|\delta B\|}{\|A^H A\|} \in \mathcal{O}(\epsilon)$$

$$|\tilde{\lambda}_{\ell_2} - \lambda_{\ell_2}| \leq \|\delta B\| \in \mathcal{O}(\epsilon \|A^H A\|) = \mathcal{O}(\epsilon \|A\|^2)$$

Compute $\tilde{\sigma}_k = \sqrt{\tilde{\lambda}_k}$ and obtain

$$\begin{aligned} |\tilde{\sigma}_k - \sigma_k| &= |\sqrt{\tilde{\lambda}_k} - \sqrt{\lambda_k}| \leq \frac{|\tilde{\lambda}_k - \lambda_k|}{\underbrace{\sqrt{\tilde{\lambda}_k}}_{\geq 0} + \underbrace{\sqrt{\lambda_k}}_{\sigma_k}} \\ &\leq \frac{|\tilde{\lambda}_k - \lambda_k|}{\sigma_k} \end{aligned}$$

$|\tilde{\sigma}_k - \sigma_k|$ behaves as $O(\varepsilon \frac{\|A\|^2}{\sigma_k})$

worse by factor $\frac{\|A\|_2}{\sigma_k}$ than previous result

\Rightarrow "no" problem for large σ_k with $\sigma_k \approx \|A\|_2$
but big problem for $\sigma_k \ll \|A\|_2$

How we actually compute the SVD

For the sake of the argument, assume that A is a square $m \times m$ (following is applicable to rectangular matrices too) \rightsquigarrow board

Instead of $A^H A$, let us use

$$H = \begin{bmatrix} 0 & A^H \\ A & 0 \end{bmatrix}$$

larger matrix
 \rightarrow ignore for now

$$\begin{bmatrix} 0 & A^H \\ A & 0 \end{bmatrix} \begin{bmatrix} v & v \\ v & -v \end{bmatrix} = \begin{bmatrix} v & v \\ v & -v \end{bmatrix} \begin{bmatrix} \Sigma & 0 \\ 0 & -\Sigma \end{bmatrix}$$

Numerical Methods I

MATH-GA 2010.001/CSCI-GA 2420.001

Benjamin Peherstorfer
Courant Institute, NYU

Based on slides by G. Stadler and A. Donev

Today

Last time

- ▶ Singular value decomposition

Today

- ▶ Singular value decomposition
- ▶ Iterative methods for solving linear systems

Announcements

- ▶ Homework 3 posted, is due Mon, Oct 21 before class

Recap

Let $A \in \mathbb{C}^{m \times n}$. A singular value decomposition of A is a factorization

$$A = U\Sigma V^H,$$

where

$$U \in \mathbb{C}^{m \times m} \text{ is unitary} \quad (12)$$

$$V \in \mathbb{C}^{n \times n} \text{ is unitary} \quad (13)$$

$$\Sigma \in \mathbb{R}^{m \times n} \text{ is diagonal.} \quad (14)$$

Additionally, the diagonal entries σ_j of Σ are non-negative and in non-decreasing order so that $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_p \geq 0$ where $p \in \min(m, n)$.

- ▶ The diagonal matrix Σ is real and has the same shape as A even when A is not square
- ▶ The matrices U and V are always square

Recap

Full SVD

$$\underbrace{\mathbf{A}}_{m \times n} = \underbrace{\mathbf{U}}_{m \times m} \underbrace{\mathbf{\Sigma}}_{m \times n} \underbrace{\mathbf{V}^H}_{n \times n}$$

Reduced SVD

$$\underbrace{\mathbf{A}}_{m \times n} = \underbrace{\mathbf{U}}_{m \times n} \underbrace{\mathbf{\Sigma}}_{n \times n} \underbrace{\mathbf{V}^H}_{n \times n}$$

Rank-revealing SVD of a rank r matrix with dimension $m \times n$

$$\underbrace{\mathbf{A}}_{m \times n} = \underbrace{\mathbf{U}}_{m \times r} \underbrace{\mathbf{\Sigma}}_{r \times r} \underbrace{\mathbf{V}^H}_{r \times n}$$

Recap: The SVD and matrix properties

In the following:

- ▶ The matrix A is of dimension $m \times n$
- ▶ $p = \min(m, n)$
- ▶ $r \leq p$ is the number of non-zero singular values of A

We now list how the SVD is related to fundamental properties of the matrix A

- ▶ The rank of A is r , the number of non-zero singular values.
- ▶ The range (column span) of A is $\text{span}(u_1, \dots, u_r)$, the kernel is $\text{span}(v_{r+1}, \dots, v_n)$
- ▶ $\|A\|_2 = \sigma_1$ and $\|A\|_F = \sqrt{\sigma_1^2 + \dots + \sigma_r^2}$
- ▶ For square A , $|\det(A)| = \prod_{i=1}^m \sigma_i$
- ▶ The non-zero singular values of A are the square roots of the non-zero eigenvalues of $A^H A$ and AA^H

Recap: SVD for low-rank approximation

Consider the SVD $A = U\Sigma V^H$, then

$$A = \sum_{j=1}^r \sigma_j u_j v_j^H$$

Let us truncate the sum after $1 \leq q \leq r$ terms and define

$$A_q = \sum_{j=1}^q \sigma_j u_j v_j^H.$$

Then, A_q is a best rank q approximation of A in the $\|\cdot\|_2$ norm

$$\|A - A_q\|_2 = \inf_{\substack{B \in \mathbb{C}^{m \times n} \\ \text{rank}(B) \leq q}} \|A - B\|_2 = \sigma_{q+1}$$

The error is the first left out singular value σ_{q+1} !

Recap

Furthermore, in the Frobenius norm $\|\cdot\|_F$, for any $0 \leq q \leq r$, the matrix A_q from the previous slide also satisfies

$$\|A - A_q\|_F = \inf_{\substack{B \in \mathbb{C}^{m \times n} \\ \text{rank}(B) \leq q}} \|A - B\|_F = \sqrt{\sigma_{q+1}^2 + \cdots + \sigma_r^2}$$

Recap: SVD and pseudo inverse

The (Moore-Penrose) pseudo inverse of an $m \times n$ matrix that is regular and square is $A^+ = A^{-1}$.

Otherwise, the pseudo inverse is given by

$$A^+ = V\Sigma^+U^H,$$

where

$$\Sigma^+ = \text{diag}(\sigma_1^{-1}, \sigma_2^{-1}, \dots, \sigma_r^{-1}, 0, \dots, 0)$$

Thus, we can solve least-squares problems $\min_x \|b - Ax\|_2$ by taking the SVD of A , computing A^+ , and setting $x = A^+b$ for the minimal norm solution w.r.t. $\|\cdot\|_2$

The approach we discussed via the QR decomposition is cheaper but the approach via the SVD is sometimes preferred because it allows to easily regularize the problem by truncating small singular values.

How to compute the SVD?

How to compute the SVD?

The SVD of an $m \times n$ ($m \geq n$) matrix A is related to the eigenvalue decomposition of $A^H A$

$$A^H A = V \Sigma^H \Sigma V^H,$$

How to compute the SVD?

The SVD of an $m \times n$ ($m \geq n$) matrix A is related to the eigenvalue decomposition of $A^H A$

$$A^H A = V \Sigma^H \Sigma V^H,$$

Since we know how to numerically compute the eigendecomposition, we could compute the SVD of A as follows

1. Form $A^H A$
2. Compute the eigendecomposition $A^H A = V \Lambda V^H$ (Notice that $Z = A^H A$ is normal because $Z^H Z = Z Z^H$)
3. Let Σ be the $m \times n$ non-negative diagonal square root of Λ
4. Solve the system $U \Sigma = A V$ for unitary U

Is this a good idea?

How **not** to compute the SVD?

The SVD of an $m \times n$ ($m \geq n$) matrix A is related to the eigenvalue decomposition of $A^H A$

$$A^H A = V \Sigma^H \Sigma V^H,$$

Since we know how to numerically compute the eigendecomposition, we could compute the SVD of A as follows

1. Form $A^H A$
2. Compute the eigendecomposition $A^H A = V \Lambda V^H$ (Notice that $Z = A^H A$ is normal because $Z^H Z = Z Z^H$)
3. Let Σ be the $m \times n$ non-negative diagonal square root of Λ
4. Solve the system $U \Sigma = A V$ for unitary U

Is this a good idea? \rightsquigarrow as we have seen before, it is typically dangerous[†] from a stability perspective to compute something of the matrix A via the matrix $A^H A$ (think of least-squares regression) \rightsquigarrow **board**

[†]<https://nhigham.com/2022/10/11/seven-sins-of-numerical-linear-algebra/>

How we actually compute the SVD

For the sake of the argument, assume that A is a square $m \times m$ (following is applicable to rectangular matrices too) \rightsquigarrow board

Instead of $A^H A$, let us use

$$H = \begin{bmatrix} 0 & A^H \\ A & 0 \end{bmatrix}$$

Because $A = U \Sigma V^H$ have

$$AV = U \Sigma$$

$$A^H U = V \Sigma^H U^H U = V \Sigma^H = V \Sigma$$

$$\begin{bmatrix} 0 & A^H \\ A & 0 \end{bmatrix} \underbrace{\begin{bmatrix} V & V \\ U & -U \end{bmatrix}}_X = \begin{bmatrix} A^H U & -A^H U \\ AV & AV \end{bmatrix}$$

$$= \begin{bmatrix} V \Sigma & -V \Sigma \\ U \Sigma & U \Sigma \end{bmatrix} = \begin{bmatrix} V & V \\ U & -U \end{bmatrix} \underbrace{\begin{bmatrix} \Sigma & 0 \\ 0 & -\Sigma \end{bmatrix}}_{\text{diagonal}}$$

$$\begin{bmatrix} V & V \\ U & -U \end{bmatrix}^H \begin{bmatrix} V & V \\ U & -U \end{bmatrix} = \begin{bmatrix} 2I & 0 \\ 0 & 2I \end{bmatrix}$$

Would like to find the eigendecomposition of
 $\begin{bmatrix} 0 & A^H \\ A & 0 \end{bmatrix}$ Algorithm that uses this idea
but avoids assembling H matrix

How we actually compute the SVD

For the sake of the argument, assume that A is a square $m \times m$ (following is applicable to rectangular matrices too) \rightsquigarrow board

Consider the $2m \times 2m$ Hermitian matrix

$$H = \begin{bmatrix} 0 & A^H \\ A & 0 \end{bmatrix}$$

Since $A = U\Sigma V^H$, we have

$$AV = U\Sigma \text{ and } A^H U = V\Sigma^H = V\Sigma \quad (\text{recall that singular values are real})$$

which we write in matrix form as

$$\begin{bmatrix} 0 & A^H \\ A & 0 \end{bmatrix} \begin{bmatrix} V & V \\ U & -U \end{bmatrix} = \begin{bmatrix} V & V \\ U & -U \end{bmatrix} \begin{bmatrix} \Sigma & 0 \\ 0 & -\Sigma \end{bmatrix}$$

which is the eigendecomposition of H \rightsquigarrow avoids using $A^H A$ and AA^H and is stable

Computing the SVD is typically a two-phase procedure:

First, reduce the matrix A to bidiagonal form, then diagonalize the bidiagonal matrix

$$\begin{bmatrix} * & * & * & * \\ * & * & * & * \\ * & * & * & * \\ * & * & * & * \\ * & * & * & * \\ * & * & * & * \end{bmatrix} \rightarrow \begin{bmatrix} * & * & & \\ & * & * & \\ & & * & * \\ & & & * \end{bmatrix} \rightarrow \begin{bmatrix} * & & & \\ & * & & \\ & & * & \\ & & & * \end{bmatrix}$$

- ▶ Phase 1 involves a finite number of operations that scale as $\mathcal{O}(mn^2)$
- ▶ Phase 2 (recall eigenvalue problems) is iterative but converges very quickly; in practice achieves convergence in $\mathcal{O}(n)$ to machine precision
- ▶ Thus, state-of-the-art computation of the SVD has costs that scale as $\mathcal{O}(mn^2)$

This is the celebrated Golub-Kahan approach from the 1960s

How can we bidiagonalize a matrix? Recall that we already know how to triangularize a matrix via (unitary) Householder reflection

$$\begin{array}{ccccccc}
 \begin{bmatrix} * & * & * & * \\ * & * & * & * \\ * & * & * & * \\ * & * & * & * \\ * & * & * & * \\ * & * & * & * \end{bmatrix} & \rightarrow & \begin{bmatrix} * & * & * & * \\ & * & * & * \\ & & * & * \\ & & & * \\ & & & * \\ & & & * \end{bmatrix} & \rightarrow & \begin{bmatrix} * & * & * & * \\ & * & * & * \\ & & * & * \\ & & & * \\ & & & * \\ & & & * \end{bmatrix} & \rightarrow & \begin{bmatrix} * & * & * & * \\ & * & * & * \\ & & * & * \\ & & & * \\ & & & * \\ & & & * \end{bmatrix} & \rightarrow & \begin{bmatrix} * & * & * & * \\ & * & * & * \\ & & * & * \\ & & & * \\ & & & * \\ & & & * \end{bmatrix} \\
 A & & Q_1 A & & Q_2 Q_1 A & & Q_3 Q_2 Q_1 A & & Q_4 Q_3 Q_2 Q_1 A
 \end{array}$$

How can we bidiagonalize a matrix? Recall that we already know how to triangularize a matrix via (unitary) Householder reflection

$$\begin{matrix}
 \begin{bmatrix} * & * & * & * \\ * & * & * & * \\ * & * & * & * \\ * & * & * & * \\ * & * & * & * \\ * & * & * & * \end{bmatrix} & \rightarrow & \begin{bmatrix} * & * & * & * \\ & * & * & * \\ & * & * & * \\ & * & * & * \\ & * & * & * \\ & * & * & * \end{bmatrix} & \rightarrow & \begin{bmatrix} * & * & * & * \\ & * & * & * \\ & & * & * \\ & & * & * \\ & & * & * \\ & & * & * \end{bmatrix} & \rightarrow & \begin{bmatrix} * & * & * & * \\ & * & * & * \\ & & * & * \\ & & & * \\ & & & * \\ & & & * \end{bmatrix} & \rightarrow & \begin{bmatrix} * & * & * & * \\ & * & * & * \\ & & * & * \\ & & & * \\ & & & * \\ & & & * \end{bmatrix} \\
 A & & Q_1 A & & Q_2 Q_1 A & & Q_3 Q_2 Q_1 A & & Q_4 Q_3 Q_2 Q_1 A
 \end{matrix}$$

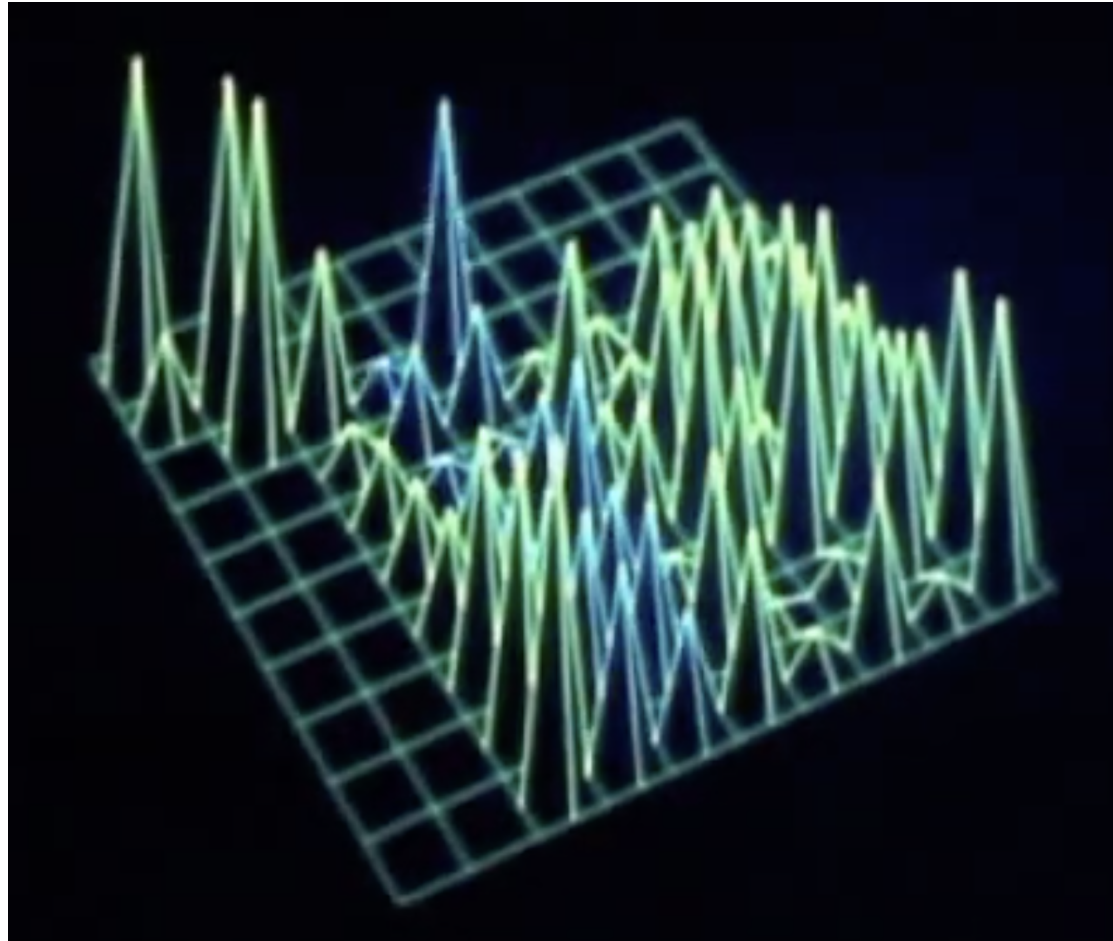
Now, we apply interleaved Householder reflection from left and right

$$\begin{matrix}
 \begin{bmatrix} * & * & * & * \\ * & * & * & * \\ * & * & * & * \\ * & * & * & * \\ * & * & * & * \\ * & * & * & * \end{bmatrix} & \rightarrow & \begin{bmatrix} * & * & * & * \\ & * & * & * \\ & * & * & * \\ & * & * & * \\ & * & * & * \\ & * & * & * \end{bmatrix} & \rightarrow & \begin{bmatrix} * & * & & \\ & * & * & * \\ & * & * & * \\ & * & * & * \\ & * & * & * \\ & * & * & * \end{bmatrix} & \rightarrow & \begin{bmatrix} * & * & & \\ & * & * & * \\ & & * & * \\ & & * & * \\ & & * & * \\ & & * & * \end{bmatrix} & \rightarrow & \begin{bmatrix} * & * & & \\ & * & * & \\ & & * & * \\ & & * & * \\ & & * & * \\ & & * & * \end{bmatrix} & \rightarrow & \dots \\
 A & & U_1^H A & & U_1^H A V_1 & & U_2^H U_1^H A V_1 & & U_2^H U_1^H A V_1 V_2
 \end{matrix}$$

At the end, n reflectors are applied from the left and $n - 2$ from the right

A variant of the QR algorithm is then applied to the bidiagonal matrix (or other eigendecomposition algorithms) \rightsquigarrow details in Golub et al., Matrix Computations.

Visualizing SVD computation



Video: <http://youtu.be/R9UoFyqJca8>

SVD saves the universe

The first Star Trek movie came out in 1979. The producers had asked Los Alamos for computer graphics to run on the displays on the bridge of the Enterprise:



[Figure: Paramount Pictures]

Link: <https://blogs.mathworks.com/cleve/2012/12/10/1976-matrix-singular-value-decomposition-film/>

Matlab

```
1: >> X = randn(10, 4);
2: >> [U, S, V] = svd(X); % full SVD
3: >> [Ur, Sr, Vr] = svd(X, 0); % reduced (economic) SVD
4: >> size(S)
5: >> size(Sr)
6: ans =
7:      10      4
8: ans =
9:       4      4
```

In most cases, we want the reduced SVD and then we should explicitly compute it:

```
1: >> X = randn(10000, 50);
2: >> tic; [U, S, V] = svd(X); toc
3: Elapsed time is 6.746243 seconds.
4: >> tic; [U, S, V] = svd(X, 0); toc
5: Elapsed time is 0.064572 seconds.
```


For very large and sparse matrices, or if the matrix is unavailable and we only have a procedure that returns the matrix-vector product, we can iteratively compute the first few singular vectors via svds

```
1: >> X = gallery('poisson', 100);
2: >> whos X
3: Name          Size          Bytes   Class
   Attributes
4:
5:  X            10000x10000      1033608  double
   sparse
6:
7: >> tic; [U, S, V] = svds(X, 1); toc
8:
9: Elapsed time is 0.307015 seconds.
```

SVD is great because its computation is (very) stable \rightsquigarrow building block for computing many basic linear algebra quantities

type 'edit rank.m' in Matlab

```
1: function r = rank(A,tol)
2: %RANK Matrix rank.
3: % RANK(A) provides an estimate of the number of linearly
4: % independent rows or columns of a matrix A.
5: %
6: % RANK(A,TOL) is the number of singular values of A
7: % that are larger than TOL. By default, TOL = max(size(A)) * eps(norm(A)
8: % ).
9: % Class support for input A:
10: % float: double, single
11:
12: % Copyright 1984–2015 The MathWorks, Inc.
13:
14: s = svd(A);
15: if nargin==1
16:     tol = max(size(A)) * eps(max(s));
17: end
18: r = sum(s > tol);
```

type 'edit pinv.m' in Matlab

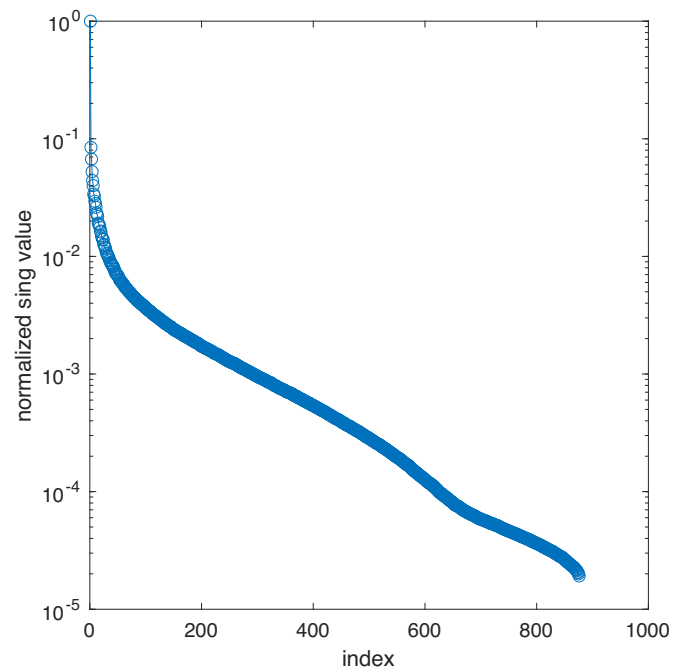
```
1: function X = pinv(A,tol)
2: %PINV   Pseudoinverse.
3: %   X = PINV(A) produces a matrix X of the same dimensions
4: %   as A' so that A*X*A = A, X*A*X = X and A*X and X*A
5: %   are Hermitian. The computation is based on SVD(A) and any
6: %   singular values less than a tolerance are treated as zero.
7: %
8: %   PINV(A,TOL) treats all singular values of A that are less than TOL as
9: %   zero. By default, TOL = max(size(A)) * eps(norm(A)).
10: %
11: %   Class support for input A:
12: %       float: double, single
13: %
14: %   See also RANK.
15:
16: %   Copyright 1984–2015 The MathWorks, Inc.
17:
18: [U,S,V] = svd(A, 'econ');
19: s = diag(S);
20: if nargin < 2
21:     tol = max(size(A)) * eps(norm(s,inf));
22: end
23: r1 = sum(s > tol)+1;
24: V(:,r1:end) = [];
25: U(:,r1:end) = [];
26: s(r1:end) = [];
27: s = 1./s(:);
28: X = (V.*s.')
```

type 'edit orth.m' in Matlab

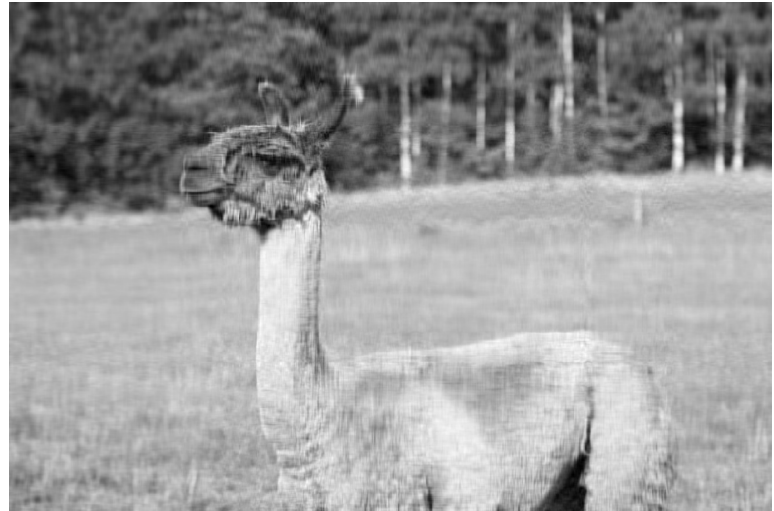
```
1: function Q = orth(A)
2: %ORTH   Orthogonalization.
3: %   Q = ORTH(A) is an orthonormal basis for the range of A.
4: %   That is,  $Q' * Q = I$ , the columns of Q span the same space as
5: %   the columns of A, and the number of columns of Q is the
6: %   rank of A.
7: %
8: %   Class support for input A:
9: %       float: double, single
10: %
11: %   See also SVD, RANK, NULL.
12:
13: %   Copyright 1984–2015 The MathWorks, Inc.
14:
15: [Q,S] = svd(A, 'econ'); %S is always square.
16: s = diag(S);
17: tol = max(size(A)) * eps(max(s));
18: r = sum(s > tol);
19: Q(:, r+1:end) = [];
```

Application of SVD for image compression

```
1: >> A = rgb2gray(imread('llama.jpg'));
2: >> figure; imshow(A);
3: >> [U, S, V] = svd(double(A));
4: >> figure; semilogy(diag(S)/S(1, 1), '-o');
5: >> xlabel('index');ylabel('normalized sing value');
6:
7: >> r = 50;
8: >> Aapprox = U(:, 1:r)*S(1:r, 1:r)*V(:, 1:r)';
9: >> figure; imshow(uint8(Aapprox));
```

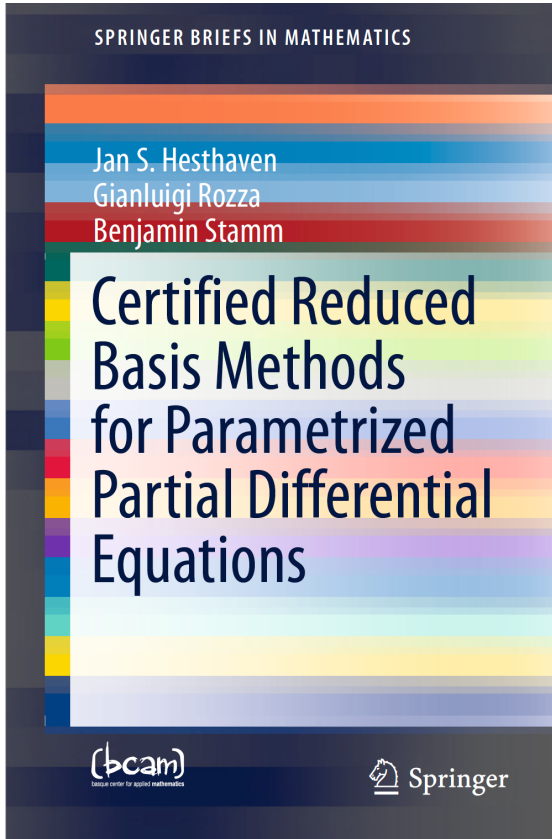


Original picture and singular values



Original picture (left) and reduced picture (right)

Outlook: Model reduction and latent dynamics



Arch Comput Methods Eng manuscript No.
(will be inserted by the editor)

G. Rozza · D.B.P. Huynh · A.T. Patera

Reduced basis approximation and a posteriori error estimation for affinely parametrized elliptic coercive partial differential equations

Application to transport and continuum mechanics

Received: August 2007 / Accepted: Date

Abstract In this paper we consider (hierarchical, Lagrange) reduced basis approximation and a *a posteriori* error estimation for linear functional outputs of affinely parametrized elliptic coercive partial differential equations. The essential ingredients are (primal-dual) Galerkin projection onto a low-dimensional space associated with a smooth “parametric manifold” — dimension reduction; efficient and effective greedy sampling methods for identification of optimal and numerically stable approximations — rapid convergence; a *a posteriori* error estimation procedure — rigorous and sharp bounds for the linear-functional outputs of interest; and Offline-Online computational decomposition strategies — minimum *marginal cost* for high performance in the real-time/embedded (e.g., parameter-estimation, control) and many-query (e.g., design optimization, multi-model/scale) contexts. We present illustrative results for heat conduction and convection-diffusion, inviscid flow, and linear elasticity; outputs include transport rates, added mass, and stress intensity factors.

Keywords Partial differential equations, parameter variation, affine geometry description, Galerkin approx-

This work was supported by DARPA/AFOSR Grants FA9550-06-1-0114 and FA-9550-07-1-0425, the Singapore-MIT Alliance, the Pappalardo MIT Mechanical Engineering Graduate Monograph Fund, and the Progetto Roberto Rozza, Politecnico di Milano-MIT. We acknowledge many helpful discussions with Professor Yvon Maday of University Paris6.

G. Rozza
Massachusetts Institute of Technology, Mechanical Engineering Department, Room 3-264, 77 Mass Avenue, Cambridge MA, 02142-4307, USA. Tel.: +1 617-452-3285; E-mail: rozza@mit.edu.

D.B.P. Huynh
National University of Singapore, Singapore-MIT Alliance, EA494-10, 4 Eng. Drive, Singapore, 117576. Tel.: +65 91324387 E-mail: baophuong@nus.edu.sg.

A.T. Patera
Massachusetts Institute of Technology, Room 3-266, 77 Mass Avenue, Cambridge MA, 02142-4307, USA. Tel.: +1 617-253-8122; E-mail: patera@mit.edu.

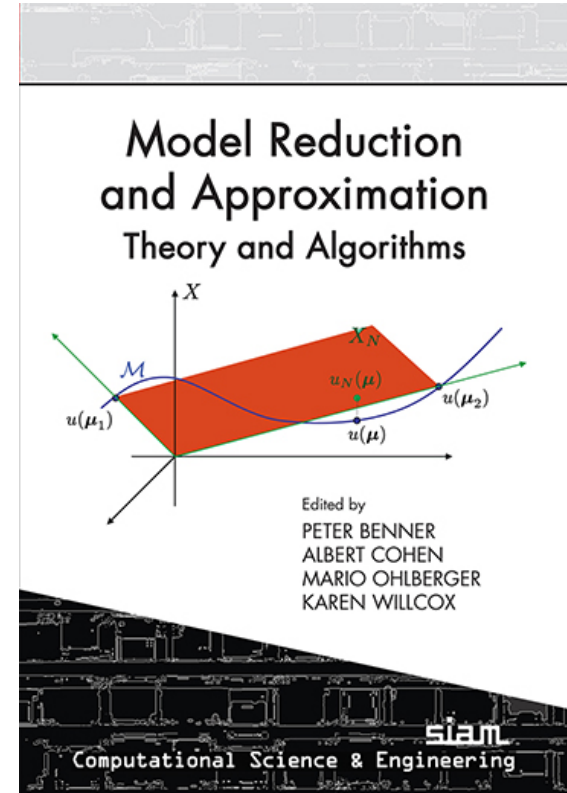
imation, a *a posteriori* error estimation, reduced basis, reduced order model, sampling strategies, POD, greedy techniques, offline-online procedures, marginal cost, coercivity lower bound, successive constraint method, real-time computation, many-query.

1 Introduction and Motivation

In this work we describe reduced basis (RB) approximation and a *a posteriori* error estimation methods for rapid and reliable evaluation of *input-output relationships* in which the *output* is expressed as a functional of a *field variable* that is the solution of an *input-parametrized* partial differential equation (PDE). In this particular paper we shall focus on linear output functionals and affinely parametrized linear elliptic coercive PDEs; however the methodology is much more generally applicable, as we discuss in Section 2.

We emphasize applications in transport and mechanics: unsteady and steady heat and mass transfer; acoustics; and solid and fluid mechanics. (Of course we do not preclude other domains of inquiry within engineering (e.g., electromagnetics) or even more broadly within the quantitative disciplines (e.g., finance).) The *input-parameter* vector typically characterizes the geometric configuration, the physical properties, and the boundary conditions and sources. The *outputs of interest* might be the maximum system temperature, an added mass coefficient, a crack stress intensity factor, an effective constitutive property, an acoustic waveguide transmission loss, or a channel flowrate or pressure drop. Finally, the *field variables* that connect the input parameters to the outputs can represent a distribution function, temperature or concentration, displacement, pressure, or velocity.

The methodology we describe in this paper is motivated by, optimized for, and applied within two particular contexts: the *real-time context* (e.g., parameter-estimation [54,96,154] or control [124]); and the *many-query context* (e.g., design optimization [107] or multi-model/scale simulation [26,49]). Both these contexts are



Iterative methods for solving systems of linear equations

Why iterative methods for linear systems?

- ▶ Costs of solving a linear system with direct method are $\mathcal{O}(n^3)$. This is too much! If n gets large, then n^3 is huge and costs become intractable.

- ▶ History of matrix computations over the years (according to Trefethen & Bau)

1950: $m = 20$

1965: $m = 200$

1980: $m = 2000$

1995: $m = 20000$

This is an increase of a factor 10^3 . However, computing power (FLOP/sec) increased by about 10^9 . Notice that $(10^3)^3 = 10^9$, which reflects the $\mathcal{O}(n^3)$ bottleneck

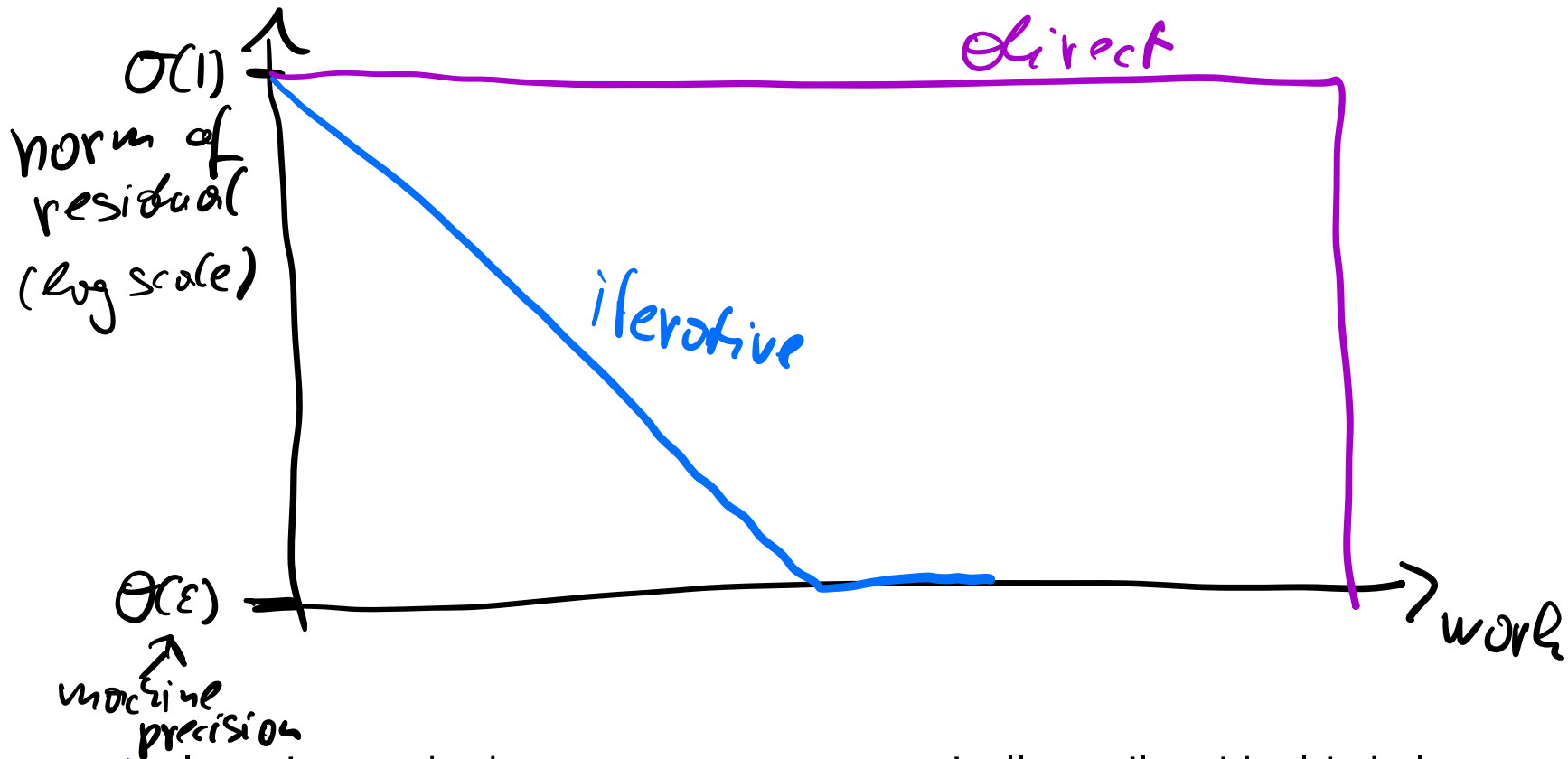
- ▶ (Side remark: There are direct methods that beat the complexity $\mathcal{O}(n^3)$ (think of Strassen's algorithm); however,

Why iterative methods for linear systems?

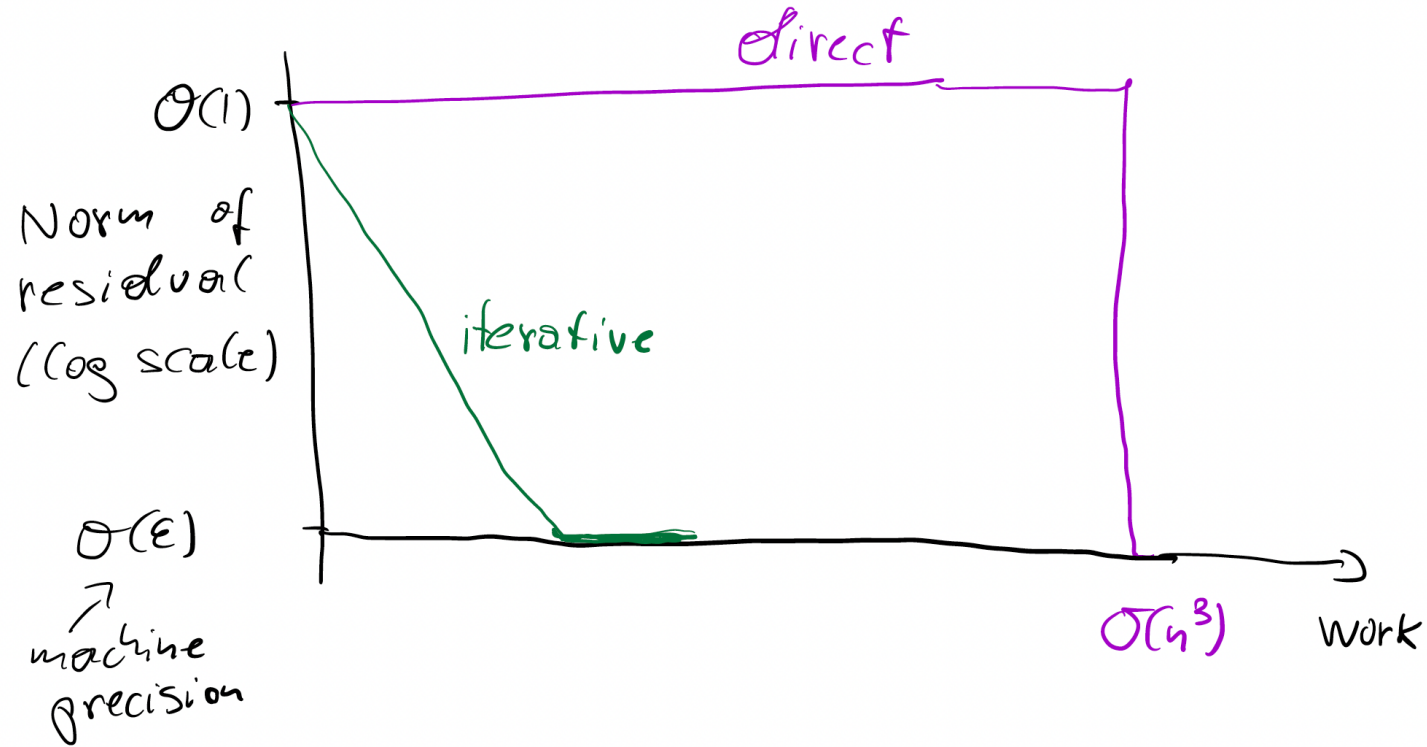
- ▶ Costs of solving a linear system with direct method are $\mathcal{O}(n^3)$. This is too much! If n gets large, then n^3 is huge and costs become intractable.
- ▶ History of matrix computations over the years (according to Trefethen & Bau)
 - 1950: $m = 20$
 - 1965: $m = 200$
 - 1980: $m = 2000$
 - 1995: $m = 20000$

This is an increase of a factor 10^3 . However, computing power (FLOP/sec) increased by about 10^9 . Notice that $(10^3)^3 = 10^9$, which reflects the $\mathcal{O}(n^3)$ bottleneck

- ▶ (Side remark: There are direct methods that beat the complexity $\mathcal{O}(n^3)$ (think of Strassen's algorithm); however, the numerical stability of these algorithms is not well understood and the constants hidden in the complexity results are huge so that we “never” see the improved rate in practice.)



- ▶ Iterative methods *can* converge geometrically until residual is below machine precision
- ▶ Direct methods make no progress at all until $O(n^3)$ work is done, and then lead to residual on the order of machine precision



- ▶ Iterative methods *can* converge geometrically until residual is below machine precision
- ▶ Direct methods make no progress at all until $O(n^3)$ work is done, and then lead to residual on the order of machine precision

Why iterative methods for linear systems? (cont'd)

- ▶ Classical direct methods (e.g., Gauss elimination) follow the pattern of taking $\mathcal{O}(n)$ steps and each step costs $\mathcal{O}(n^2)$ but don't exploit properties of the matrix
- ▶ Iterative methods reduce the number of steps and the costs of each step, depending on properties of the problem at hand (e.g., spectral properties of A in $Ax = b$)
- ▶ The ideal iterative method requires $\mathcal{O}(1)$ steps and $\mathcal{O}(n)$ costs per step to reach machine precision (think of multigrid and pre-conditioned CG)
- ▶ Iterative methods, even in the absence of rounding errors, do not deliver the exact answer (recall that LU, QR give us exact answer in finite number of steps if we are in exact arithmetic).
- ▶ However, we only can be as accurate as machine precision anyway if we do calculations on a computer. Furthermore, often matrices are stemming from discretizations of PDEs and then we need the solution only up to the discretization error.

Iterative solution of linear systems

Target problems: very **large** ($n = 10^5, 10^6, \dots$), A is usually **sparse** and has specific **properties**.

To solve

$$Ax = b$$

we construct a sequence

$$\mathbf{x}_1, \mathbf{x}_2, \dots$$

of iterates that converges fast to the solution \mathbf{x} , where \mathbf{x}_{k+1} can be cheaply computed from $\{\mathbf{x}_1, \dots, \mathbf{x}_k\}$ (e.g., one matrix-vector multiplication).

Thought experiment: If we can compute one iteration with cost $\mathcal{O}(n)$ (e.g., one matrix-vector multiplication with a sparse matrix) and need a constant $\mathcal{O}(1)$ number of iterations to reach desired precision, then we solve $Ax = b$ with costs $\mathcal{O}(n)$.

Intuitively, we cannot do better than that because we solve for n quantities and thus need to touch each at least once.

Prototype of iterative method

Let's start by trying to write $\mathbf{Ax} = \mathbf{b}$ more generally as

$$\mathbf{x} = f(\mathbf{x})$$

For example set $f(\mathbf{x}) = (\mathbf{I} - \mathbf{A})\mathbf{x} + \mathbf{b}$ to obtain

$$\mathbf{x} = f(\mathbf{x}) = (\mathbf{I} - \mathbf{A})\mathbf{x} + \mathbf{b}$$

We can now try to solve this via a fixed-point iteration

$$\mathbf{x}_{k+1} = f(\mathbf{x}_k)$$

which is in our case

$$\mathbf{x}_{k+1} = (\mathbf{I} - \mathbf{A})\mathbf{x}_k + \mathbf{b}$$

Instead of picking a specific f , let's look at the prototype

$$\mathbf{x}_{k+1} = \mathbf{G}\mathbf{x}_k + \mathbf{c}$$

where \mathbf{G} is an iteration matrix somehow related to \mathbf{A} and \mathbf{c} is related to \mathbf{b} .

We must have the actual solution $\mathbf{x} = \mathbf{A}^{-1}\mathbf{b}$ as a unique fixed point of the iteration:

$$\mathbf{x} = \mathbf{G}\mathbf{x} + \mathbf{c}$$

And we need that the sequence (\mathbf{x}_k) converges

Theorem: The fixed point method $\mathbf{x}_{k+1} = G\mathbf{x}_k + \mathbf{c}$ with an invertible G converges for each starting point \mathbf{x}_o if and only if

$$\rho(G) < 1,$$

where $\rho(G) = \max_j |\lambda_j|$ is the largest magnitude of all eigenvalue of G (i.e., the spectral radius).

~> **proof**

Symmetric G , real

There is an orthogonal Q

$$Q G Q^T = D = \begin{bmatrix} \lambda_1 & & \\ & \ddots & \\ & & \lambda_n \end{bmatrix}$$

with eigenvalues

Because $|\lambda_i| \leq \rho(G) < 1$ for all $i=1, \dots, n$ we have

$$\lim_{h \rightarrow \infty} \underbrace{\begin{bmatrix} \lambda_1^h & & \\ & \ddots & \\ & & \lambda_n^h \end{bmatrix}}_{D^h} = 0$$

Thus

$$\lim_{h \rightarrow \infty} G^h = \lim_{h \rightarrow \infty} Q D^h Q^T = 0$$

$$\begin{aligned} \|x_{k+m} - x_k\| &= \|G x_{k+m-1} + c - G x_{k-1} - c\| \\ &= \|G(x_{k+m-1} - x_{k-1})\| \\ &= \dots \\ &= \|G^m(x_k - x_0)\| \end{aligned}$$

$$\lim_{h \rightarrow \infty} \|x_{k+m} - x_k\| = 0, \quad m \geq 1$$

\rightarrow Cauchy sequence

We often ask for $\|G\| < 1$, which implies

$$\rho(G) < 1$$

because $\rho(G) \leq \|G\|$ for $\|\cdot\|$ induced by a vector norm.

$$\|x_k - x\| \leq \underline{\|G\|^k} \|x_0 - x\|$$

estimates how quickly error decays

Theorem: The fixed point method $\mathbf{x}_{k+1} = G\mathbf{x}_k + \mathbf{c}$ with an invertible G converges for each starting point \mathbf{x}_0 if and only if

$$\rho(G) < 1,$$

where $\rho(G) = \max_j |\lambda_j|$ is the largest magnitude of all eigenvalue of G (i.e., the spectral radius).

~> **proof**

In particular, for any induced matrix norm $\|\cdot\|$ we have $\rho(G) \leq \|G\|$ and thus $\|G\|$ is a sufficient criterion for convergence. We then obtain

$$\|\mathbf{x}_k - \mathbf{x}\| \leq \|G\|^k \|\mathbf{x}_0 - \mathbf{x}\|$$

Let Q be invertible, then

$$\begin{aligned} \mathbf{Ax} = \mathbf{b} &\Leftrightarrow Q^{-1}(\mathbf{b} - \mathbf{Ax}) = \mathbf{0} \\ &\Leftrightarrow (\mathbf{I} - Q^{-1}\mathbf{A})\mathbf{x} + Q^{-1}\mathbf{b} = \mathbf{x} \\ &\Leftrightarrow \mathbf{Gx} + \mathbf{c} = \mathbf{x} \end{aligned}$$

Leads to fixed-point iteration with $\mathbf{G} = (\mathbf{I} - Q^{-1}\mathbf{A})$ and $\mathbf{c} = Q^{-1}\mathbf{b}$

$$\mathbf{x}_{k+1} = \mathbf{Gx}_k + \mathbf{c}$$

and $\mathbf{x} = \mathbf{A}^{-1}\mathbf{b}$ is stationary point

$$\mathbf{Gx} + \mathbf{c} = (\mathbf{I} - Q^{-1}\mathbf{A})\mathbf{x} + Q^{-1}\mathbf{b} = \mathbf{x} - Q^{-1}\mathbf{b} + Q^{-1}\mathbf{b} = \mathbf{x}$$

Extreme cases for selecting Q

What are two extreme cases for selecting Q (recall: Q needs to be invertible)?

Extreme cases for selecting Q

What are two extreme cases for selecting Q (recall: Q needs to be invertible)?

Choose $Q = A$, then our iteration

$$\mathbf{x}_{k+1} = \mathbf{G}\mathbf{x}_k + \mathbf{c}, \quad \mathbf{G} = \mathbf{I} - \mathbf{Q}^{-1}\mathbf{A}, \quad \mathbf{c} = \mathbf{Q}^{-1}\mathbf{b}$$

becomes

$$(\mathbf{I} - \mathbf{A}^{-1}\mathbf{A})\mathbf{x}_k + \underbrace{\mathbf{A}^{-1}\mathbf{b}}_{\mathbf{x}} = \mathbf{x}_{k+1}$$

$$\mathbf{0} + \mathbf{x} = \mathbf{x}_{k+1}$$

and we are done in just a single step

$$\mathbf{x}_{k+1} = \mathbf{x}$$

Thus, if we “know the solution” (in form of having the inverse \mathbf{A}^{-1}) then no further work is needed here because we already did all the work when finding \mathbf{A}^{-1}

The other extreme is setting $Q = I$, this leads to the Richardson method

$$\mathbf{x}_{k+1} = (I - \mathbf{A})\mathbf{x}_k + \mathbf{b}$$

We have invested zero costs in finding Q (and Q^{-1}) and so we can expect that $Q = I$ will require high costs in terms of number of iterations to converge in general, if it converges at all

When does Richardson method converge? \rightsquigarrow board

Let A be spd, then

$$\rho(G) = \rho(I - A) = \max \{ |1 - \lambda_{\max}|, |1 - \lambda_{\min}| \}$$

thus need

$$\lambda_{\max}(A) < 2$$

\Rightarrow rarely usable, but important building block

Numerical Methods I

MATH-GA 2010.001/CSCI-GA 2420.001

Benjamin Peherstorfer
Courant Institute, NYU

Based on slides by G. Stadler and A. Donev

Today

Last time

- ▶ Iterative methods for systems of linear equations

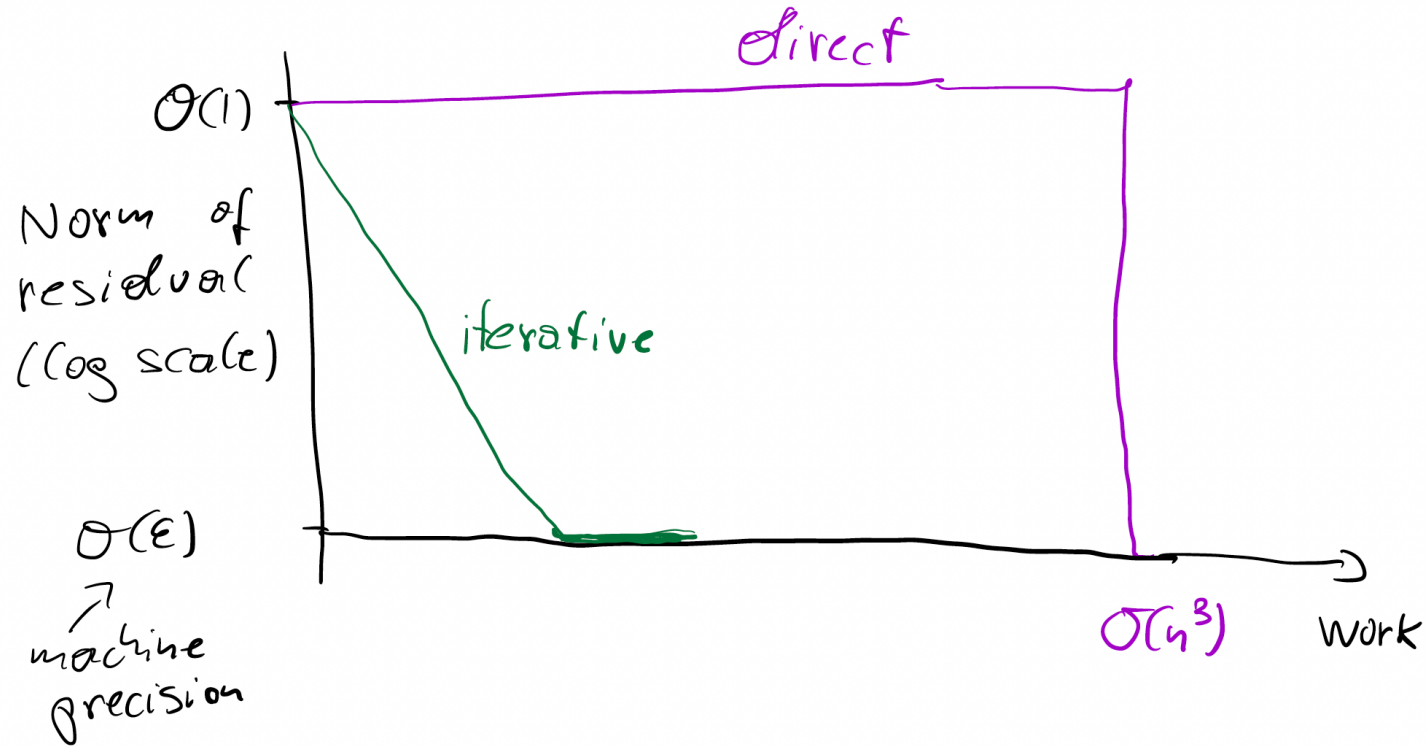
Today

- ▶ Iterative methods for systems of linear equations

Announcements

- ▶ Homework 4 is due Mon, Nov 4, 2024 before class

Recap



- ▶ Iterative methods *can* converge geometrically until residual is below machine precision
- ▶ Direct methods make no progress at all until $O(n^3)$ work is done, and then lead to residual on the order of machine precision

Recap: Iterative solution of linear systems

Target problems: very **large** ($n = 10^5, 10^6, \dots$), A is usually **sparse** and has specific **properties**.

To solve

$$Ax = b$$

we construct a sequence

$$\mathbf{x}_1, \mathbf{x}_2, \dots$$

of iterates that converges fast to the solution \mathbf{x} , where \mathbf{x}_{k+1} can be cheaply computed from $\{\mathbf{x}_1, \dots, \mathbf{x}_k\}$ (e.g., one matrix-vector multiplication).

Thought experiment: If we can compute one iteration with cost $\mathcal{O}(n)$ (e.g., one matrix-vector multiplication with a sparse matrix) and need a constant $\mathcal{O}(1)$ number of iterations to reach desired precision, then we solve $Ax = b$ with costs $\mathcal{O}(n)$.

Intuitively, we cannot do better than that because we solve for n quantities and thus need to touch each at least once.

Recap

Instead of picking a specific f , let's look at the prototype

$$\mathbf{x}_{k+1} = \mathbf{G}\mathbf{x}_k + \mathbf{c}$$

where \mathbf{G} is an iteration matrix somehow related to \mathbf{A} and \mathbf{c} is related to \mathbf{b} .

We must have the actual solution $\mathbf{x} = \mathbf{A}^{-1}\mathbf{b}$ as a unique fixed point of the iteration:

$$\mathbf{x} = \mathbf{G}\mathbf{x} + \mathbf{c}$$

And we need that the sequence (\mathbf{x}_k) converges

Recap

Let Q be invertible, then

$$\begin{aligned} \mathbf{Ax} = \mathbf{b} &\Leftrightarrow Q^{-1}(\mathbf{b} - \mathbf{Ax}) = \mathbf{0} \\ &\Leftrightarrow (\mathbf{I} - Q^{-1}\mathbf{A})\mathbf{x} + Q^{-1}\mathbf{b} = \mathbf{x} \\ &\Leftrightarrow \mathbf{Gx} + \mathbf{c} = \mathbf{x} \end{aligned}$$

Leads to fixed-point iteration with $\mathbf{G} = (\mathbf{I} - Q^{-1}\mathbf{A})$ and $\mathbf{c} = Q^{-1}\mathbf{b}$

$$\mathbf{x}_{k+1} = \mathbf{Gx}_k + \mathbf{c}$$

and $\mathbf{x} = \mathbf{A}^{-1}\mathbf{b}$ is stationary point

$$\mathbf{Gx} + \mathbf{c} = (\mathbf{I} - Q^{-1}\mathbf{A})\mathbf{x} + Q^{-1}\mathbf{b} = \mathbf{x} - Q^{-1}\mathbf{b} + Q^{-1}\mathbf{b} = \mathbf{x}$$

Recap: Extreme cases for selecting Q

What are two extreme cases for selecting Q (recall: Q needs to be invertible)?

Choose $Q = A$, then our iteration

$$\mathbf{x}_{k+1} = \mathbf{G}\mathbf{x}_k + \mathbf{c}, \quad \mathbf{G} = \mathbf{I} - \mathbf{Q}^{-1}\mathbf{A}, \quad \mathbf{c} = \mathbf{Q}^{-1}\mathbf{b}$$

becomes

$$(\mathbf{I} - \mathbf{A}^{-1}\mathbf{A})\mathbf{x}_k + \underbrace{\mathbf{A}^{-1}\mathbf{b}}_{\mathbf{x}} = \mathbf{x}_{k+1}$$

$$\mathbf{0} + \mathbf{x} = \mathbf{x}_{k+1}$$

and we are done in just a single step

$$\mathbf{x}_{k+1} = \mathbf{x}$$

Thus, if we “know the solution” (in form of having the inverse \mathbf{A}^{-1}) then no further work is needed here because we already did all the work when finding \mathbf{A}^{-1}

Recap

The other extreme is setting $Q = I$, this leads to the Richardson method

$$\mathbf{x}_{k+1} = (\mathbf{I} - \mathbf{A})\mathbf{x}_k + \mathbf{b}$$

We have invested zero costs in finding Q (and Q^{-1}) and so we can expect that $Q = I$ will require high costs in terms of number of iterations to converge in general, if it converges at all

When does Richardson method converge? \rightsquigarrow if $\lambda_{\max} < 2$, very restrictive

The Richardson method is consistent (solution is a stationary point) but it may not converge or converge very slowly

$$\mathbf{x}_{k+1} = (\mathbf{I} - \mathbf{A})\mathbf{x}_k + \mathbf{b}$$

What could we do instead of Richardson method?

The Richardson method is consistent (solution is a stationary point) but it may not converge or converge very slowly

$$\mathbf{x}_{k+1} = (\mathbf{I} - \mathbf{A})\mathbf{x}_k + \mathbf{b}$$

What could we do instead of Richardson method?

Let us think of \mathbf{Q} as a preconditioner of the Richardson method:

$$\mathbf{Q}^{-1} \approx \mathbf{A}^{-1}$$

and transform

$$\mathbf{Q}^{-1}\mathbf{A}\mathbf{x} = \mathbf{Q}^{-1}\mathbf{b}$$

Applying the Richardson iteration to this modified system gives

$$\mathbf{x}_{k+1} = (\mathbf{I} - \mathbf{Q}^{-1}\mathbf{A})\mathbf{x}_k + \mathbf{Q}^{-1}\mathbf{b} = \mathbf{G}\mathbf{x}_k + \mathbf{c}$$

Thus, with an iteration matrix $\mathbf{I} - \mathbf{Q}^{-1}\mathbf{A} \approx 0$, expect more rapid convergence

Common choices for Q

Requirements on a good Q are

Common choices for Q

Requirements on a good Q are

- ▶ Q^{-1} is representative of A^{-1}
- ▶ We can numerically solve $Qy = z$ much quicker than $Ax = b$ because in each step we solve

$$Q\mathbf{x}_{k+1} = (Q - A)\mathbf{x}_k + b$$

Common choices for Q

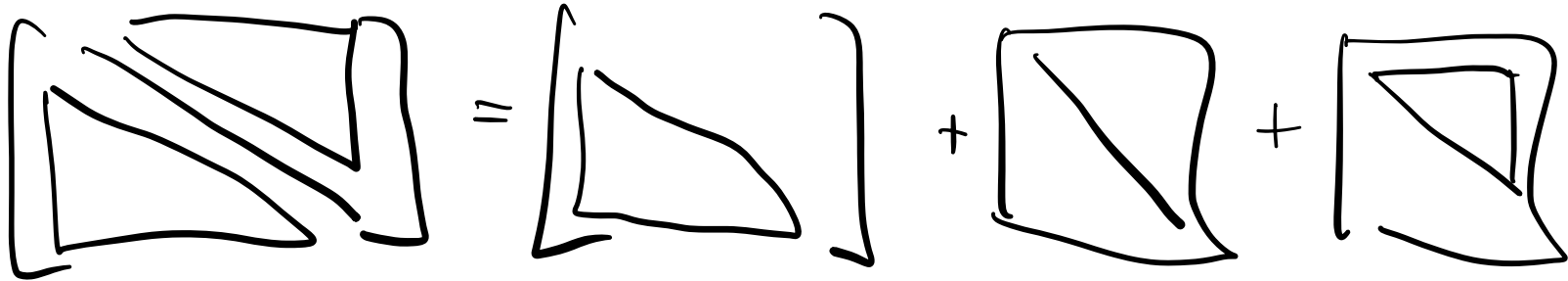
Requirements on a good Q are

- ▶ Q^{-1} is representative of A^{-1}
- ▶ We can numerically solve $Qy = z$ much quicker than $Ax = b$ because in each step we solve

$$Qx_{k+1} = (Q - A)x_k + b$$

Split the matrix A as follows

$$A = L + D + U$$



Jacobi method

Theorem: Select now $Q = D \dots$ **Jacobi method**. The Jacobi method converges for any starting point \mathbf{x}_0 to the solution of $A\mathbf{x} = \mathbf{b}$ if A is strictly diagonal dominant, i.e.,

$$|a_{ii}| > \sum_{j \neq i} |a_{ij}|, \quad \text{for } i = 1, \dots, n.$$

↪ board

Jacobi method

Theorem: Select now $Q = D$... **Jacobi method**. The Jacobi method converges for any starting point \mathbf{x}_0 to the solution of $A\mathbf{x} = \mathbf{b}$ if A is strictly diagonal dominant, i.e.,

$$|a_{ii}| > \sum_{j \neq i} |a_{ij}|, \quad \text{for } i = 1, \dots, n.$$

↪ board

Notice that $Q = D$ is a good choice in terms of computational costs because we can very quickly solve $D\mathbf{y} = \mathbf{z}$ for the diagonal matrix D

Jacobi method

$$x_{k+1} = (I - D^{-1}A)x_k + D^{-1}b$$

$$= (I - D^{-1}(L+D+R))x_k + D^{-1}b$$

$$= \underbrace{-D^{-1}(L+R)}_G x_k + D^{-1}b$$

not 12:1

||

$$\rho(-D^{-1}(L+R)) = \rho(D^{-1}(L+R))$$

$$\left. \begin{array}{l} \rho(A) \leq \|A\| \\ \text{if } \|\cdot\| \text{ is induced} \\ \text{by vector norm} \end{array} \right\} \leq \|D^{-1}(L+R)\|_{\infty} \xrightarrow{L+R}$$
$$= \max_i \sum_{j \neq i} \left| \frac{a_{ij}}{d_{ii}} \right| < 1$$

D^{-1} \uparrow diagonal dominant

Gauss-Seidel method

Theorem: Choose $Q = D + L$... **Gauss-Seidel method**. The Gauss-Seidel method converges for any starting point \mathbf{x}_0 if A is symmetric positive definite (spd).

↪ **board**

Gauss-Seidel convergence

$$\begin{aligned}x_{q+1} &= (I - (D+L)^{-1}A)x_q + (D+L)^{-1}b \\&= ((D+L)^{-1}(D+L) - (D+L)^{-1}(D+L+R))x_q + (D+L)^{-1}b \\&= -(D+L)^{-1}R x_q + (D+L)^{-1}b\end{aligned}$$

For any Spd matrix A , we have a scalar product

$$\langle x, y \rangle_A = \langle x, Ay \rangle, \quad \text{on } \mathbb{R}^n$$

and for any matrix $B \in \mathbb{R}^{n \times n}$, we have its adjoint w.r.t. $\langle \cdot, \cdot \rangle_A$ given by

$$B^* = A^{-1} B^T A$$

so that

$$\langle Bx, y \rangle_A = \langle x, B^*y \rangle_A \quad \forall x, y \in \mathbb{R}^n$$

$$\left[\downarrow \right. \\ \left. x^T B^T A y = x^T A^{-1} B^T A y = x^T A B^* y \right]$$

A self-adjoint matrix $B = B^*$ is positive

w.r.t. $\langle \cdot, \cdot \rangle_A$ iff $\langle Bx, x \rangle_A > 0 \quad \forall x \neq 0$

First we show: Let $G \in \mathbb{R}^{n \times n}$ with adjoint G^* w.r.t. $\langle \cdot, \cdot \rangle$. Then, if $B = \underline{I - G^*G}$ is positive w.r.t. $\langle \cdot, \cdot \rangle$, it follows $\rho(G) < 1$.

$\hookrightarrow B$ positive \Rightarrow

$$\forall x \neq 0$$

$$\begin{aligned} 0 < \langle Bx, x \rangle &= \langle x, x \rangle - \langle G^*Gx, x \rangle \\ &= \langle x, x \rangle - \langle Gx, Gx \rangle \end{aligned}$$

$$\|x\| > \|Gx\| \quad \forall x \neq 0$$

$$\|G\| = \sup_{\|x\|=1} \frac{\|Gx\|}{\|x\|} < 1$$

implies $\rho(G) \leq \|G\| < 1$.

Now, we can go back to our original problem of showing convergence of Gauss-Seidel for any spd matrix

\hookrightarrow We show that $B = I - G^*G$ with

$$G = I - (D+L)^{-1}A$$

is a positive matrix w.r.t. $\langle \cdot, \cdot \rangle_A$

with A spd.

\hookrightarrow Because A spd, have

↳ Because A spd, $A = L + D + R$, we have

$$\boxed{L^T = R}$$

$$G^* = A^{-1} G^T A =$$

$$= I - \underbrace{A^{-1} A^T}_{(D+R)^{-1}} \underbrace{(D+R)^{-1} A}_{(D+L)^{-T}}$$

$$= I - (D+R)^{-1} A$$

$$\underline{B = I - G^* G} = \dots = (D+R)^{-1} D (D+L)^{-1} A$$

⇒ now show that B pos. wrt. $\langle \cdot, \cdot \rangle_A$

$$\langle Bx, x \rangle_{\underline{A}} = \langle \underbrace{(D+R)^{-1} D (D+L)^{-1} A}_{(D+L)^{-T}} x, \underline{Ax} \rangle_{\underline{2}}$$

$$= \langle \underline{D^{1/2} (D+L)^{-1} Ax}, \underline{D^{1/2} (D+L)^{-1} Ax} \rangle_{\underline{2}}$$

$$= \| D^{1/2} (D+L)^{-1} Ax \|_2^2 > 0 \quad \forall x \neq 0$$

⇒ B is positive and thus

$$\rho(G) < 1$$

Gauss-Seidel method

Theorem: Choose $Q = D + L$... **Gauss-Seidel method**. The Gauss-Seidel method converges for any starting point x_0 if A is symmetric positive definite (spd).

↪ board

Notice that $Q = D + L$ is a good choice in terms of computational costs because we can very quickly solve $(D + L)y = z$ for the lower triangular matrix $D + L$ (forward substitution)

Relaxation methods:

Use linear combination between new and previous iterate:

$$\mathbf{x}_{k+1} = \omega \underbrace{(G\mathbf{x}_k + \mathbf{c})}_{\mathbf{x}'_{k+1}} + (1 - \omega)\mathbf{x}_k = G_\omega \mathbf{x}_k + \omega \mathbf{c},$$

where $\omega > 0$ is a **damping/relaxation parameter** (sometimes, $\omega > 1$ is used, leading to overrelaxation). Target is to choose ω such that $\rho(G_\omega)$ is minimal.

Def: A fixed point method $\mathbf{x}_{k+1} = G\mathbf{x}_k + \mathbf{c}$ with $G = G(A)$ is called *symmetrizable* if for any symmetric positive definite (spd) matrix A , the matrix $I - G$ is similar to an spd matrix, i.e., there is a regular W such that $W(I - G)W^{-1}$ spd.

Symmetrizable: need that $I-G$ is similar to an spd matrix:

↳ Richardson: $G = I - A$

$$\Rightarrow I - G = A \quad \begin{array}{l} \text{is spd if} \\ A \text{ is spd} \end{array}$$

↳ Jacobi: $G = I - D^{-1}A$, set $\underline{W} = D^{1/2}$
($W(I-G)W^{-1}$ spd)

$$\begin{aligned} D^{1/2}(I-G)D^{-1/2} &= I - I + D^{1/2}D^{-1}AD^{-1/2} \\ &= D^{-1/2}AD^{-1/2} \end{aligned}$$

which is spd because D has non-zero diagonal entries if A spd

We have the following result for sym. schemes:

$$x_{n+1} = Gx_n + c, \quad G = G(A), \quad A \text{ spd}$$

$$\Rightarrow \sigma(G) \subset (-\infty, 1) \quad (p(G) = \det \Sigma(\lambda I))$$

$\lambda \in \sigma(G) \Leftrightarrow \lambda$ eigenvalue of A

Because $I-G$ is similar to spd matrix,
all eigenvalues are real and positive

\Rightarrow eigenvalues of G are

$$1 - \lambda_1(G) > 0$$

$$1 > \lambda_2(G)$$

Finding the optimal damping parameter: \rightsquigarrow board

θ

Finding the optimal damping parameter: \rightsquigarrow board

We obtain that

$$\bar{\omega} = \frac{2}{2 - \lambda_{\max}(G) - \lambda_{\min}(G)}$$

is the optimal damping parameter for symmetrizable iteration methods that minimizes the spectral radius. The spectral radius is

$$\rho(G_{\bar{\omega}}) < 1$$

This means that for a suitable choice $\bar{\omega}$ we can make *any* symmetrizable iteration method convergent for an spd A !

Finding the optimal damping parameter: \rightsquigarrow board

We obtain that

$$\bar{\omega} = \frac{2}{2 - \lambda_{\max}(G) - \lambda_{\min}(G)}$$

is the optimal damping parameter for symmetrizable iteration methods that minimizes the spectral radius. The spectral radius is

$$\rho(G_{\bar{\omega}}) < 1$$

This means that for a suitable choice $\bar{\omega}$ we can make *any* symmetrizable iteration method convergent for an spd A !

Optimal damping parameter for Richardson iteration \rightsquigarrow board



Numerical Methods I

MATH-GA 2010.001/CSCI-GA 2420.001

Benjamin Peherstorfer
Courant Institute, NYU

Based on slides by G. Stadler and A. Donev

Today

Last time

- ▶ Iterative methods for systems of linear equations

Today

- ▶ Iterative methods for systems of linear equations
- ▶ Conjugate gradient method

Announcements

- ▶ Homework 4 is due Mon, Nov 4, 2024 before class

Recap: Iterative methods

Let Q be invertible, then

$$\begin{aligned} \mathbf{Ax} = \mathbf{b} &\Leftrightarrow Q^{-1}(\mathbf{b} - \mathbf{Ax}) = \mathbf{0} \\ &\Leftrightarrow (\mathbf{I} - Q^{-1}\mathbf{A})\mathbf{x} + Q^{-1}\mathbf{b} = \mathbf{x} \\ &\Leftrightarrow \mathbf{Gx} + \mathbf{c} = \mathbf{x} \end{aligned}$$

Leads to fixed-point iteration with $\mathbf{G} = (\mathbf{I} - Q^{-1}\mathbf{A})$ and $\mathbf{c} = Q^{-1}\mathbf{b}$

$$\mathbf{x}_{k+1} = \mathbf{Gx}_k + \mathbf{c}$$

and $\mathbf{x} = \mathbf{A}^{-1}\mathbf{b}$ is stationary point

$$\mathbf{Gx} + \mathbf{c} = (\mathbf{I} - Q^{-1}\mathbf{A})\mathbf{x} + Q^{-1}\mathbf{b} = \mathbf{x} - Q^{-1}\mathbf{b} + Q^{-1}\mathbf{b} = \mathbf{x}$$

Recap: Common choices for Q

Requirements on a good Q are

Recap: Common choices for Q

Requirements on a good Q are

- ▶ Q^{-1} is representative of A^{-1}
- ▶ We can numerically solve $Qy = z$ much quicker than $Ax = b$ because in each step we solve

$$Q\mathbf{x}_{k+1} = (Q - A)\mathbf{x}_k + b$$

Recap: Common choices for Q

Requirements on a good Q are

- ▶ Q^{-1} is representative of A^{-1}
- ▶ We can numerically solve $Qy = z$ much quicker than $Ax = b$ because in each step we solve

$$Q\mathbf{x}_{k+1} = (Q - A)\mathbf{x}_k + b$$

Split the matrix A as follows

$$A = L + D + U$$

Recap: Jacobi method

Theorem: Select now $Q = D \dots$ **Jacobi method**. The Jacobi method converges for any starting point \mathbf{x}_0 to the solution of $A\mathbf{x} = \mathbf{b}$ if A is strictly diagonal dominant, i.e.,

$$|a_{ii}| > \sum_{j \neq i} |a_{ij}|, \quad \text{for } i = 1, \dots, n.$$

Recap: Jacobi method

Theorem: Select now $Q = D$... **Jacobi method**. The Jacobi method converges for any starting point \mathbf{x}_0 to the solution of $A\mathbf{x} = \mathbf{b}$ if A is strictly diagonal dominant, i.e.,

$$|a_{ii}| > \sum_{j \neq i} |a_{ij}|, \quad \text{for } i = 1, \dots, n.$$

Notice that $Q = D$ is a good choice in terms of computational costs because we can very quickly solve $D\mathbf{y} = \mathbf{z}$ for the diagonal matrix D

Recap: Gauss-Seidel method

Theorem: Choose $Q = D + L$... **Gauss-Seidel method**. The Gauss-Seidel method converges for any starting point \mathbf{x}_0 if A is symmetric positive definite (spd).

Recap: Gauss-Seidel method

Theorem: Choose $Q = D + L$... **Gauss-Seidel method**. The Gauss-Seidel method converges for any starting point \mathbf{x}_0 if A is symmetric positive definite (spd).

Notice that $Q = D + L$ is a good choice in terms of computational costs because we can very quickly solve $(D + L)y = z$ for the lower triangular matrix $D + L$ (forward substitution)

Recap: Relaxation methods

Use linear combination between new and previous iterate:

$$\mathbf{x}_{k+1} = \omega \underbrace{(G\mathbf{x}_k + \mathbf{c})}_{\mathbf{x}'_{k+1}} + (1 - \omega)\mathbf{x}_k = G_\omega \mathbf{x}_k + \omega \mathbf{c},$$

where $\omega > 0$ is a **damping/relaxation parameter** (sometimes, $\omega > 1$ is used, leading to overrelaxation). Target is to choose ω such that $\rho(G_\omega)$ is minimal.

Def: A fixed point method $\mathbf{x}_{k+1} = G\mathbf{x}_k + \mathbf{c}$ with $G = G(A)$ is called *symmetrizable* if for any symmetric positive definite (spd) matrix A , the matrix $I - G$ is similar to an spd matrix, i.e., there is a regular W such that $W(I - G)W^{-1}$ spd.

Finding the optimal damping parameter: \rightsquigarrow board

$$x_{t+1} = w(Gx_t + c) + (1-w)x_t$$

$$= G_w x_t + wc, \quad G_w = wG + (1-w)I$$

Goal: find $0 < w \leq 1$ with minimal $\rho(G_w)$

The EV of G_w are

$$\begin{aligned} \lambda_i(G_w) &= w \lambda_i(G) + 1-w \\ &= 1 - w(1 - \lambda_i(G)) \end{aligned}$$

Because $\lambda_{\min}(G) \leq \lambda_{\max}(G) < 1$

$$\lambda_i(G_w) < 1$$

and

$$\rho(G_w) = \max \left\{ |1 - w(1 - \lambda_{\min}(G))|, |1 - w(1 - \lambda_{\max}(G))| \right\}$$

Let's look at $w_{\min} = \frac{1}{1 - \lambda_{\min}(G)} > 0$

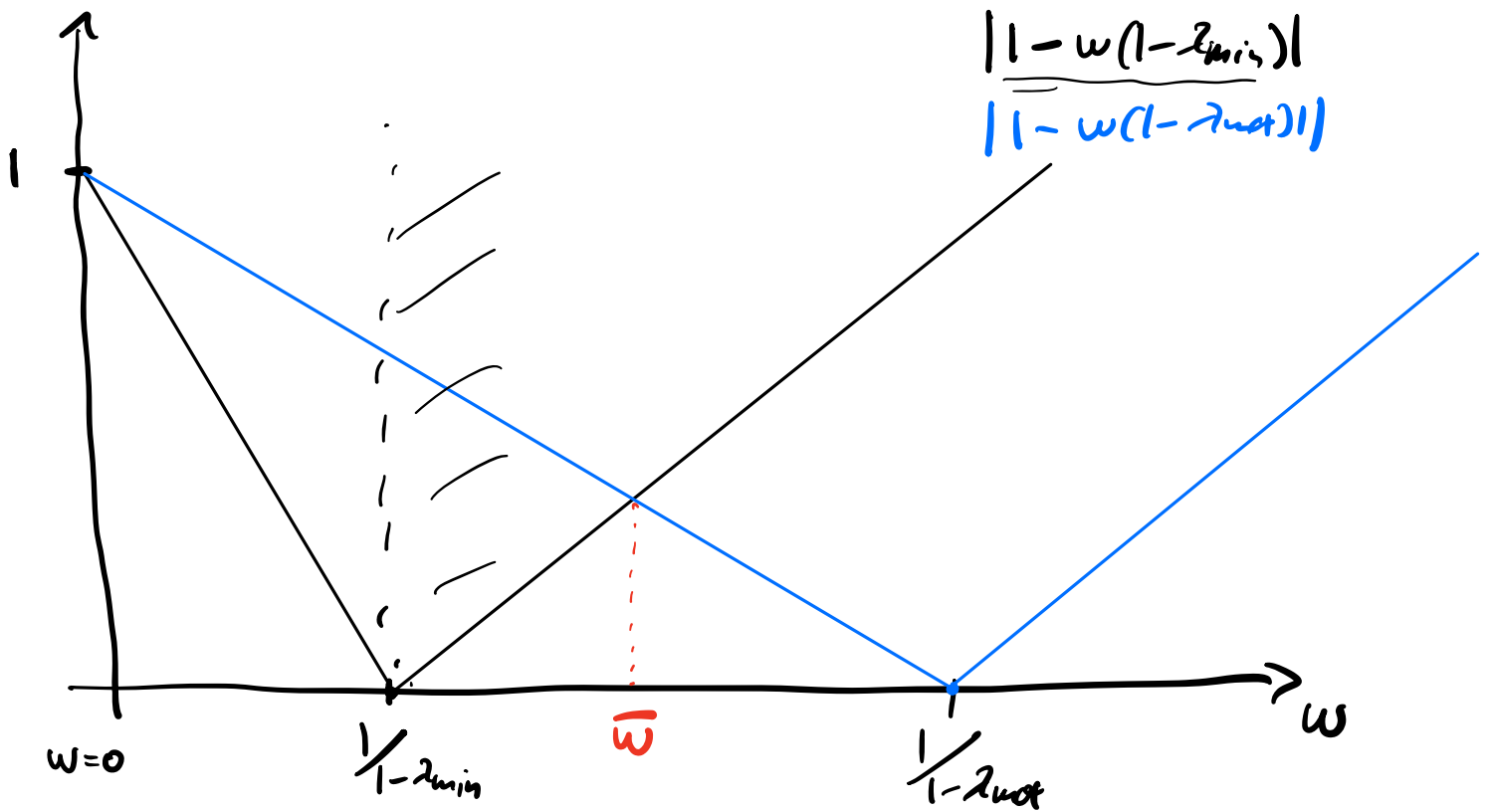
because $0 < 1 - \lambda_{\max}(G) \leq 1 - \lambda_{\min}(G)$

with w_{\min} obtain

$$|1 - w_{\min}(1 - \lambda_{\min}(G))| = 0$$

Similarly, for $w_{\max} = \frac{1}{1 - \lambda_{\max}(G)}$ so that

$$|1 - w_{\max}(1 - \lambda_{\max}(G))| = 0$$



optimal \bar{w} satisfies

$$- [1 - \bar{w}(1 - \lambda_{\min}(G))] = 1 - \bar{w}(1 - \lambda_{\max}(G))$$

solve for \bar{w}

$$\bar{w} = \frac{2}{2 - \lambda_{\max}(G) - \lambda_{\min}(G)}$$

Example: Richardson iteration with SPD matrix

↳ for SPD matrix A have

$$G = I - A$$

$$\lambda_{\min}(G) = 1 - \lambda_{\max}(A)$$

$$\lambda_{\max}(G) = 1 - \lambda_{\min}(A)$$

$$\bar{w} = \frac{2}{\lambda_{\max}(A) + \lambda_{\min}(A)}$$

Spectral radius

$$\begin{aligned}\rho(G_w) &= \rho(I - \bar{w}A) = \\ &= \frac{\lambda_{\max}(A) - \lambda_{\min}(A)}{\lambda_{\max}(A) + \lambda_{\min}(A)} \\ &= \frac{\frac{1}{\lambda_{\min}}}{\frac{1}{\lambda_{\min}}} \frac{\lambda_{\max} - \lambda_{\min}}{\lambda_{\max} + \lambda_{\min}} = \frac{\kappa_2(A) - 1}{\kappa_2(A) + 1} < 1\end{aligned}$$

\Rightarrow notice that the condition number of A enters

\Rightarrow different from direct methods

\Rightarrow if $\kappa_2(A)$ is high, then

$$\frac{\kappa_2(A) - 1}{\kappa_2(A) + 1} \approx 1$$

\Rightarrow very slow convergence

Finding the optimal damping parameter: \rightsquigarrow board

We obtain that

$$\bar{\omega} = \frac{2}{2 - \lambda_{\max}(G) - \lambda_{\min}(G)}$$

is the optimal damping parameter for symmetrizable iteration methods that minimizes the spectral radius. The spectral radius is

$$\rho(G_{\bar{\omega}}) < 1$$

This means that for a suitable choice $\bar{\omega}$ we can make *any* symmetrizable iteration method convergent for an spd A !

Finding the optimal damping parameter: \rightsquigarrow board

We obtain that

$$\bar{\omega} = \frac{2}{2 - \lambda_{\max}(G) - \lambda_{\min}(G)}$$

is the optimal damping parameter for symmetrizable iteration methods that minimizes the spectral radius. The spectral radius is

$$\rho(G_{\bar{\omega}}) < 1$$

This means that for a suitable choice $\bar{\omega}$ we can make *any* symmetrizable iteration method convergent for an spd A !

Optimal damping parameter for Richardson iteration \rightsquigarrow board

Algorithmic perspective on iterative methods

Another interpretation of Richardson iterations?

$$\mathbf{x}_{k+1} = (I - A)\mathbf{x}_k + b$$

Algorithmic perspective on iterative methods

Another interpretation of Richardson iterations?

$$\mathbf{x}_{k+1} = (I - A)\mathbf{x}_k + b = \mathbf{x}_k + (b - A\mathbf{x}_k)$$

Algorithmic perspective on iterative methods

Another interpretation of Richardson iterations?

$$\mathbf{x}_{k+1} = (I - A)\mathbf{x}_k + \mathbf{b} = \mathbf{x}_k + (\mathbf{b} - A\mathbf{x}_k)$$

for $k = 0, 1, \dots$

for $i = 0, 1, \dots, n-1$: $x_{k+1}[i] = x_k[i] + r_k[i]$

where the residual \mathbf{r}_k at iteration k is given by

$$\mathbf{r}_k = \mathbf{b} - A\mathbf{x}_k$$

What would we like to update \mathbf{x}_k with?

Algorithmic perspective on iterative methods

Another interpretation of Richardson iterations?

$$\mathbf{x}_{k+1} = (I - A)\mathbf{x}_k + \mathbf{b} = \mathbf{x}_k + (\mathbf{b} - A\mathbf{x}_k)$$

for $k = 0, 1, \dots$

for $i = 0, 1, \dots, n-1$: $x_{k+1}[i] = x_k[i] + r_k[i]$

where the residual \mathbf{r}_k at iteration k is given by

$$\mathbf{r}_k = \mathbf{b} - A\mathbf{x}_k$$

What would we like to update \mathbf{x}_k with? We would like to update \mathbf{x}_k in the direction of the error $\mathbf{e}_k = \mathbf{x} - \mathbf{x}_k$ because then

$$\mathbf{x}_{k+1} = \mathbf{x}_k + \mathbf{e}_k = \mathbf{x}$$

However, we don't have the error \mathbf{e}_k and therefore it is reasonable to use the next best thing which is the residual in many situations

$$\mathbf{r}_k = \mathbf{b} - A\mathbf{x}_k = A\mathbf{x} - A\mathbf{x}_k = A(\mathbf{x} - \mathbf{x}_k) = A\mathbf{e}_k$$

↪ Richardson iteration updates \mathbf{x}_k in the direction of the residual \mathbf{r}_k

Jacobi iterations

$$\mathbf{x}_{k+1} = (I - D^{-1}A)\mathbf{x}_k + D^{-1}b$$

for $k = 0, 1, \dots$

for $i = 0, 1, \dots, n-1$: $y[i] = 1/a_{ii}r_k[i]$

for $i = 0, 1, \dots, n-1$: $x_{k+1}[i] = x_k[i] + y[i]$

- ▶ In every substep i of iteration k , an update $y[i]$ is computed and stored
- ▶ Applied immediately, this would lead to the (momentary) disappearance of the i -th component of the residual r_k
- ▶ Thus, with this current approximation, equation i would be solved exactly—an improvement that would be lost immediately in the following substep for the equation $i + 1$
- ▶ However, the updates of a component are not applied immediately but only at the end of an iteration step (second i -loop)

Gauss-Seidel iteration

$$\mathbf{x}_{k+1} = (I - (L + D)^{-1}A)\mathbf{x}_k + (L + D)^{-1}b$$

for $k = 0, 1, \dots$

$$\begin{aligned} \text{for } i = 0, 1, \dots, n-1: \quad & r_k[i] = b[i] - \sum_{j=1}^{i-1} a_{ij}x_{k+1}[j] - \sum_{j=i}^n a_{ij}x_k[j] \\ & y[i] = 1/a_{ii}r_k[i], \quad x_{k+1}[i] = x_k[i] + y[i] \end{aligned}$$

- ▶ In contrast to Jacobi method, the update is performed immediately
- ▶ Therefore the new modified values for components $1, \dots, i - 1$ are already available for updating component i

Damping (mostly Jacobi) or over-relaxation (mostly Gauss-Seidel) means to take step lengths different from 1 in the direction of the residual

The spectral radius of typical iterative matrices

- ▶ The spectral radius ρ determines convergence and speed; the smaller ρ , the faster the error decays. In practice, ρ is often very close to 1 so that even though $\rho < 1$ it takes an unreasonable amount of iterations to get a reasonable answer
- ▶ An important sample scenario is the discretization of PDEs: It is typical that ρ depends on the dimension n of the matrix A , and thus in terms of PDE discretization it depends on the mesh width h of the underlying grid. For example

$$\rho \in \mathcal{O}(1 - h_l^2) = \mathcal{O}(1 - \frac{1}{4^l})$$

with mesh width $h_l = 2^{-l}$ in one dimension ($\mathcal{O}(1 - 16^{-l})$ in two spatial dimensions)

- ▶ This is a huge disadvantage: **Why?**

The spectral radius of typical iterative matrices

- ▶ The spectral radius ρ determines convergence and speed; the smaller ρ , the faster the error decays. In practice, ρ is often very close to 1 so that even though $\rho < 1$ it takes an unreasonable amount of iterations to get a reasonable answer
- ▶ An important sample scenario is the discretization of PDEs: It is typical that ρ depends on the dimension n of the matrix A , and thus in terms of PDE discretization it depends on the mesh width h of the underlying grid. For example

$$\rho \in \mathcal{O}(1 - h_l^2) = \mathcal{O}(1 - \frac{1}{4^l})$$

with mesh width $h_l = 2^{-l}$ in one dimension ($\mathcal{O}(1 - 16^{-l})$ in two spatial dimensions)

- ▶ This is a huge disadvantage: **Why?** the finer the grid is (and therefore the more accurate our approximation of the PDE solution should be), the slower these iterative methods get.

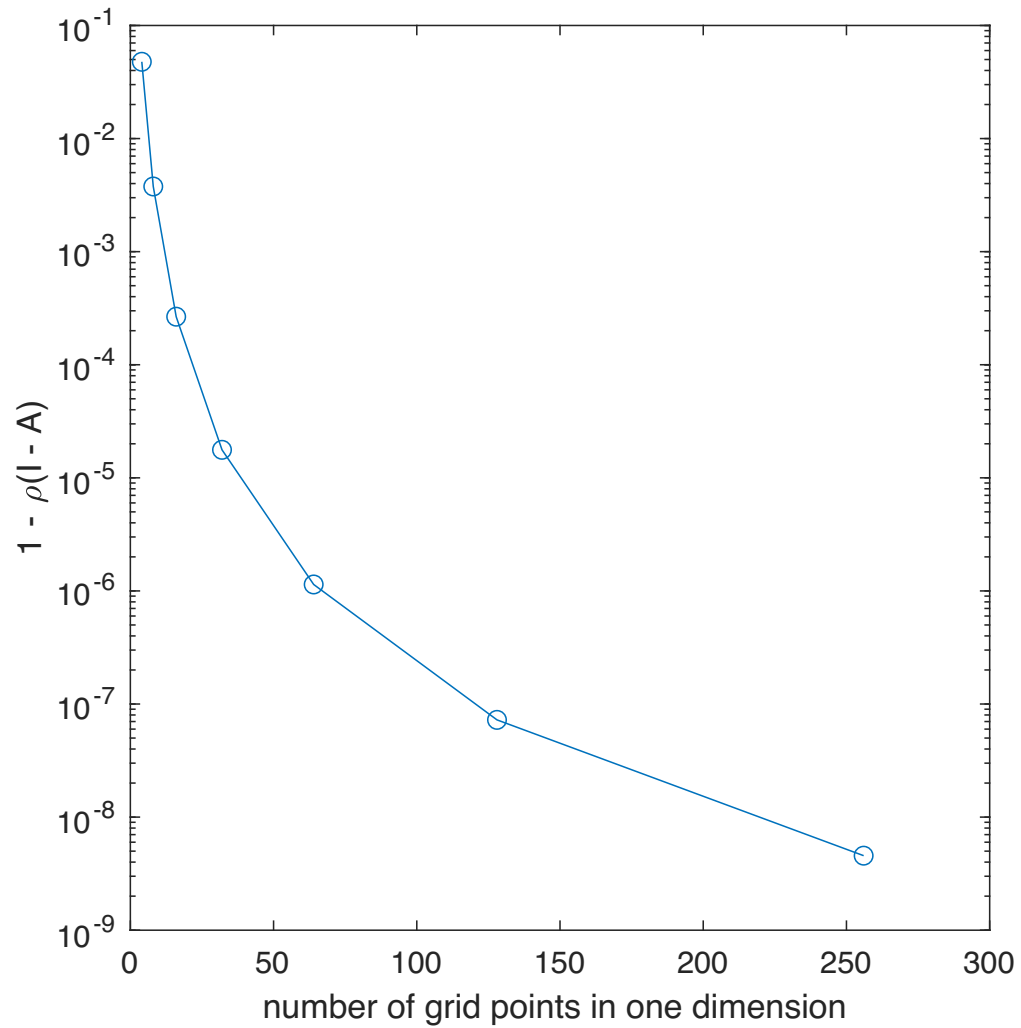
Consider the Laplace equation

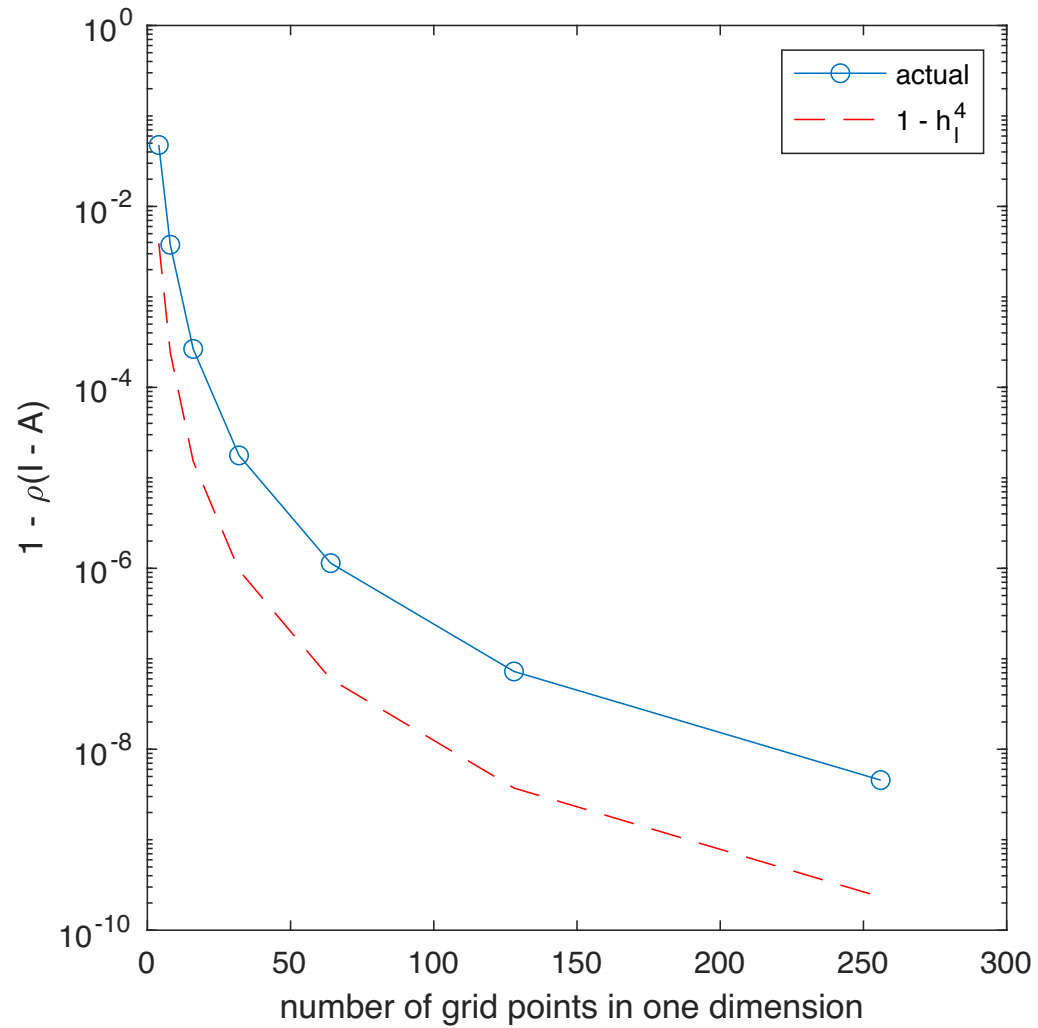
$$\Delta u(x_1, x_2) = 0,$$

in two dimensions and discretize with five-point second-order finite-difference stencil on $N \times N$ grid points (\rightsquigarrow NM2)

Plot the spectral radius $\rho(I - A)$ (Richardson interpolation) w.r.t. number of grid points

```
1: Nlist = 2.^(2:8);
2: eList = [];
3: for N=Nlist
4:     X = gallery('poisson', N)*(1/N^2);
5:     eList(end + 1) = eigs(speye(N^2) - X, 1);
6: end
```



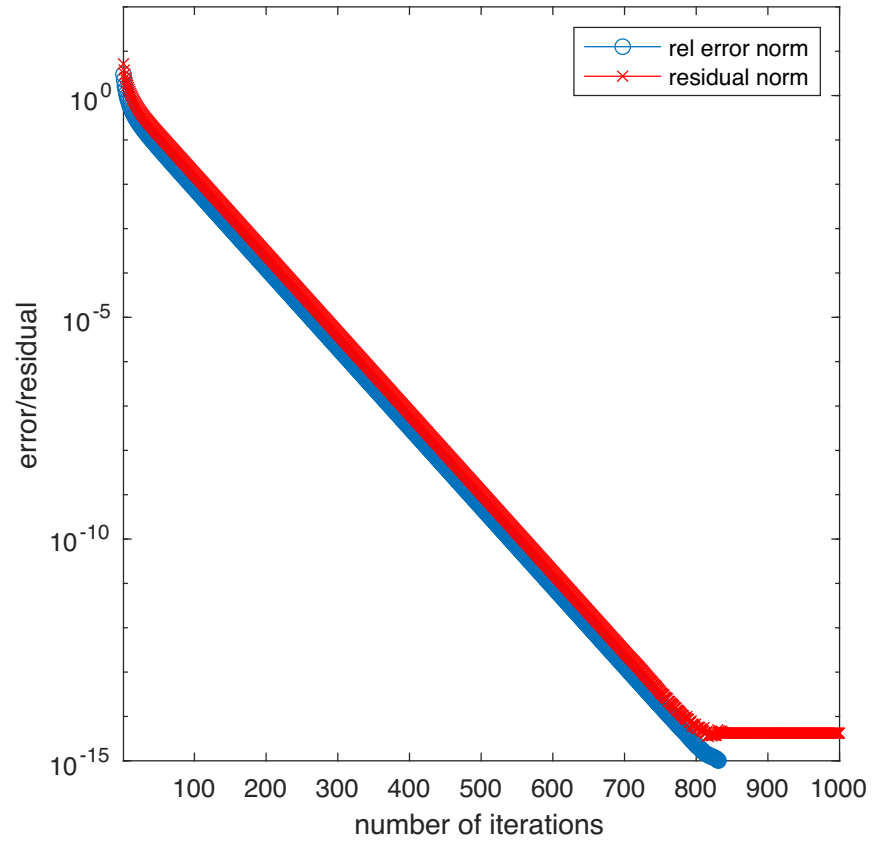


Running Jacobi on a Poisson matrix with $N = 10$ and $N = 100$ grid points in each dimension:

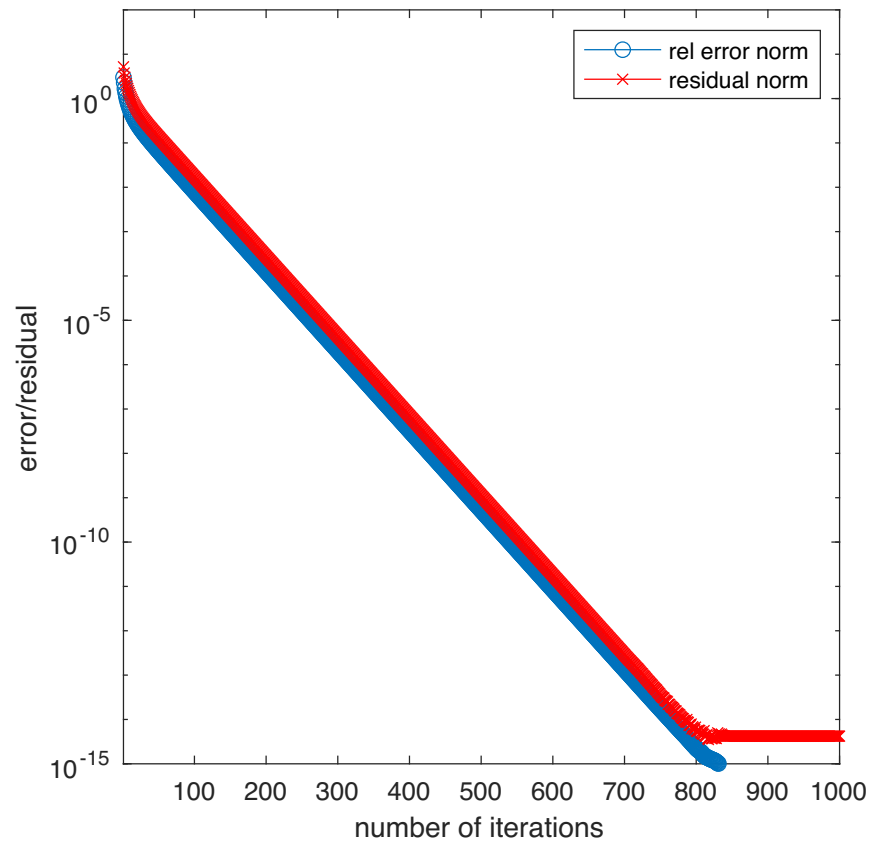
```
1: A = gallery('poisson', N)*(1/N^2);
2: b = randn(N^2, 1);
3: x = randn(N^2, 1);
4: xTrue = A\b;
5: M = speye(N^2) - spdiags(1./diag(A), 0, N^2, N^2)*A;
6: c = spdiags(1./diag(A), 0, N^2, N^2)*b;
7:
8: hist = [];
9: for iter=1:1000
10:     x = M*x + c;
11:     hist(iter, 1) = norm(x - xTrue)/norm(x);
12:     hist(iter, 2) = norm(A*x - b);
13: end
```

$N = 10$

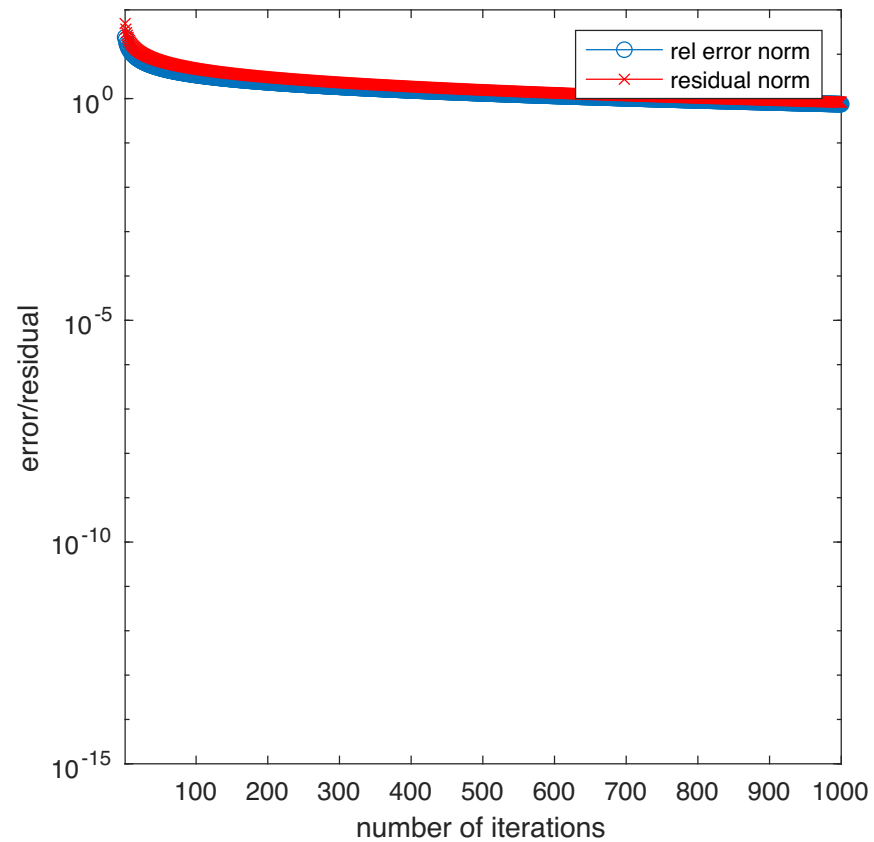
$N = 100$



$N = 10$

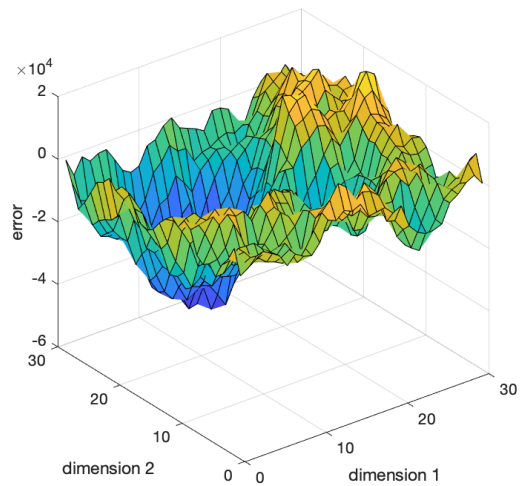


$N = 100$

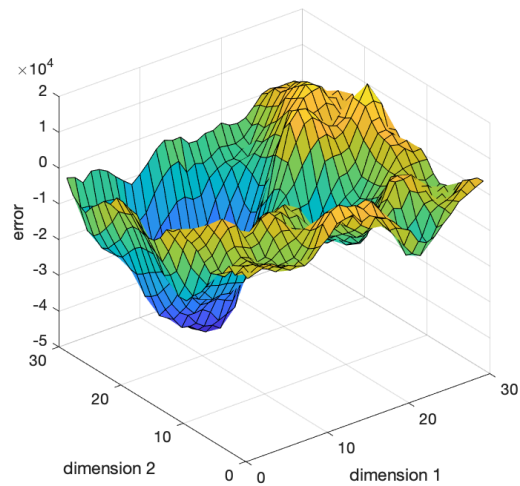


First few iterations of Jacobi relaxation

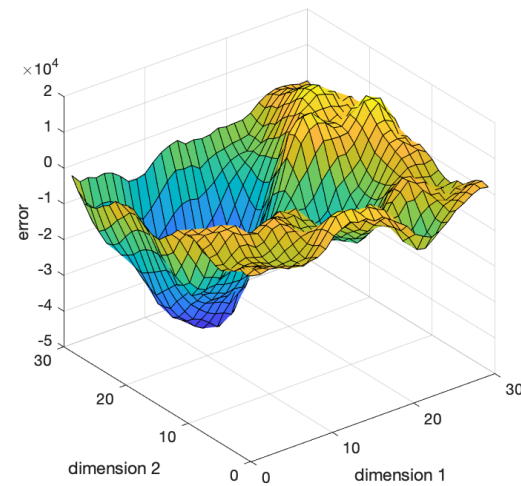
```
1:
2: N = 30; A = gallery('poisson', N)*(1/N^2);
3: [X, Y] = meshgrid(linspace(1, N, N), linspace(1, N, N));
4: b = 10*randn(N^2, 1);
5: x = randn(N^2, 1);
6: xTrue = A\b;
7:
8: for i=1:5
9:     x = (speye(N^2) - spdiags(1./diag(A), 0, N^2, N^2)*A)*x + ...
10:         spdiags(1./diag(A), 0, N^2, N^2)*b;
11: end
```



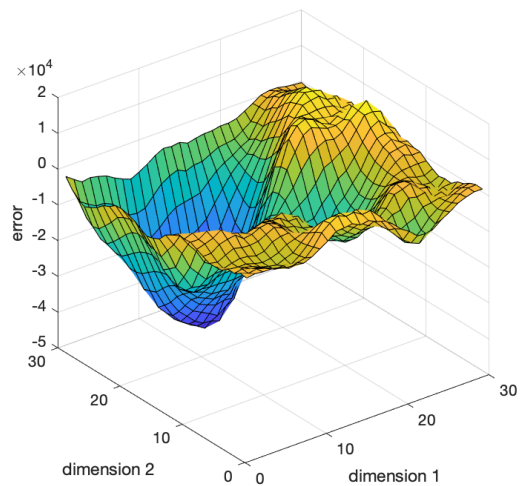
iteration 0



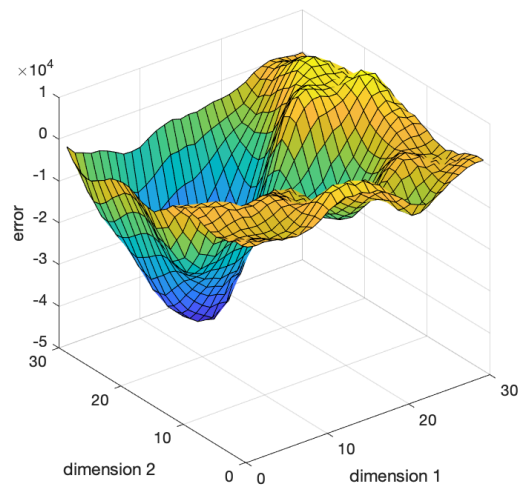
iteration 1



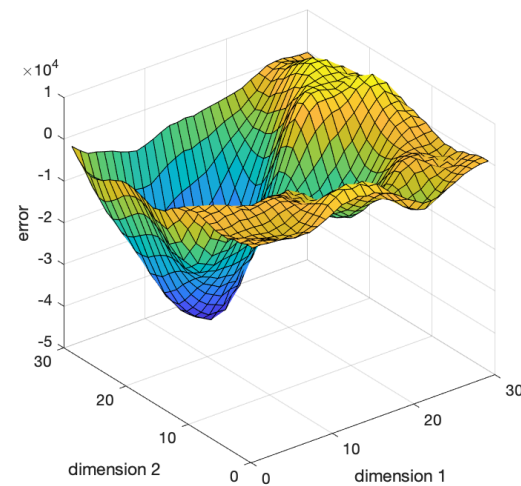
iteration 2



iteration 3

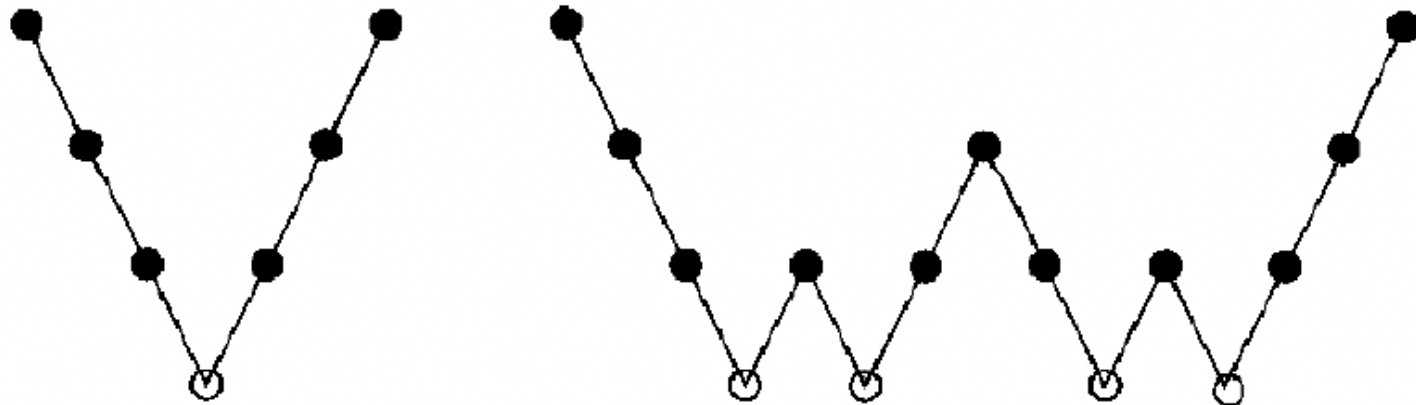


iteration 4



iteration 5

- ▶ Relaxation methods such as Jacobi and Gauss Seidel have a smoothing effect on the error
- ▶ Even if $\rho \approx 1$, only a few iterations are necessary to obtain a smooth error; this means that high-frequency error is reduced very quickly whereas the low-frequency error is reduced slowly
- ▶ Multigrid methods exploit this effect and represent the smoothed error on a coarse grid, where it becomes high-frequency error again, which can be smoothed quickly



⇒ Multigrid methods are among the most efficient solvers for PDE problems ⇒

Conjugate gradient method

In the following A is symmetric positive definite.

Formulate solving $Ax = b$ as an optimization problem: Define

$$f(x) = \frac{1}{2}x^T Ax - b^T x,$$

and minimize

$$\min_{x \in \mathbb{R}^n} f(x)$$

Because A is positive definite, the function f is convex. It is sufficient to look at the gradient

$$\nabla f(x) = \frac{1}{2}A^T x + \frac{1}{2}Ax - b = Ax - b = -r(x) = 0 \iff Ax = b$$

What is the benefit of this point of view?

In the following A is symmetric positive definite.

Formulate solving $Ax = b$ as an optimization problem: Define

$$f(x) = \frac{1}{2}x^T Ax - b^T x,$$

and minimize

$$\min_{x \in \mathbb{R}^n} f(x)$$

Because A is positive definite, the function f is convex. It is sufficient to look at the gradient

$$\nabla f(x) = \frac{1}{2}A^T x + \frac{1}{2}Ax - b = Ax - b = -r(x) = 0 \iff Ax = b$$

What is the benefit of this point of view? We now can let loose all what we know about optimization to solve $Ax = b$

Our first try is applying the method of steepest descent in the direction of the negative gradient

$$-\nabla f = r$$

which happens to be the residual

$$r_k = b - Ax_k$$

$$\alpha_k = \frac{r_k^T r_k}{r_k^T A r_k}$$

$$x_{k+1} = x_k + \alpha_k r_k$$

The step length α_k minimizes $f(x_k + \alpha_k r_k)$ as a function of $\alpha_k \rightsquigarrow$ board

(This is the same as Richardson iterations with damping α_k ; what is the difference between α_k and $\bar{\omega}$? \rightsquigarrow
)

Step length

$$f(x) = \frac{1}{2} x^T A x - b^T x, \quad \underline{\nabla f(x) = Ax - b}$$

$$\min_{\alpha_k} f(x_k + \alpha_k r_k) = \frac{1}{2} (x_k + \alpha_k r_k)^T A (x_k + \alpha_k r_k) - b^T (x_k + \alpha_k r_k)$$

$$\begin{aligned} \frac{\partial}{\partial \alpha_k} f(x_k + \alpha_k r_k) &= r_k^T [\nabla f(x_k + \alpha_k r_k)] \\ &= r_k^T [Ax_k + A\alpha_k r_k - b] \stackrel{!}{=} 0 \end{aligned}$$

$$r_k^T Ax_k + \alpha_k r_k^T A r_k - r_k^T b = 0$$

$$r_k^T \underbrace{(Ax_k - b)}_{r_k} + \alpha_k r_k^T A r_k = 0$$

$$\Rightarrow \alpha_k^* = \frac{r_k^T r_k}{r_k^T A r_k}$$

Our first try is applying the method of steepest descent in the direction of the negative gradient

$$-\nabla f = r$$

which happens to be the residual

$$r_k = b - Ax_k$$

$$\alpha_k = \frac{r_k^T r_k}{r_k^T A r_k}$$

$$x_{k+1} = x_k + \alpha_k r_k$$

The step length α_k minimizes $f(x_k + \alpha_k r_k)$ as a function of $\alpha_k \rightsquigarrow$ board

(This is the same as Richardson iterations with damping α_k ; **what is the difference between α_k and $\bar{\omega}$?** \rightsquigarrow α_k is computable with just mat-vec with A , whereas $\bar{\omega}$ depends on eigenvalues of A)

For steepest descent, if A is spd, we obtain

$$\|x^* - x_k\|_A \leq \left(\frac{\kappa_2(A) - 1}{\kappa_2(A) + 1} \right)^k \|x^* - x_0\|_A,$$

where $\langle x, y \rangle_A = x^T A y$ and $\|\cdot\|_A = \sqrt{\langle \cdot, \cdot \rangle_A}$.

Proof \rightsquigarrow board

$$\|x_k - x^*\|_A \leq \left(\frac{\rho_2(A) - 1}{\rho_2(A) + 1} \right)^k \|x^0 - x^*\|_A$$

Consider

$$\begin{aligned} x_{k+1}(\alpha) &= x_k + \alpha r_k & r_k &= b - Ax_k \\ & & &= Ax^* - Ax_k \\ &= x_k + \alpha A(x^* - x_k) \end{aligned}$$

$$\begin{aligned} \underbrace{x^* - x_{k+1}(\alpha)} &= x^* - x_k - \alpha A(x^* - x_k) \\ &= (1 - \alpha A) \underbrace{(x^* - x_k)} \end{aligned}$$

$$e_{k+1}(\alpha) = (1 - \alpha A) e_k$$

$$\|e_{k+1}(\alpha)\|_A^2 = e_k^T (1 - \alpha A)^T A (1 - \alpha A) e_k$$

Because spd matrix, have orthogonal eigensystem v.r.f. λ_1, λ_2

$$\{z_1, \dots, z_N\} \text{ and } \{\lambda_1, \dots, \lambda_N\}$$

so that we find $\{\alpha_j\}_{j=1}^N \in \mathbb{R}$

$$e_k = \sum_{j=1}^N \alpha_j z_j$$

$$\|e_k\|_A^2 = \dots = \sum_{j=1}^N \alpha_j^2 \lambda_j$$

$$\|e_{k+1}(\alpha)\|_A^2 = \dots = \sum_{j=1}^N \lambda_j \alpha_j^2 (1 - \alpha \lambda_j)^2$$

Numerical Methods I

MATH-GA 2010.001/CSCI-GA 2420.001

Benjamin Peherstorfer
Courant Institute, NYU

Based on slides by G. Stadler and A. Donev

Today

Last time

- ▶ Iterative methods for systems of linear equations

Today

- ▶ Conjugate gradient method

Announcements

- ▶ Homework 4 is due Mon, Nov 4, 2024 before class

Recap

In the following A is symmetric positive definite.

Formulate solving $Ax = b$ as an optimization problem: Define

$$f(x) = \frac{1}{2}x^T Ax - b^T x,$$

and minimize

$$\min_{x \in \mathbb{R}^n} f(x)$$

Because A is positive definite, the function f is convex. It is sufficient to look at the gradient

$$\nabla f(x) = \frac{1}{2}A^T x + \frac{1}{2}Ax - b = Ax - b = -r(x) = 0 \iff Ax = b$$

What is the benefit of this point of view?

Recap

In the following A is symmetric positive definite.

Formulate solving $Ax = b$ as an optimization problem: Define

$$f(x) = \frac{1}{2}x^T Ax - b^T x,$$

and minimize

$$\min_{x \in \mathbb{R}^n} f(x)$$

Because A is positive definite, the function f is convex. It is sufficient to look at the gradient

$$\nabla f(x) = \frac{1}{2}A^T x + \frac{1}{2}Ax - b = Ax - b = -r(x) = 0 \iff Ax = b$$

What is the benefit of this point of view? We now can let loose all what we know about optimization to solve $Ax = b$

Recap

Our first try is applying the method of steepest descent in the direction of the negative gradient

$$-\nabla f = r$$

which happens to be the residual

$$r_k = b - Ax_k$$

$$\alpha_k = \frac{r_k^T r_k}{r_k^T A r_k}$$

$$x_{k+1} = x_k + \alpha_k r_k$$

The step length α_k minimizes $f(x_k + \alpha_k r_k)$ as a function of $\alpha_k \rightsquigarrow$ **last time**

(This is the same as Richardson iterations with damping α_k ; **what is the difference between α_k and $\bar{\omega}$?** \rightsquigarrow
)

Recap

Our first try is applying the method of steepest descent in the direction of the negative gradient

$$-\nabla f = r$$

which happens to be the residual

$$r_k = b - Ax_k$$

$$\alpha_k = \frac{r_k^T r_k}{r_k^T A r_k}$$

$$x_{k+1} = x_k + \alpha_k r_k$$

The step length α_k minimizes $f(x_k + \alpha_k r_k)$ as a function of $\alpha_k \rightsquigarrow$ **last time**

(This is the same as Richardson iterations with damping α_k ; **what is the difference between α_k and $\bar{\omega}$?** \rightsquigarrow α_k is computable with just mat-vec with A , whereas $\bar{\omega}$ depends on eigenvalues of A)

Recap

For steepest descent, if A is spd, we obtain

$$\|x^* - x_k\|_A \leq \left(\frac{\kappa_2(A) - 1}{\kappa_2(A) + 1} \right)^k \|x^* - x_0\|_A,$$

where $\langle x, y \rangle_A = x^T A y$ and $\|\cdot\|_A = \sqrt{\langle \cdot, \cdot \rangle_A}$.

Proof \rightsquigarrow board

Steepest descent chooses α_k such that it minimizes

$$f(x_{k+1}) = f(x_k + \alpha_k r_k)$$

at each step

For the minimizer x^* we have

$$f(x) = \underline{f(x^*)} + \frac{1}{2} \|x^* - x\|_A^2 \quad \leadsto \text{HW4}$$

Thus

$$\min_{\alpha_k} f(x_k + \alpha_k r_k) \Leftrightarrow \min_{\alpha_k} \frac{1}{2} \underbrace{\|x^* - (x_k + \alpha_k r_k)\|_A^2}_{e_{k+1}(\alpha_k)}$$

$$\|e_{k+1}(\alpha_k)\|_A = \min_{\alpha} \|e_{k+1}(\alpha)\|_A$$

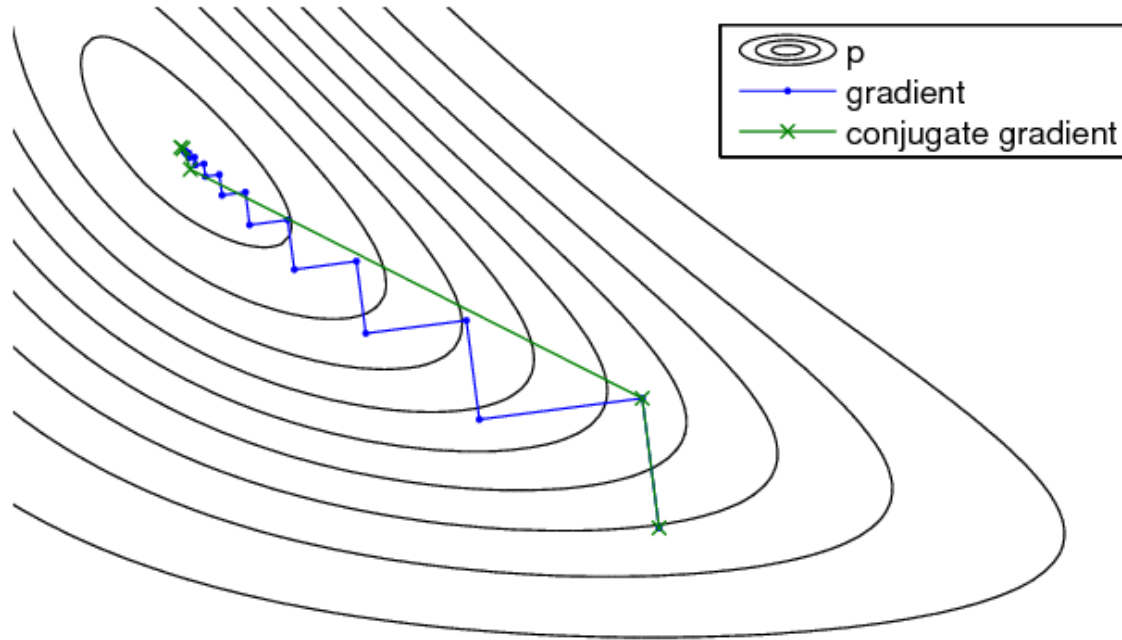
$$\leq \|e_{k+1}(\bar{\alpha})\|_A$$

$$\leq \left(\frac{\lambda_1 - \lambda_N}{\lambda_1 + \lambda_N} \right) \|e_k\|_A$$

By induction

$$\|e_k\|_A \leq \left(\frac{\lambda_1 - \lambda_N}{\lambda_1 + \lambda_N} \right)^k \|e_0\|_A$$

$$\frac{h_2(A) - 1}{h_2(A) + 1} = \frac{\lambda_1/\lambda_N - 1}{\lambda_1/\lambda_N + 1} = \frac{\lambda_1 - \lambda_N}{\lambda_1 + \lambda_N}$$



[Figure: Kuusela et al., 2009]

The convergence behavior of steepest descent in this context can be poor: we eventually get arbitrarily close to the minimum but we can always destroy something of the already achieved when applying the update \rightsquigarrow can we find better search directions?

Conjugate gradient method

- ▶ What do all iterative methods we looked at so far have in common?

Conjugate gradient method

- ▶ What do all iterative methods we looked at so far have in common?
- ▶ All methods so far use information about x_{k-1} to get x_k . All information about earlier iterations is ignored.

Conjugate gradient method

- ▶ What do all iterative methods we looked at so far have in common?
- ▶ All methods so far use information about x_{k-1} to get x_k . All information about earlier iterations is ignored.
- ▶ The conjugate gradient (CG) method is a variation of steepest descent that *has a memory*.

- ▶ Let p_1, \dots, p_k be the directions up to step k , then CG uses the space

$$x_0 + \text{span}\{p_1, \dots, p_k\}, \quad x_0 \text{ starting point}$$

to find the next iterate x_k and thus

$$x_k = x_0 + \sum_{i=1}^k \alpha_i p_i$$

- ▶ (Recall that steepest descent uses only the search direction $p_k = r_{k-1} = -\nabla f(x_{k-1})$ to find the iterate x_k)

We want the following

- a The search directions p_1, \dots, p_k should be linearly independent ("we don't destroy what we have achieved")
- b We have ("we do the best we can at each step")

$$f(x_k) = \min_{x \in x_0 + \text{span}(p_1, \dots, p_k)} f(x)$$

- c The step x_k can be calculated easily from x_{k-1}

What do conditions (a) and (b) guarantee?

We want the following

- a The search directions p_1, \dots, p_k should be linearly independent ("we don't destroy what we have achieved")
- b We have ("we do the best we can at each step")

$$f(x_k) = \min_{x \in x_0 + \text{span}(p_1, \dots, p_k)} f(x)$$

- c The step x_k can be calculated easily from x_{k-1}

What do conditions (a) and (b) guarantee? Convergence in N steps because at the N -th step we have $x_0 + \text{span}(p_1, \dots, p_N) = \mathbb{R}^N$ and thus we minimize f over \mathbb{R}^N

Let's start by writing

$$x_k = x_0 + P_{k-1}y + \alpha p_k,$$

where $P_{k-1} = [p_1, \dots, p_{k-1}] \in \mathbb{R}^{N \times (k-1)}$, $y \in \mathbb{R}^{k-1}$, $\alpha \in \mathbb{R}$.

Our aim is to determine y and α . So let's look at minimizing $f(x_k)$ w.r.t. y and α

$$f(x_k) = \dots = \underbrace{f(x_0 + P_{k-1}y)}_{\text{only depends on } y \text{ not on } \alpha} + \alpha p_k^T A P_{k-1} y + \underbrace{\frac{\alpha^2}{2} p_k^T A p_k - \alpha p_k^T r_0}_{\text{only depends on } \alpha \text{ not on } y}$$

(recall that $f(x) = \frac{1}{2}x^T A x - b^T x$).

Let's start by writing

$$x_k = x_0 + P_{k-1}y + \alpha p_k,$$

where $P_{k-1} = [p_1, \dots, p_{k-1}] \in \mathbb{R}^{N \times (k-1)}$, $y \in \mathbb{R}^{k-1}$, $\alpha \in \mathbb{R}$.

Our aim is to determine y and α . So let's look at minimizing $f(x_k)$ w.r.t. y and α

$$f(x_k) = \dots = \underbrace{f(x_0 + P_{k-1}y)}_{\text{only depends on } y \text{ not on } \alpha} + \alpha p_k^T A P_{k-1} y + \underbrace{\frac{\alpha^2}{2} p_k^T A p_k - \alpha p_k^T r_0}_{\text{only depends on } \alpha \text{ not on } y}$$

(recall that $f(x) = \frac{1}{2}x^T A x - b^T x$).

The mixed term in the middle depends on α and y , otherwise we could optimize separately for y and α . **How should we choose p_k ?**

Let's start by writing

$$x_k = x_0 + P_{k-1}y + \alpha p_k,$$

where $P_{k-1} = [p_1, \dots, p_{k-1}] \in \mathbb{R}^{N \times (k-1)}$, $y \in \mathbb{R}^{k-1}$, $\alpha \in \mathbb{R}$.

Our aim is to determine y and α . So let's look at minimizing $f(x_k)$ w.r.t. y and α

$$f(x_k) = \dots = \underbrace{f(x_0 + P_{k-1}y)}_{\text{only depends on } y \text{ not on } \alpha} + \alpha p_k^T A P_{k-1} y + \underbrace{\frac{\alpha^2}{2} p_k^T A p_k - \alpha p_k^T r_0}_{\text{only depends on } \alpha \text{ not on } y}$$

(recall that $f(x) = \frac{1}{2}x^T A x - b^T x$).

The mixed term in the middle depends on α and y , otherwise we could optimize separately for y and α . **How should we choose p_k ?**

Let's choose the search direction p_k such that

$$p_k^T A P_{k-1} = 0$$

which means

$$p_k \in \text{span}\{A p_1, \dots, A p_{k-1}\}^\perp$$

Thus, with $p_k^T A P_{k-1} = 0$ we get

$$\min_{x_k \in x_0 + \text{span}\{p_1, \dots, p_k\}} f(x_k) = \min_{y \in \mathbb{R}^{k-1}} f(x_0 + P_{k-1}y) + \min_{\alpha \in \mathbb{R}} \left(\frac{\alpha^2}{2} p_k^T A p_k - \alpha p_k^T r_0 \right)$$

- ▶ The first minimization problem is solved for $y = y_{k-1}$ computed from step $k - 1$ and then $x_{k-1} = x_0 + P_{k-1}y_{k-1}$ satisfies

$$f(x_{k-1}) = \min_{x_0 + \text{span}\{p_1, \dots, p_{k-1}\}} f(x)$$

- ▶ The solution to the second minimization problem is just a scalar

$$\alpha_k = \frac{p_k^T r_0}{p_k^T A p_k}$$

↪ satisfy conditions (b) and (c) from above.

- ▶ We said the search directions p_1, \dots, p_k have to be conjugate, i.e., orthogonal w.r.t. A

$$p_i^T A p_j = 0, \quad i, j = 1, \dots, k, i \neq j \quad (1)$$

- ▶ One can show that (1) implies that p_1, \dots, p_k are linearly independent (w.r.t. $\langle \cdot, \cdot \rangle$), which satisfies condition (a)
- ▶ To find the search direction p_k , we want to combine positive aspects of steepest descent and conjugate gradients. In steepest descent we have $p_k = r_{k-1}$. So let's stay close to r_{k-1} but additionally enforce that p_k is A -conjugate to previous search directions p_1, \dots, p_{k-1}

How can we achieve this?

- ▶ We said the search directions p_1, \dots, p_k have to be conjugate, i.e., orthogonal w.r.t. A

$$p_i^T A p_j = 0, \quad i, j = 1, \dots, k, i \neq j \quad (1)$$

- ▶ One can show that (1) implies that p_1, \dots, p_k are linearly independent (w.r.t. $\langle \cdot, \cdot \rangle$), which satisfies condition (a)
- ▶ To find the search direction p_k , we want to combine positive aspects of steepest descent and conjugate gradients. In steepest descent we have $p_k = r_{k-1}$. So let's stay close to r_{k-1} but additionally enforce that p_k is A -conjugate to previous search directions p_1, \dots, p_{k-1}

How can we achieve this? \rightsquigarrow Gram-Schmidt orthogonalization

Apply Gram-Schmidt to r_{k-1} so that we obtain p_k that is A -conjugate to p_1, \dots, p_{k-1}

$$p_k = r_{k-1} - \sum_{j=1}^{k-1} \frac{\langle r_{k-1}, p_j \rangle_A}{\|p_j\|_A^2} p_j$$

We need following technical statements \rightsquigarrow board:

- ▶ If $r_{k-1} = b - Ax_{k-1} \neq 0$, then there exists $p_k \in \text{span}\{Ap_1, \dots, Ap_{k-1}\}^\perp$ such that $p_k^T r_{k-1} \neq 0$ and $p_k^T r_{k-1} = p_k^T r_0$
- ▶ It then follows that (why is this helpful?)

$$\alpha_k = \frac{p_k^T r_0}{p_k^T Ap_k} = \frac{p_k^T r_{k-1}}{p_k^T Ap_k}$$

- ▶ If $r_j \neq 0$ for $j < k$, then

$$\langle r_{k-1}, p_j \rangle_A = 0, \quad j < k - 1.$$

CG technical lemma:

$$r_{k-1} = b - Ax_{k-1} \neq 0$$

$$\Rightarrow \exists p_k \in \text{span} \{A p_1, \dots, A p_{k-1}\}^\perp$$

with $p_k^T r_{k-1} \neq 0$

for $k=1$: Set $p_1 = r_0$

for $k > 1$: $r_{k-1} \neq 0$, then

$$x^* = A^{-1}b \notin x_0 + \text{span} \{p_1, \dots, p_{k-1}\}$$

$$b \notin Ax_0 + \text{span} \{A p_1, \dots, A p_{k-1}\}$$

$$\underbrace{b - Ax_0}_{r_0} \notin \text{span} \{A p_1, \dots, A p_{k-1}\}$$

Thus, there exists $p_k \in \text{span} \{A p_1, \dots, A p_{k-1}\}^\perp$
with $p_k^T r_0 \neq 0$

Select such a $p_k \in \text{span} \{ \dots \}^\perp$ with $p_k^T r_0 \neq 0$, then

$$\begin{aligned} p_k^T r_{k-1} &= p_k^T (b - Ax_{k-1}) \\ &= p_k^T (b - A(x_0 + p_{k-1} \gamma_{k-1})) \\ &= \underbrace{p_k^T (b - Ax_0)}_{p_k^T r_0} - \underbrace{p_k^T A p_{k-1}}_{=0} \gamma_{k-1} \neq 0 \end{aligned}$$

Additionally

$$P_2^T r_{k-1} = P_2^T r_0$$

$$\alpha_k = \frac{P_2^T r_0}{P_2^T A P_2} = \frac{P_2^T r_{k-1}}{P_2^T A P_2}$$

• Now want to convince us that it holds

$$\forall p \in \text{span}\{P_1, \dots, P_2\}: p^T r_k = 0$$

Recall from lecture about lsg problems

$$r = b - Ax \perp \text{columnspace}(A)$$

Condition (b) of CG tells us

$$x_k = \arg \min_{x \in \text{span}\{P_1, \dots, P_2\}} [f(x) = \|Ax - b\|_2^2]$$

$$(*) \quad p^T r_k = 0 \quad \forall p \in \{P_1, \dots, P_2\}$$

• Statement

$$\langle r_{k-1}, P_j \rangle_A = 0 \quad j < k-1$$

Consider

$$x_j = x_{j-1} + \alpha_j P_j$$

$$r_j = b - Ax_j = \underbrace{b - Ax_{j-1}}_{r_{j-1}} - \alpha_j AP_j$$

$$\alpha_j AP_j = r_j - r_{j-1} \quad j < k-1$$

From the Gram-Schmidt process

$$r_j = P_{j+1} + \sum_{i=1}^j \frac{\langle r_j, P_i \rangle_A}{\|P_i\|_A^2} P_i \in \text{span}\{P_1, \dots, P_j\}$$

Further

$$\text{span}\{P_1, \dots, P_{j+1}\} \subseteq \text{span}\{P_1, \dots, P_{k-1}\}$$

$$\alpha_j AP_j = r_{j-1} - r_j \in \text{span}\{P_1, \dots, P_{k-1}\} \quad (**)$$

By def of α_j

$$\begin{aligned} \alpha_j &= \frac{r_{j-1}^T P_j}{\|P_j\|_A^2} = \frac{1}{\|P_j\|_A^2} \left[\|r_{j-1}\|_2^2 - \sum_{i=1}^{j-1} \frac{\langle r_{j-1}, P_i \rangle_A}{\|P_i\|_A^2} \underbrace{r_{j-1}^T P_i}_{=0} \right] \\ &= \frac{\|r_{j-1}\|_2^2}{\|P_j\|_A^2} > 0 \quad \text{because } r_j \neq 0 \quad j < k \end{aligned}$$

$\Rightarrow \alpha_j \neq 0$; with (**) we obtain

$$A p_j \in \text{span}\{p_1, \dots, p_{q-1}\}, \quad j < q-1$$

$$\begin{aligned} \langle r_{q-1}, p_j \rangle_A &= r_{q-1}^T A p_j \\ &= \langle r_{q-1}, \underbrace{A p_j}_{\in \text{span}\{p_1, \dots, p_{q-1}\}} \rangle_{\mathbb{R}^n} \end{aligned}$$

$$\left[\langle r_{q-1}, p \rangle = 0 \quad \forall p \in \text{span}\{p_1, \dots, p_{q-1}\} \right]$$

$$= \langle r_{q-1}, p_j \rangle_A = 0 \quad j < q-1$$

This is important because

$$p_h = r_{q-1} - \sum_{j=1}^{q-1} \frac{\langle r_{q-1}, p_j \rangle_A}{\|p_j\|_A^2} p_j$$

$$= r_{q-1} - \frac{\langle r_{q-1}, p_{q-1} \rangle_A}{\|p_{q-1}\|_A^2} p_{q-1}$$



Apply Gram-Schmidt to r_{k-1} so that we obtain p_k that is A -conjugate to p_1, \dots, p_{k-1}

$$p_k = r_{k-1} - \sum_{j=1}^{k-1} \frac{\langle r_{k-1}, p_j \rangle_A}{\|p_j\|_A^2} p_j$$

We need following technical statements \rightsquigarrow board:

- ▶ If $r_{k-1} = b - Ax_{k-1} \neq 0$, then there exists $p_k \in \text{span}\{Ap_1, \dots, Ap_{k-1}\}^\perp$ such that $p_k^T r_{k-1} \neq 0$ and $p_k^T r_{k-1} = p_k^T r_0$
- ▶ It then follows that (why is this helpful?)

$$\alpha_k = \frac{p_k^T r_0}{p_k^T Ap_k} = \frac{p_k^T r_{k-1}}{p_k^T Ap_k}$$

- ▶ If $r_j \neq 0$ for $j < k$, then

$$\langle r_{k-1}, p_j \rangle_A = 0, \quad j < k - 1.$$

to obtain that (why is this useful?)

$$p_k = r_{k-1} - \frac{\langle r_{k-1}, p_{k-1} \rangle_A}{\|p_{k-1}\|_A^2} p_{k-1}$$

The conjugate gradient method

Choose $x_0 \in \mathbb{R}^N$ and set $p_0 = 0$. For $k = 1, 2, 3, \dots$, stop if $r_{k-1} = b - Ax_{k-1}$ small

1. Set

$$\beta_{k-1} = \frac{\langle r_{k-1}, p_{k-1} \rangle_A}{\|p_{k-1}\|_A^2}$$

2. Set

$$p_k = r_{k-1} - \beta_{k-1} p_{k-1}$$

3. Set

$$\alpha_k = \frac{r_{k-1}^T p_k}{\|p_k\|_A^2}$$

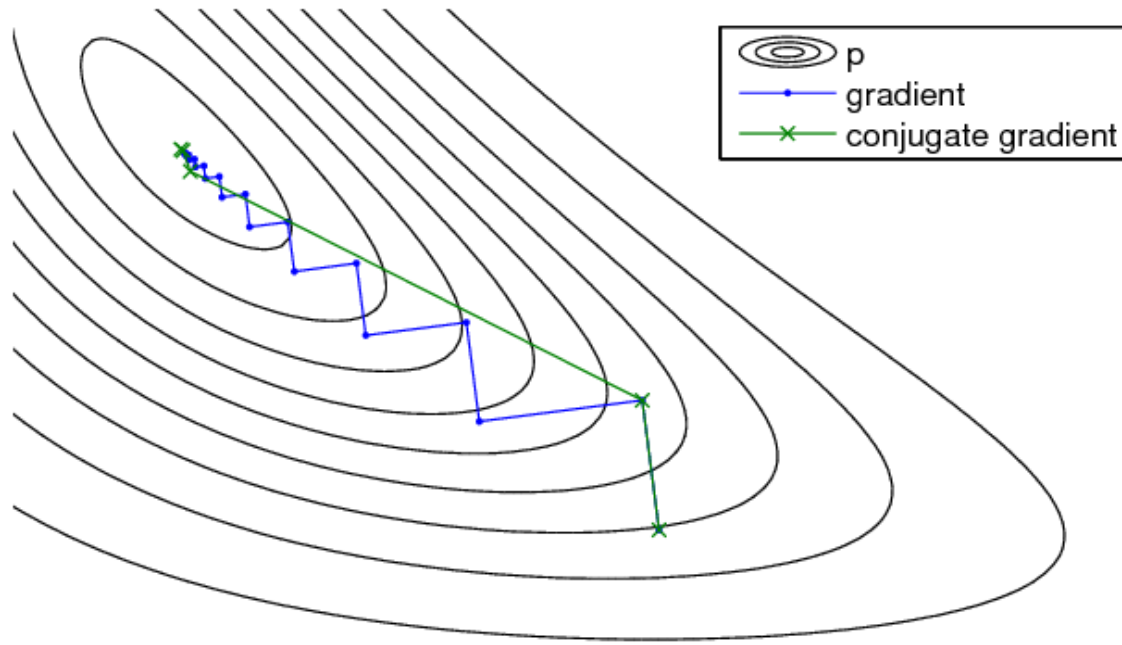
4. Set

$$x_k = x_{k-1} + \alpha_k p_k$$

5. Set

$$r_k = b - Ax_k$$

and check for convergence



[Figure: Kuusela et al., 2009]

It can be shown that for $k \geq 1$ and $e_j \neq 0, j < k$ it holds

$$\|e_k\|_A \leq 2 \left(\frac{\sqrt{\kappa_2(A)} - 1}{\sqrt{\kappa_2(A)} + 1} \right)^k \|e_0\|_A$$

for spd matrices A . \rightsquigarrow Trefethen & Bau

Krylov subspace

Given an spd matrix $A \in \mathbb{R}^{N \times N}$, the Krylov subspace of order k is

$$\mathcal{K}_k(A, r_0) = \text{span} \left\{ r_0, Ar_0, \dots, A^{k-1} r_0 \right\}$$

where, e.g., $r_0 = b - Ax_0$

All search directions of CG are in $\mathcal{K}_k(A, r_0)$ and all iterates x_1, x_2, \dots, x_k are in $x_0 + \mathcal{K}_k(A, r_0)$

There is a range of other methods that apply to more general matrices than spd that build on approximations in Krylov subspaces to accelerate convergence (e.g., GMRES (general residual method))

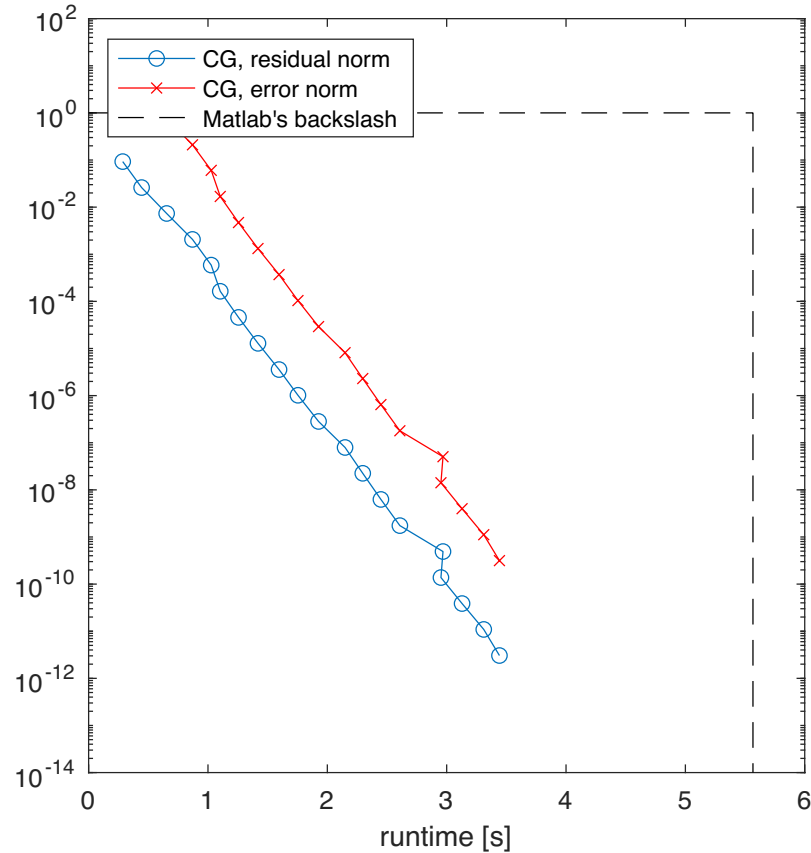
- ▶ There are also methods for finding eigenvalues via Krylov methods (Lanczos, Arnoldi iterations)
- ▶ Think of Krylov methods has having a memory of previous iterations, whereas, e.g., a power method only looks at the previous iteration (if you like stochastic processes, think of Markovian vs. non-Markovian dynamics)

Matlab implementation

```
1: function x = conjgrad(A, b, maxIter)
2:
3: [N, ~] = size(A);
4: x = zeros(N, 1);
5: r = b - A*x;
6: p = r;
7: alpha = (r'*p)/(p'*A*p);
8: x = x + alpha*p;
9: r = b - A*x;
10:
11: for i=1:maxIter
12:     beta = (r'*A*p)/(p'*A*p);
13:     p = r - beta*p;
14:     alpha = (r'*p)/(p'*A*p);
15:     x = x + alpha*p;
16:     r = b - A*x;
17: end
```

Experiment with 10000×10000 spd matrix

Condition number of this matrix is ≈ 5 (very! well conditioned)



Discussion of the CG method

- ▶ In principle, the CG algorithm is a direct solver because it converges after N steps; however, it is mostly used as an iterative method because we don't want to wait for N steps
- ▶ The convergence speed of the CG method depends on matrix properties as well. Fast convergence if the spectrum is clustered.
- ▶ However, similarly slow convergence can be expected for matrices coming from PDE discretizations and therefore preconditioning is necessary

$$Q^{-1}Ax = Q^{-1}b$$

- ▶ Preconditioned CG methods (for example multigrid can act as a preconditioner) are among the fastest solvers and achieve $\mathcal{O}(N)$ in ideal settings.

Numerical Methods I

MATH-GA 2010.001/CSCI-GA 2420.001

Benjamin Peherstorfer
Courant Institute, NYU

Based on slides by G. Stadler and A. Donev

Interpolation

Function approximation

Consider a function $f \in \mathcal{V}$ in a function space \mathcal{V} . Let now ϕ_1, \dots, ϕ_n be a basis of an n -dimensional space \mathcal{V}_n .

Function approximation

Consider a function $f \in \mathcal{V}$ in a function space \mathcal{V} . Let now ϕ_1, \dots, ϕ_n be a basis of an n -dimensional space \mathcal{V}_n .

The task that we are interested in is finding a function $f^* \in \mathcal{V}_n$ that approximates f with coefficients c_1, \dots, c_n :

$$f^*(x) = \sum_{i=1}^n c_i \phi_i(x).$$

Function approximation

Consider a function $f \in \mathcal{V}$ in a function space \mathcal{V} . Let now ϕ_1, \dots, ϕ_n be a basis of an n -dimensional space \mathcal{V}_n .

The task that we are interested in is finding a function $f^* \in \mathcal{V}_n$ that approximates f with coefficients c_1, \dots, c_n :

$$f^*(x) = \sum_{i=1}^n c_i \phi_i(x).$$

If we have an inner product, what is the best approximation?

Function approximation

Consider a function $f \in \mathcal{V}$ in a function space \mathcal{V} . Let now ϕ_1, \dots, ϕ_n be a basis of an n -dimensional space \mathcal{V}_n .

The task that we are interested in is finding a function $f^* \in \mathcal{V}_n$ that approximates f with coefficients c_1, \dots, c_n :

$$f^*(x) = \sum_{i=1}^n c_i \phi_i(x).$$

If we have an inner product, what is the best approximation? The best-approximation of f in \mathcal{V}_n w.r.t. the induced norm is given by the projection

$$f^* = \Pi_n f,$$

where Π_n is the orthogonal projection onto \mathcal{V}_n .

Function approximation

Consider a function $f \in \mathcal{V}$ in a function space \mathcal{V} . Let now ϕ_1, \dots, ϕ_n be a basis of an n -dimensional space \mathcal{V}_n .

The task that we are interested in is finding a function $f^* \in \mathcal{V}_n$ that approximates f with coefficients c_1, \dots, c_n :

$$f^*(x) = \sum_{i=1}^n c_i \phi_i(x).$$

If we have an inner product, what is the best approximation? The best-approximation of f in \mathcal{V}_n w.r.t. the induced norm is given by the projection

$$f^* = \Pi_n f,$$

where Π_n is the orthogonal projection onto \mathcal{V}_n .

How can we compute the projection?

Function approximation

Consider a function $f \in \mathcal{V}$ in a function space \mathcal{V} . Let now ϕ_1, \dots, ϕ_n be a basis of an n -dimensional space \mathcal{V}_n .

The task that we are interested in is finding a function $f^* \in \mathcal{V}_n$ that approximates f with coefficients c_1, \dots, c_n :

$$f^*(x) = \sum_{i=1}^n c_i \phi_i(x).$$

If we have an inner product, what is the best approximation? The best-approximation of f in \mathcal{V}_n w.r.t. the induced norm is given by the projection

$$f^* = \Pi_n f,$$

where Π_n is the orthogonal projection onto \mathcal{V}_n .

How can we compute the projection? In many cases, we cannot directly compute the projection of f onto \mathcal{V}_n because we have “too little knowledge about f ” Instead?

Function approximation

Consider a function $f \in \mathcal{V}$ in a function space \mathcal{V} . Let now ϕ_1, \dots, ϕ_n be a basis of an n -dimensional space \mathcal{V}_n .

The task that we are interested in is finding a function $f^* \in \mathcal{V}_n$ that approximates f with coefficients c_1, \dots, c_n :

$$f^*(x) = \sum_{i=1}^n c_i \phi_i(x).$$

If we have an inner product, what is the best approximation? The best-approximation of f in \mathcal{V}_n w.r.t. the induced norm is given by the projection

$$f^* = \Pi_n f,$$

where Π_n is the orthogonal projection onto \mathcal{V}_n .

How can we compute the projection? In many cases, we cannot directly compute the projection of f onto \mathcal{V}_n because we have “too little knowledge about f ” Instead? \rightsquigarrow interpolation (/regression)

Interpolation

Consider n pairs of data samples $(x_i, y_i), i = 1, \dots, n$ with

$$y_i = f(x_i)$$

Based on $\{(x_i, y_i)\}_{i=1}^n$, we now would like to find an approximation $\tilde{f} \in \mathcal{V}_n$ that is “close” to f .

For example, we could enforce the interpolation condition, namely that it holds

$$\tilde{f}(x_i) = f(x_i), \quad i = 1, \dots, n$$

We could also use regression ($m > n$) and minimize, e.g.,

$$\frac{1}{m} \sum_{i=1}^m |y_i - \tilde{f}(x_i)|^2$$

The error of \tilde{f} w.r.t. f can then typically be split into two components (we will formalize this moving forward): **which?**

The error of \tilde{f} w.r.t. f can then typically be split into two components (we will formalize this moving forward): **which?**

$$\|\tilde{f} - f\| \leq \Lambda(x_1, \dots, x_n) \|f^* - f\|$$

The projection error $\|f^* - f\|$ describes the best we can do in the space \mathcal{V}_n . Even if we had “full knowledge” of f so that we could compute $f^* = \Pi_n f$, we are limited by the expressiveness of the space \mathcal{V}_n

Intuitively, we'd also expect that the error of \tilde{f} depends on the points x_1, \dots, x_n at which we have samples of f . This is captured by the “constant” $\Lambda(x_1, \dots, x_n)$ that is independent of f but depends on x_1, \dots, x_n .

Polynomial interpolation

Consider $n + 1$ pairs (x_i, y_i) , $i = 0, \dots, n$ of a function f with

$$y_i = f(x_i)$$

Polynomial interpolation

Consider $n + 1$ pairs $(x_i, y_i), i = 0, \dots, n$ of a function f with

$$y_i = f(x_i)$$

Let now \mathbb{P}_n be the set of all polynomials up to degree n over \mathbb{R} so that we have for all $P \in \mathbb{P}_n$

$$P(x) = a_n x^n + a_{n-1} x^{n-1} + \dots + a_1 x + a_0, \quad a_n, \dots, a_0 \in \mathbb{R}$$

Polynomial interpolation

Consider $n + 1$ pairs $(x_i, y_i), i = 0, \dots, n$ of a function f with

$$y_i = f(x_i)$$

Let now \mathbb{P}_n be the set of all polynomials up to degree n over \mathbb{R} so that we have for all $P \in \mathbb{P}_n$

$$P(x) = a_n x^n + a_{n-1} x^{n-1} + \dots + a_1 x + a_0, \quad a_n, \dots, a_0 \in \mathbb{R}$$

We would like to find a $P \in \mathbb{P}_n$ such that

$$P(x_i) = y_i, \quad i = 0, \dots, n$$

- ▶ The P is what \tilde{f} was on the previous slide
- ▶ By saying P is a polynomial of degree n , we fixed the space \mathcal{V}_{n+1} with the notation of the previous slide

Theorem: Given $n + 1$ points (x_i, y_i) with pairwise distinct x_0, \dots, x_n , there exists a unique polynomial $P \in \mathbb{P}_n$ such that

$$P(x_i) = y_i, \quad i = 0, \dots, n$$

We sometimes refer to this unique polynomial as $P_f(\cdot | x_0, \dots, x_n) \rightsquigarrow$ **board**

Given $n+1$ nodes (t_i, y_i) with pairwise distinct t_0, \dots, t_n . There exists a unique polynomial $P_n \in \mathcal{P}_n$ such that

$$P(t_i) = y_i, \quad i = 0, \dots, n$$

If $P, Q \in \mathcal{P}_n$ are two interpolatory polynomials with

$$P(t_i) = Q(t_i), \quad i = 0, \dots, n$$

then $P - Q$ is a polynomial of at most degree n and has $n+1$ roots t_0, \dots, t_n

$$(P - Q)(t_i) = 0$$

and thus $P - Q$ has to be the zero polynomial.

Furthermore, the space \mathcal{P}_n and \mathbb{R}^{n+1} have $n+1$ dimensions and the map

$$P \mapsto [P(t_0), \dots, P(t_n)]$$

is linear.

Uniqueness \Rightarrow injective $\xRightarrow{\text{linear over } \mathbb{R}^{n+1} \text{ dim space}}$ surjectivity.

Theorem: Given $n + 1$ points (x_i, y_i) with pairwise distinct x_0, \dots, x_n , there exists a unique polynomial $P \in \mathbb{P}_n$ such that

$$P(x_i) = y_i, \quad i = 0, \dots, n$$

We sometimes refer to this unique polynomial as $P_f(\cdot | x_0, \dots, x_n) \rightsquigarrow$ **board**

Let's try to construct $P_f(\cdot | x_0, \dots, x_n)$. What do we need to construct $P \in \mathbb{P}_n$ for a data set $\{(x_i, y_i)\}_{i=0}^n$?

Theorem: Given $n + 1$ points (x_i, y_i) with pairwise distinct x_0, \dots, x_n , there exists a unique polynomial $P \in \mathbb{P}_n$ such that

$$P(x_i) = y_i, \quad i = 0, \dots, n$$

We sometimes refer to this unique polynomial as $P_f(\cdot | x_0, \dots, x_n) \rightsquigarrow$ **board**

Let's try to construct $P_f(\cdot | x_0, \dots, x_n)$. What do we need to construct $P \in \mathbb{P}_n$ for a data set $\{(x_i, y_i)\}_{i=0}^n$? **A basis of \mathbb{P}_n .**

Theorem: Given $n + 1$ points (x_i, y_i) with pairwise distinct x_0, \dots, x_n , there exists a unique polynomial $P \in \mathbb{P}_n$ such that

$$P(x_i) = y_i, \quad i = 0, \dots, n$$

We sometimes refer to this unique polynomial as $P_f(\cdot | x_0, \dots, x_n) \rightsquigarrow$ **board**

Let's try to construct $P_f(\cdot | x_0, \dots, x_n)$. What do we need to construct $P \in \mathbb{P}_n$ for a data set $\{(x_i, y_i)\}_{i=0}^n$? **A basis of \mathbb{P}_n** . Let's give the monomial basis $1, x, x^2, \dots, x^n$ a chance \rightsquigarrow **board**

poly interpolation: monomials as basis

$$P(x) = a_n x^n + \dots + a_1 x + a_0$$

We have $n+1$ unknowns

$$a_n, \dots, a_0$$

and $n+1$ equations

$$P(x_0) = y_0$$

\vdots

$$P(x_n) = y_n$$

Linear system

$$\begin{bmatrix} x_0^n & x_0^{n-1} & \dots & x_0 & 1 \\ \vdots & \vdots & \dots & \vdots & \vdots \\ x_n^n & x_n^{n-1} & \dots & x_n & 1 \end{bmatrix} \begin{bmatrix} a_n \\ \vdots \\ a_1 \\ a_0 \end{bmatrix} = \begin{bmatrix} y_0 \\ \vdots \\ y_n \end{bmatrix}$$

$$= V_n \quad \begin{matrix} (n+1) \times (n+1) & n+1 & n+1 \end{matrix}$$

It can be shown that

$$\det(V_n) = \prod_{j=0}^n \prod_{i=j+1}^n (x_i - x_j)$$

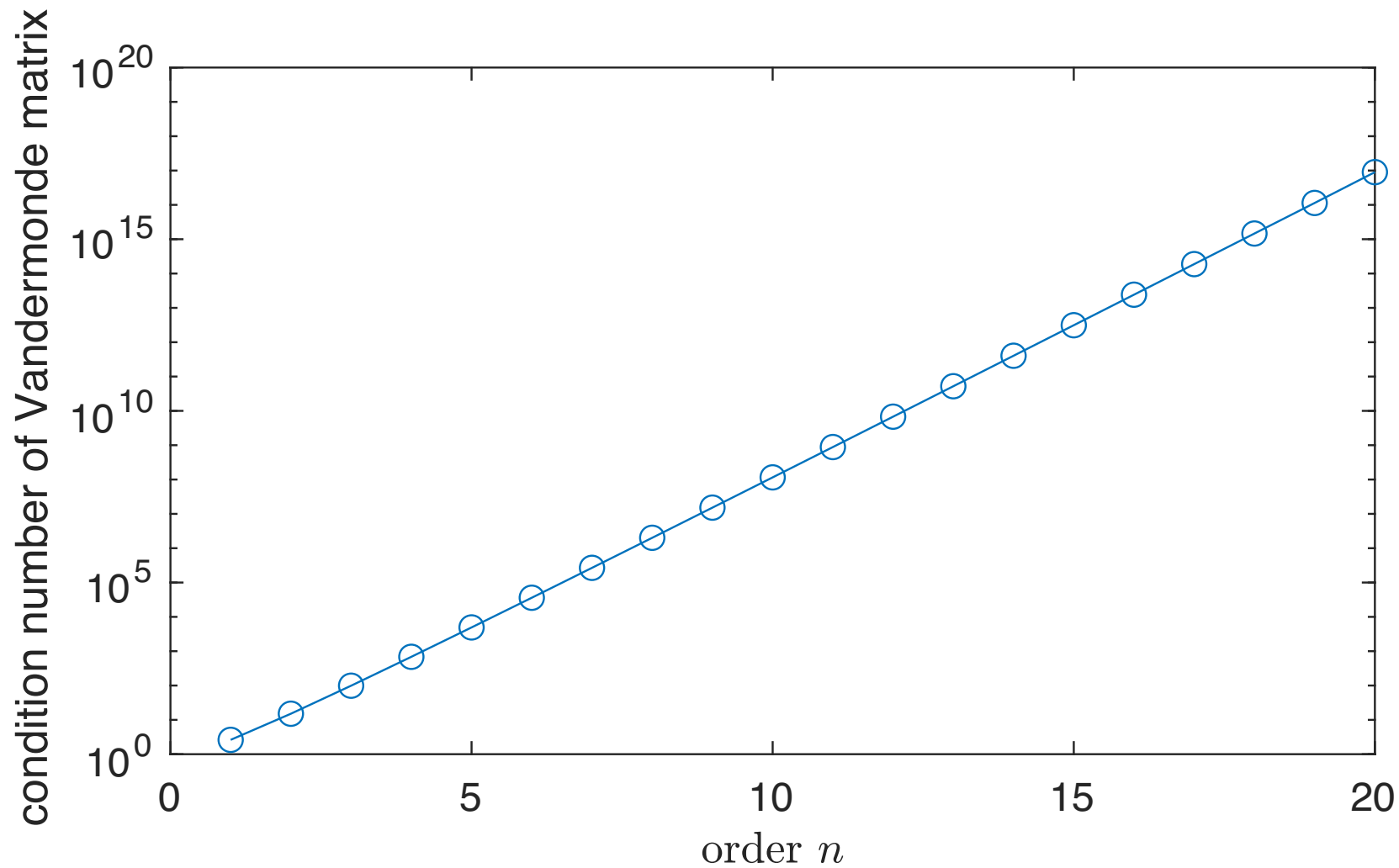
\det is non-zero iff $x_i \neq x_j$ for $i \neq j$

\Rightarrow agrees with assumption that nodes x_0, \dots, x_n distinct

What could go wrong?

↳ condition number of V_n is high

↳ Do we really need to invest $O(n^3)$ costs to compute $P_p(\cdot | x_0, \dots, x_n)$?

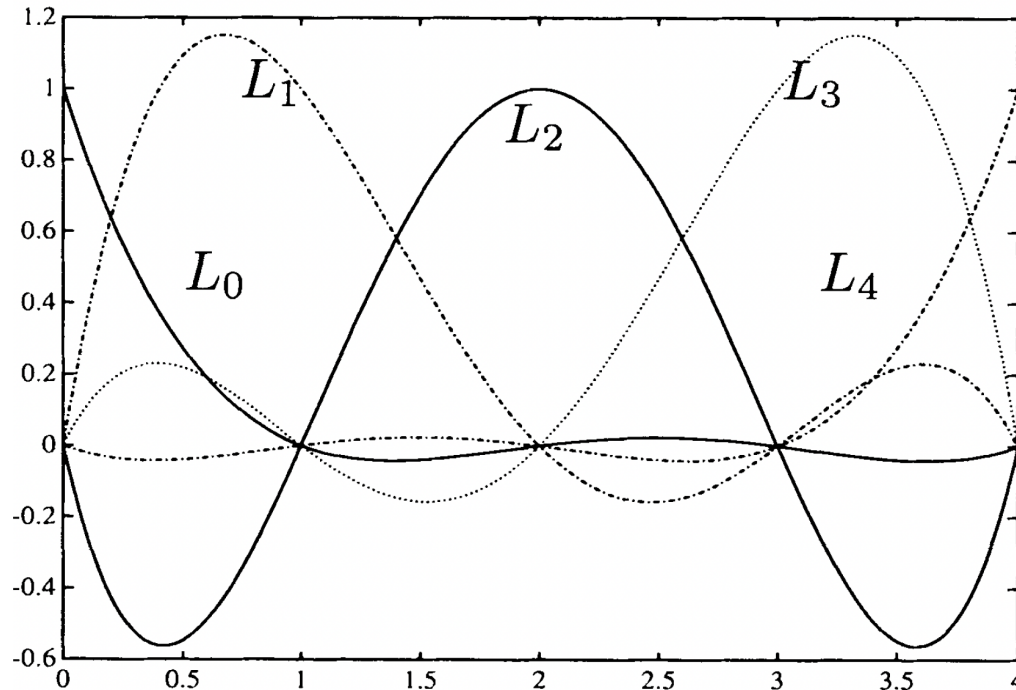


⇒ we really should look for another basis

Lagrange basis

The Lagrange polynomials $L_0, \dots, L_n \in \mathbb{P}_n$ are uniquely defined for distinct x_0, \dots, x_n

$$L_i(x_j) = \delta_{ij}, \quad L_i \in \mathbb{P}_n.$$



Lagrange polynomials up to order $n = 4$ for equidistant x_0, \dots, x_4 . [Figure: Deuffhard]

The corresponding explicit formula is

$$L_i(x) = \prod_{\substack{j=0 \\ j \neq i}}^n \frac{x - x_j}{x_i - x_j}, \quad i = 0, \dots, n$$

What are the coefficients a_n, \dots, a_0 so that $P(x_i) = y_i$ for $i = 0, \dots, n$?

The corresponding explicit formula is

$$L_i(x) = \prod_{\substack{j=0 \\ j \neq i}}^n \frac{x - x_j}{x_i - x_j}, \quad i = 0, \dots, n$$

What are the coefficients a_n, \dots, a_0 so that $P(x_i) = y_i$ for $i = 0, \dots, n$?

$$P(x) = \sum_{i=0}^n y_i L_i(x)$$

because

$$P(x_j) = \sum_{i=0}^n y_i L_i(x_j) = \sum_{i=0}^n y_i \delta_{ij} = y_j$$

If we have the basis L_0, \dots, L_n , we obtain the polynomial P for free. **Drawback?**

The corresponding explicit formula is

$$L_i(x) = \prod_{\substack{j=0 \\ j \neq i}}^n \frac{x - x_j}{x_i - x_j}, \quad i = 0, \dots, n$$

What are the coefficients a_n, \dots, a_0 so that $P(x_i) = y_i$ for $i = 0, \dots, n$?

$$P(x) = \sum_{i=0}^n y_i L_i(x)$$

because

$$P(x_j) = \sum_{i=0}^n y_i L_i(x_j) = \sum_{i=0}^n y_i \delta_{ij} = y_j$$

If we have the basis L_0, \dots, L_n , we obtain the polynomial P for free. **Drawback?** but the cost of evaluating the polynomial is too high for practical computations

The Lagrange polynomials are orthogonal w.r.t. the following inner product over \mathbb{P}_n

$$\langle P, Q \rangle = \sum_{i=0}^n P(x_i)Q(x_i), \quad P, Q \in \mathbb{P}_n$$

Let's try to generalize this to other scalar products to find other orthogonal bases

Orthogonal polynomials

Define an **inner product between functions**:

$$(f, g) = \int_a^b \omega(x) f(x) g(x) dx,$$

where $\omega(x) > 0$ for $a \leq x \leq b$ is a **weight function**. The **induced norm** is $\|f\| := \sqrt{(f, f)}$.

Let $P_0, P_1, P_2, \dots, P_K$ be polynomials of $0, 1, 2, \dots, K$ order, respectively. They are called orthogonal polynomials on $[a, b]$ with respect to the weight function $\omega(x)$ if it holds

$$(P_i, P_j) = \int_a^b \omega(x) P_i(x) P_j(x) dx = \delta_{ij} \gamma_i, \quad i, j = 0, \dots, K,$$

with $\gamma_i = \|P_i\|^2 > 0$.

To define orthogonal polynomials uniquely, we normalize them so that the leading coefficient is one, i.e.,

$$P_k(x) = x^k + \dots$$

Theorem: There exist uniquely determined orthogonal polynomials $P_k \in \mathbb{P}_k$ with leading coefficient 1. These polynomials satisfy the 3-term recurrence relation:

$$P_k(x) = (x + a_k)P_{k-1}(x) + b_k P_{k-2}(x), \quad k = 2, 3, \dots$$

with starting values $P_0 = 1$, $P_1 = x + a_1$, where

$$a_k = -\frac{(xP_{k-1}, P_{k-1})}{(P_{k-1}, P_{k-1})}, \quad b_k = -\frac{(P_{k-1}, P_{k-1})}{(P_{k-2}, P_{k-2})}$$

Proof: \rightsquigarrow board

3-term recurrence; show this by induction

↳ only polynomial of degree 0 with leading coefficient 1 is

$$P_0 = 1 \in P_0$$

↳ $P_1 = x + a_1$; it is unique because a_1 is unique is computed as

$$a_1 = - \frac{(xP_0, P_0)}{(P_0, P_0)}$$

and $(P_0, P_1) = 0$

↳ Suppose P_0, \dots, P_{k-1} have been constructed

↳ P_j has degree j

↳ P_i, P_j are orthogonal if $i \neq j$ w.r.t. (\cdot, \cdot)

↳ leading coefficient is 1

If $P_k \in P_k$ is a polynomial of degree k and normalized (leading coefficient is 1), then

$$\underbrace{P_k}_{x^k + \dots} - x \underbrace{P_{k-1}}_{x(x^{k-1} + \dots)}$$

is of degree $\leq k-1$

the P_0, \dots, P_{r-1} form orthogonal basis of \mathbb{P}_{r-1}
w.r.t. weighted inner product so that

$$P_r - x P_{r-1} = \sum_{j=1}^{r-1} c_j P_j$$

with coefficients

$$c_j = \frac{(P_r - x P_{r-1}, P_j)}{(P_j, P_j)} \left. \vphantom{\frac{(P_r - x P_{r-1}, P_j)}{(P_j, P_j)}} \right\} \begin{array}{l} \text{projection of} \\ P_r - x P_{r-1} \\ \text{onto basis fun } P_j \end{array}$$

If P_r orthogonal to P_0, \dots, P_{r-1} , then

$$(P_r, P_j) = 0, \quad j = 0, \dots, r-1$$

and thus

$$\begin{aligned} c_j &= - \frac{(x P_{r-1}, P_j)}{(P_j, P_j)} = - \frac{\int w(x) (x P_{r-1}(x) P_j(x))}{(P_j, P_j)} \\ &= - \frac{(P_{r-1}, x P_j)}{(P_j, P_j)} \end{aligned}$$

We obtain

$$c_{k-1} = - \frac{(x P_{k-1}, P_{k-1})}{(P_{k-1}, P_{k-1})}$$

$$c_{k-2} = - \frac{(P_{k-1}, x P_{k-2})}{(P_{k-2}, P_{k-2})} \begin{cases} P_{k-1} = (x + a_{k-1}) P_{k-2} + b_k P_{k-3} \\ x P_{k-2} = P_{k-1} - a_{k-1} P_{k-2} - b_k P_{k-3} \end{cases}$$

$$= - \frac{(P_{k-1}, P_{k-1} - a_{k-1} P_{k-2} - b_k P_{k-3})}{(P_{k-2}, P_{k-2})}$$

$$= - \frac{(P_{k-1}, P_{k-1})}{(P_{k-2}, P_{k-2})}$$

$c_0 = \dots = c_{k-3} = 0$ because P_0, \dots, P_{k-1} orthogonal

$$P_k - x P_{k-1} = c_{k-1} P_{k-1} + c_{k-2} P_{k-2}$$

$$P_k = (x + \underbrace{c_{k-1}}_{a_k}) P_{k-1} + \underbrace{c_{k-2}}_{b_k} P_{k-2}$$

Numerics with 3-term recurrence

- ▶ 3-term recurrence can be used to compute polynomials P_k completely, or
- ▶ evaluate P_k at a point x_0 via pre-computing the a 's and b 's
- ▶ However, simple application of 3-term recurrence might not always be stable due to cancellation (coefficients a_k, b_k can be negative)
- ▶ Cancellation errors can be avoided with clever numerics (Sec 6.3 in Deuffhard)

Numerical Methods I

MATH-GA 2010.001/CSCI-GA 2420.001

Benjamin Peherstorfer
Courant Institute, NYU

Based on slides by G. Stadler and A. Donev

Today

Last time

- ▶ Polynomial interpolation

Today

- ▶ Interpolation

Announcements

- ▶ Homework 5 is due Mon, Nov 18 before class

Recap: Function approximation

Consider a function $f \in \mathcal{V}$ in a function space \mathcal{V} . Let now ϕ_1, \dots, ϕ_n be a basis of an n -dimensional space \mathcal{V}_n .

The task that we are interested in is finding a function $f^* \in \mathcal{V}_n$ that approximates f with coefficients c_1, \dots, c_n :

$$f^*(x) = \sum_{i=1}^n c_i \phi_i(x).$$

If we have an inner product, what is the best approximation? The best-approximation of f in \mathcal{V}_n w.r.t. the induced norm is given by the projection

$$f^* = \Pi_n f,$$

where Π_n is the orthogonal projection onto \mathcal{V}_n .

How can we compute the projection? In many cases, we cannot directly compute the projection of f onto \mathcal{V}_n because we have “too little knowledge about f ” Instead? \rightsquigarrow interpolation (/regression)

Recap: Interpolation

Consider n pairs of data samples $(x_i, y_i), i = 1, \dots, n$ with

$$y_i = f(x_i)$$

Based on $\{(x_i, y_i)\}_{i=1}^n$, we now would like to find an approximation $\tilde{f} \in \mathcal{V}_n$ that is “close” to f .

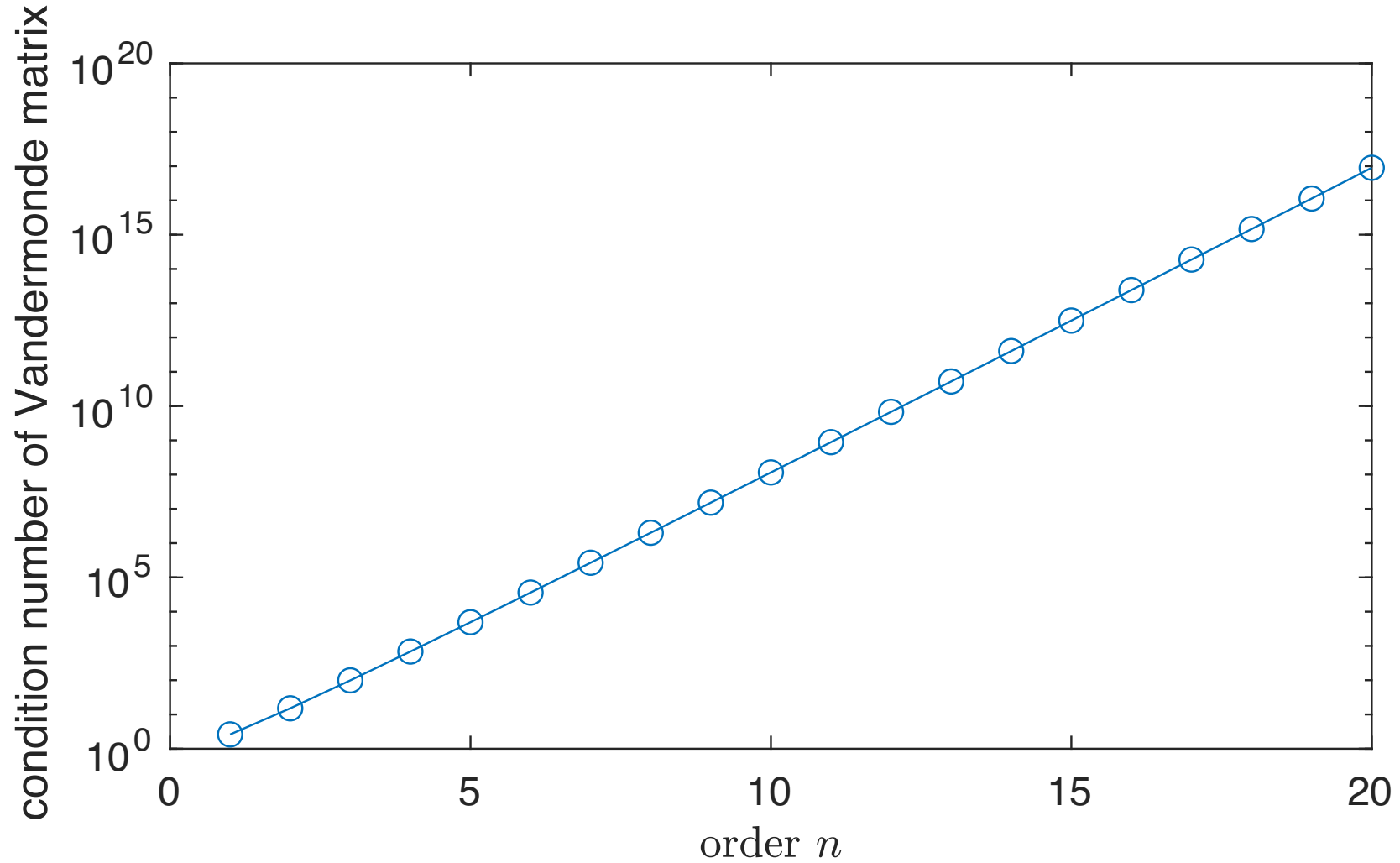
For example, we could enforce the interpolation condition, namely that it holds

$$\tilde{f}(x_i) = f(x_i), \quad i = 1, \dots, n$$

We could also use regression ($m > n$) and minimize, e.g.,

$$\frac{1}{m} \sum_{i=1}^m |y_i - \tilde{f}(x_i)|^2$$

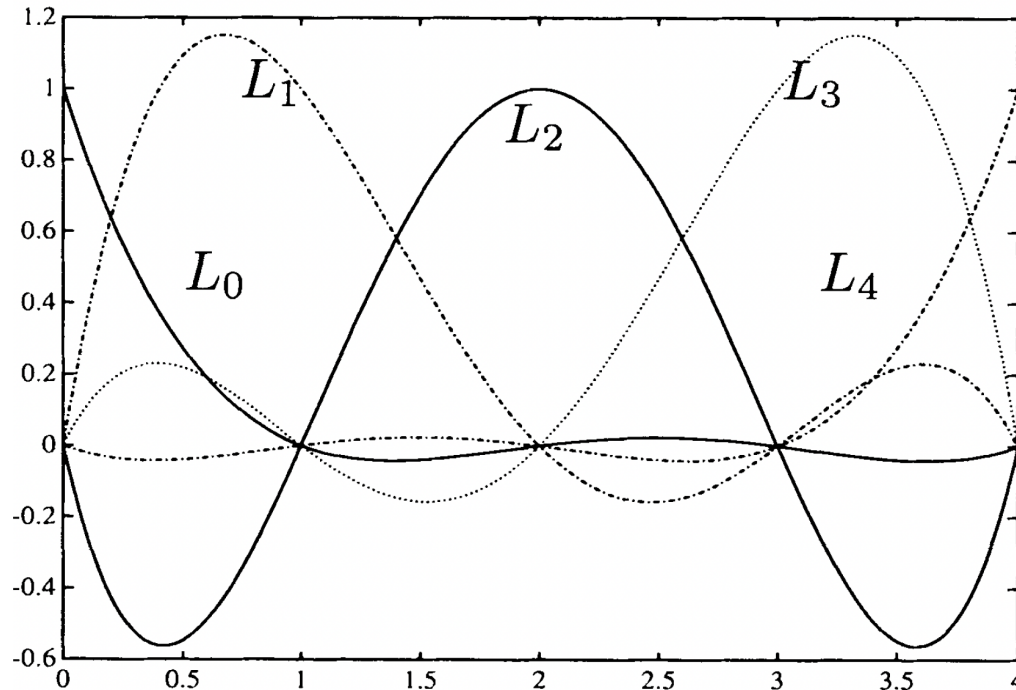
Condition number of Vandermonde matrix



Recap: Lagrange basis

The Lagrange polynomials $L_0, \dots, L_n \in \mathbb{P}_n$ are uniquely defined for distinct x_0, \dots, x_n

$$L_i(x_j) = \delta_{ij}, \quad L_i \in \mathbb{P}_n.$$



Lagrange polynomials up to order $n = 4$ for equidistant x_0, \dots, x_4 . [Figure: Deuffhard]

Recap: Orthogonal polynomials

Define an **inner product between functions**:

$$(f, g) = \int_a^b \omega(x) f(x) g(x) dx,$$

where $\omega(x) > 0$ for $a \leq x \leq b$ is a **weight function**. The **induced norm** is $\|f\| := \sqrt{(f, f)}$.

Let $P_0, P_1, P_2, \dots, P_K$ be polynomials of $0, 1, 2, \dots, K$ order, respectively. They are called orthogonal polynomials on $[a, b]$ with respect to the weight function $\omega(x)$ if it holds

$$(P_i, P_j) = \int_a^b \omega(x) P_i(x) P_j(x) dx = \delta_{ij} \gamma_i, \quad i, j = 0, \dots, K,$$

with $\gamma_i = \|P_i\|^2 > 0$.

Recap

To define orthogonal polynomials uniquely, we normalize them so that the leading coefficient is one, i.e.,

$$P_k(x) = x^k + \dots$$

Theorem: There exist uniquely determined orthogonal polynomials $P_k \in \mathbb{P}_k$ with leading coefficient 1. These polynomials satisfy the 3-term recurrence relation:

$$P_k(x) = (x + a_k)P_{k-1}(x) + b_k P_{k-2}(x), \quad k = 2, 3, \dots$$

with starting values $P_0 = 1$, $P_1 = x + a_1$, where

$$a_k = -\frac{(xP_{k-1}, P_{k-1})}{(P_{k-1}, P_{k-1})}, \quad b_k = -\frac{(P_{k-1}, P_{k-1})}{(P_{k-2}, P_{k-2})}$$

Chebyshev polynomials

Chebyshev polynomials for $-1 \leq x \leq 1$ are given by the 3-term recurrence: $T_0(x) = 1$, $T_1(x) = x$,

$$T_k(x) = 2xT_{k-1}(x) - T_{k-2}(x) \quad \text{for } k \geq 2.$$

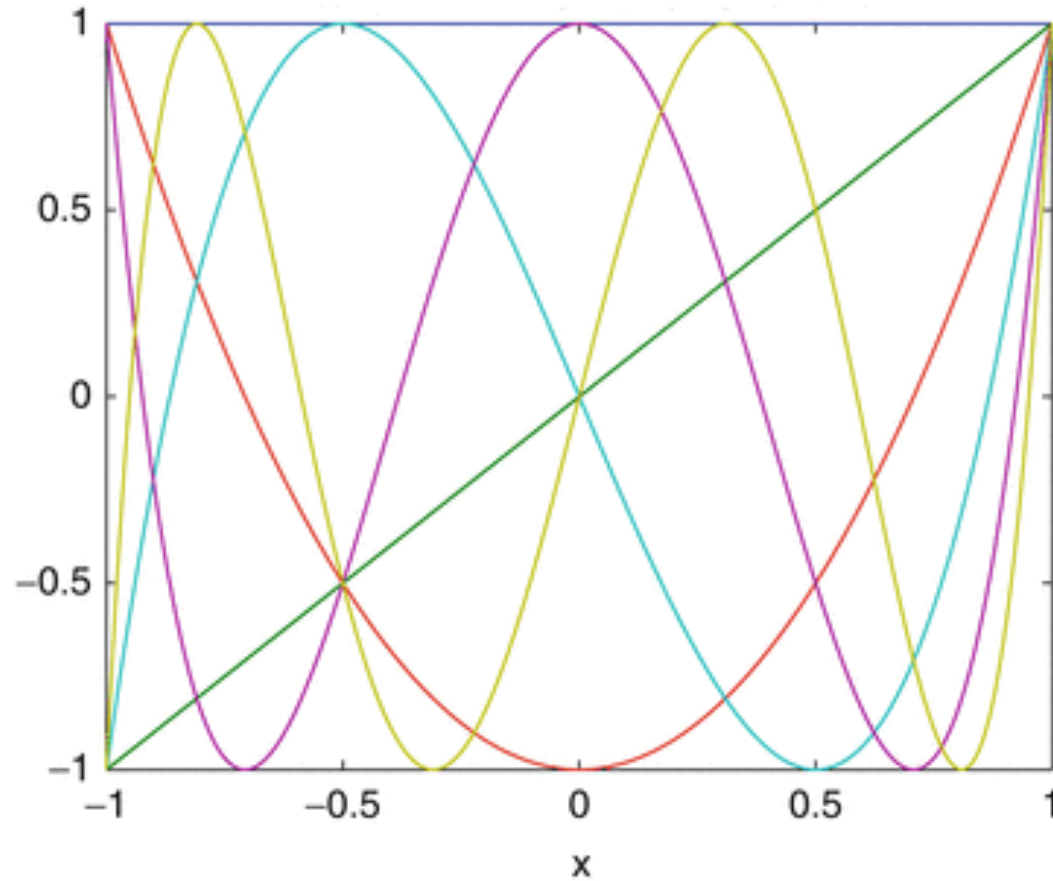
Chebyshev polynomials are orthogonal in the following inner product

$$\int_{-1}^1 \frac{1}{\sqrt{1-x^2}} T_n(x) T_m(x) dx = \begin{cases} 0 & n \neq m \\ \pi & n = m = 0 \\ \pi/2 & n = m \neq 0 \end{cases}$$

The roots of Chebyshev polynomials play an important role in interpolation and numerical quadrature.

Chebyshev polynomials are also given by $T_k(x) = \cos(k \arccos(x))$ with roots

$$x_i = \cos\left(\frac{\pi(i-1/2)}{k}\right)$$



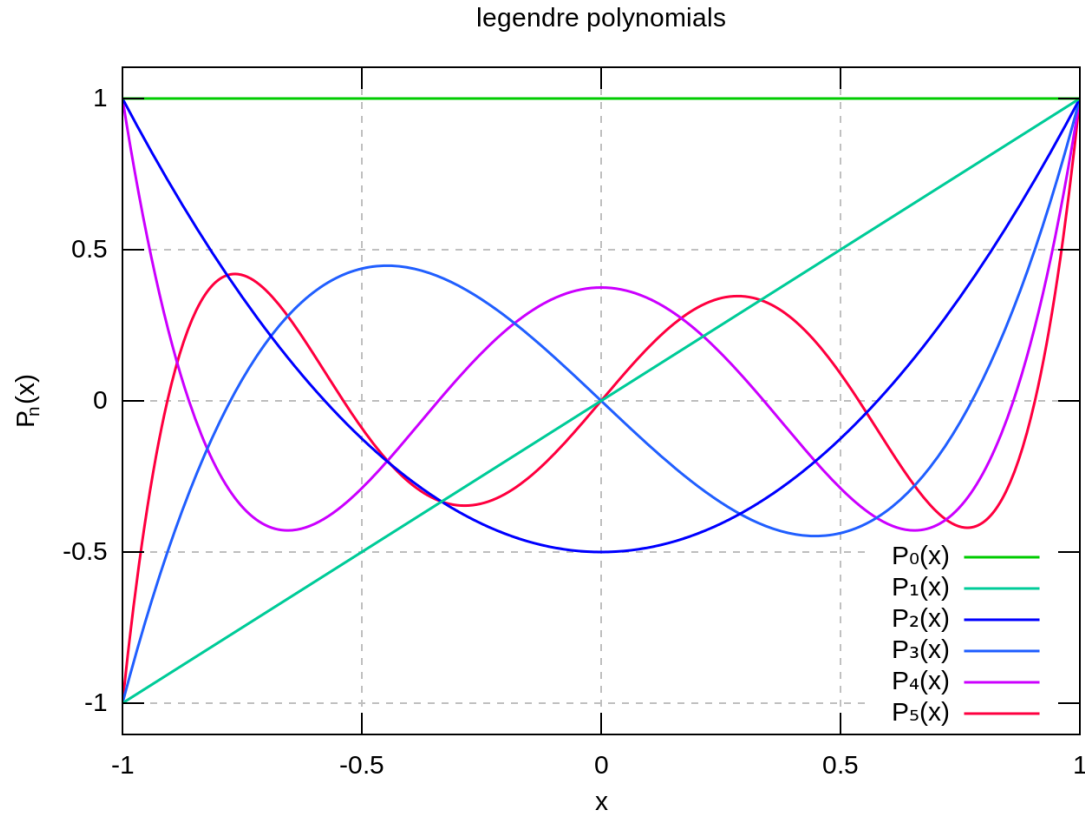
Source: Springer, Encyclopedia of Applied and Computational Mathematics.

Question: What property of the roots of the polynomials do you observe?

Legendre polynomials

Orthogonal polynomials for weight function $\omega \equiv 1$, satisfy $L_0 = 1$, $L_1 = x$, and

$$L_{k+1}(x) = \frac{2k+1}{k+1}xL_k(x) - \frac{k}{k+1}L_{k-1}(x)$$



Approximation error of polynomial interpolation

Let $f : I \rightarrow \mathbb{R}$ and let $x_0, x_1, \dots, x_n \in I$ be $n + 1$ distinct nodes. Assume $f \in C^{n+1}(I)$. Then the interpolation error at point $x \in I$ is

$$E_n(x) = f(x) - P_f(x|x_0, \dots, x_n) = \frac{f^{(n+1)}(\xi)}{(n+1)!} \omega_{n+1}(x),$$

where $\xi \in I$ and

$$\omega_{n+1}(x) = \prod_{i=0}^n (x - x_i),$$

is the nodal polynomial.

Proof: \rightsquigarrow textbook

What are observations can we make?

Approximation error of polynomial interpolation

Let $f : I \rightarrow \mathbb{R}$ and let $x_0, x_1, \dots, x_n \in I$ be $n + 1$ distinct nodes. Assume $f \in C^{n+1}(I)$. Then the interpolation error at point $x \in I$ is

$$E_n(x) = f(x) - P_f(x|x_0, \dots, x_n) = \frac{f^{(n+1)}(\xi)}{(n+1)!} \omega_{n+1}(x),$$

where $\xi \in I$ and

$$\omega_{n+1}(x) = \prod_{i=0}^n (x - x_i),$$

is the nodal polynomial.

Proof: \rightsquigarrow textbook

What are observations can we make?

- ▶ The error bound requires increasing smoothness with the degree n \rightsquigarrow can we derive bounds for $f \in C^k$ with k fixed independent of degree n ?
- ▶ The bound critically depends on ω_{n+1} \rightsquigarrow contains information about the nodes

Let's say the function f is continuous, i.e., $f \in C^0([a, b])$, but not even differentiable.

Define the maximum norm

$$\|f\|_\infty = \max_{x \in [a, b]} |f(x)|$$

Given points $X = \{x_0, \dots, x_n\} \subset [a, b]$, define the interpolation error

$$E_{n, \infty}(X) = \|f - P_f(\cdot | X)\|_\infty$$

and the best-approximation error with $f^* \in \mathbb{P}_n$

$$E_n^* = \|f - f^*\|_\infty \leq \|f - \tilde{f}\|_\infty, \quad \forall \tilde{f} \in \mathbb{P}_n$$

Side remark: What does the Stone-Weierstrass theorem tell us about the best-approximation of continuous functions with polynomials?

Let's say the function f is continuous, i.e., $f \in C^0([a, b])$, but not even differentiable.

Define the maximum norm

$$\|f\|_\infty = \max_{x \in [a, b]} |f(x)|$$

Given points $X = \{x_0, \dots, x_n\} \subset [a, b]$, define the interpolation error

$$E_{n, \infty}(X) = \|f - P_f(\cdot | X)\|_\infty$$

and the best-approximation error with $f^* \in \mathbb{P}_n$

$$E_n^* = \|f - f^*\|_\infty \leq \|f - \tilde{f}\|_\infty, \quad \forall \tilde{f} \in \mathbb{P}_n$$

Side remark: What does the Stone-Weierstrass theorem tell us about the best-approximation of continuous functions with polynomials? Universal approximation: Uniform convergence (i.e., converge in $\|\cdot\|_\infty$). However, the sequence of polynomials that converges is not necessarily obtained via interpolation.

Theorem Let $f \in C^0([a, b])$ and $X = \{x_0, \dots, x_n\} \subset [a, b]$. Then

$$E_{n,\infty}(X) = E_n^*(1 + \Lambda_n(X)),$$

where $\Lambda_n(X)$ is the *Lebesgue* constant of X , defined as

$$\Lambda_n(X) = \left\| \sum_{j=0}^n |L_j^{(n)}| \right\|_{\infty},$$

where $L_j^{(n)} \in \mathbb{P}_n$ is the j -th Lagrange polynomial with

$$L_j^{(n)}(x_i) = \begin{cases} 1, & i = j, \\ 0, & \text{otherwise} \end{cases}$$

Notice the decomposition of the error into component E_n^* (independent of X but dependent on f) and $\Lambda_n(X)$ (independent of f but dependent on X)

We control the best-approximation error E_n^* with the space \mathbb{P}_n

We control the Lebesgue constant $\Lambda_n(X)$ with the nodes x_0, \dots, x_n

What are two important questions to ask about the Lebesgue constant $\Lambda_n(X)$?

*Not necessarily nested

We control the best-approximation error E_n^* with the space \mathbb{P}_n

We control the Lebesgue constant $\Lambda_n(X)$ with the nodes x_0, \dots, x_n

What are two important questions to ask about the Lebesgue constant $\Lambda_n(X)$?

- ▶ What are “good” nodes x_0, \dots, x_n that keep $\Lambda_n(X)$ low or even minimize it?
- ▶ In the best case, how does $\Lambda_n(X)$ behave with respect to $n \rightarrow \infty$

*Not necessarily nested

We control the best-approximation error E_n^* with the space \mathbb{P}_n

We control the Lebesgue constant $\Lambda_n(X)$ with the nodes x_0, \dots, x_n

What are two important questions to ask about the Lebesgue constant $\Lambda_n(X)$?

- ▶ What are “good” nodes x_0, \dots, x_n that keep $\Lambda_n(X)$ low or even minimize it?
- ▶ In the best case, how does $\Lambda_n(X)$ behave with respect to $n \rightarrow \infty$

Famous result by Faber (1914): Given a sequence of any nodes*

$X_n = \{x_{n,0}, x_{n,1}, \dots, x_{n,n}\} \subset [a, b]$, then there always exists a continuous function f so that $P_f(\cdot | x_0, \dots, x_n)$ does not converge to f in $\|\cdot\|_\infty$ for $n \rightarrow \infty$

Thus, polynomial interpolation does not allow for approximating any continuous function

*Not necessarily nested

We control the best-approximation error E_n^* with the space \mathbb{P}_n

We control the Lebesgue constant $\Lambda_n(X)$ with the nodes x_0, \dots, x_n

What are two important questions to ask about the Lebesgue constant $\Lambda_n(X)$?

- ▶ What are “good” nodes x_0, \dots, x_n that keep $\Lambda_n(X)$ low or even minimize it?
- ▶ In the best case, how does $\Lambda_n(X)$ behave with respect to $n \rightarrow \infty$

Famous result by Faber (1914): Given a sequence of any nodes*

$X_n = \{x_{n,0}, x_{n,1}, \dots, x_{n,n}\} \subset [a, b]$, then there always exists a continuous function f so that $P_f(\cdot | x_0, \dots, x_n)$ does not converge to f in $\|\cdot\|_\infty$ for $n \rightarrow \infty$

Thus, polynomial interpolation does not allow for approximating any continuous function

However, interpolation works fantastic for “most” functions \rightsquigarrow “Six myths of polynomial interpolation and quadrature,” Trefethen, *Approximation Theory and Approximation Practice* and also chebfun.org

*Not necessarily nested

It also has been shown that for any possible choice $X_n = \{x_{n,0}, \dots, x_{n,n}\}$, there exists a constant $C > 0$ such that

$$\Lambda_n(X_n) > \frac{2}{\pi} \log(n+1) - C, \quad n = 0, 1, \dots$$

Thus $\Lambda_n(X_n) \rightarrow \infty$ for $n \rightarrow \infty$

- ▶ This is an instability statement in the sense that “investing more time” (by selecting more grid points), leads to a larger Lebesgue constant
- ▶ However, the Lebesgue constant might grow very slowly with n (even slower than E_n^* decreases!)

The Lebesgue constant is also the absolute condition number of polynomial interpolation on $[a, b]$ with points X

$$\kappa_{\text{abs}} = \Lambda_n(X)$$

↪ textbook by Deuffhard

Let's approximate

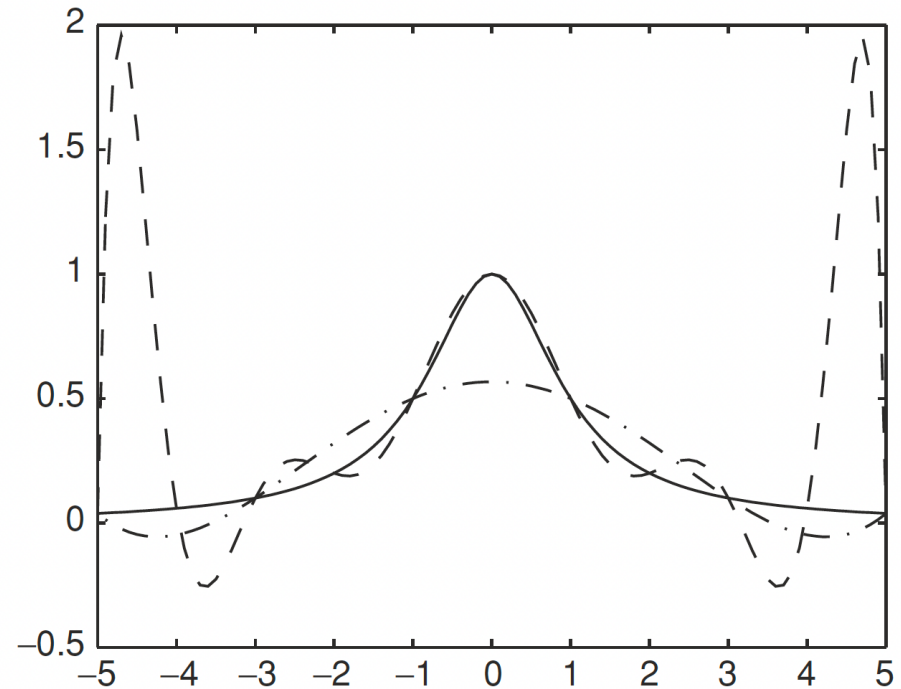
$$f(x) = \frac{1}{1+x^2}, \quad -5 \leq x \leq 5,$$

using polynomial interpolation on equally spaced nodes in $[-5, 5]$.

Let's approximate

$$f(x) = \frac{1}{1+x^2}, \quad -5 \leq x \leq 5,$$

using polynomial interpolation on equally spaced nodes in $[-5, 5]$.



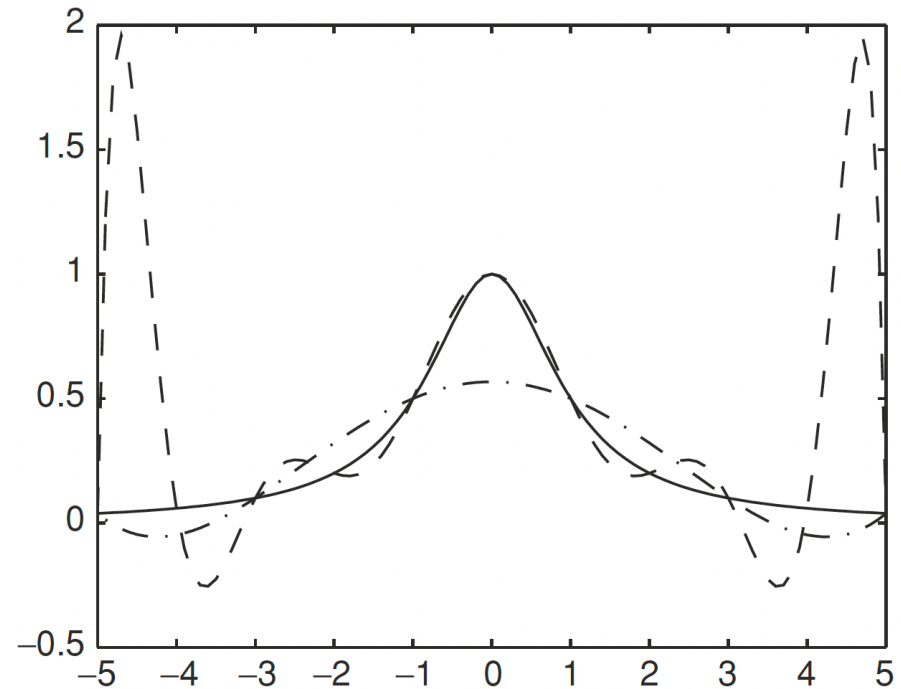
[Figure: Quarteroni]

Interpolants of degree $n = 5$ and $n = 10$ of $f(x) = 1/(1+x^2)$ on equidistant nodes.

Let's approximate

$$f(x) = \frac{1}{1+x^2}, \quad -5 \leq x \leq 5,$$

using polynomial interpolation on equally spaced nodes in $[-5, 5]$.

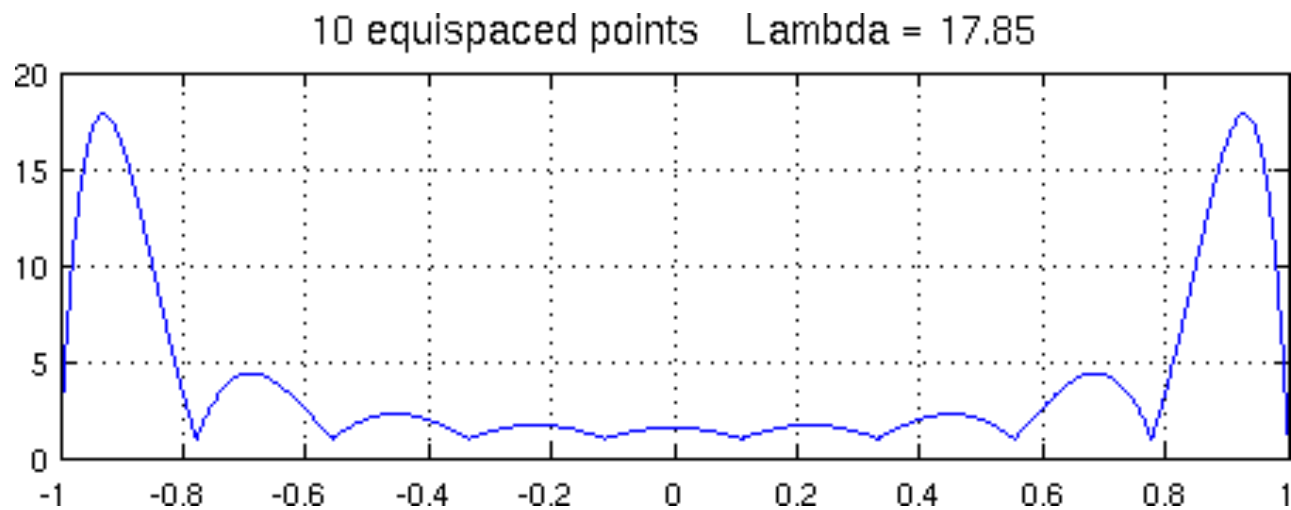


[Figure: Quarteroni]

Interpolants of degree $n = 5$ and $n = 10$ of $f(x) = 1/(1+x^2)$ on equidistant nodes. It can be shown that polynomial interpolation does not converge for $|x| > 4$ for this f

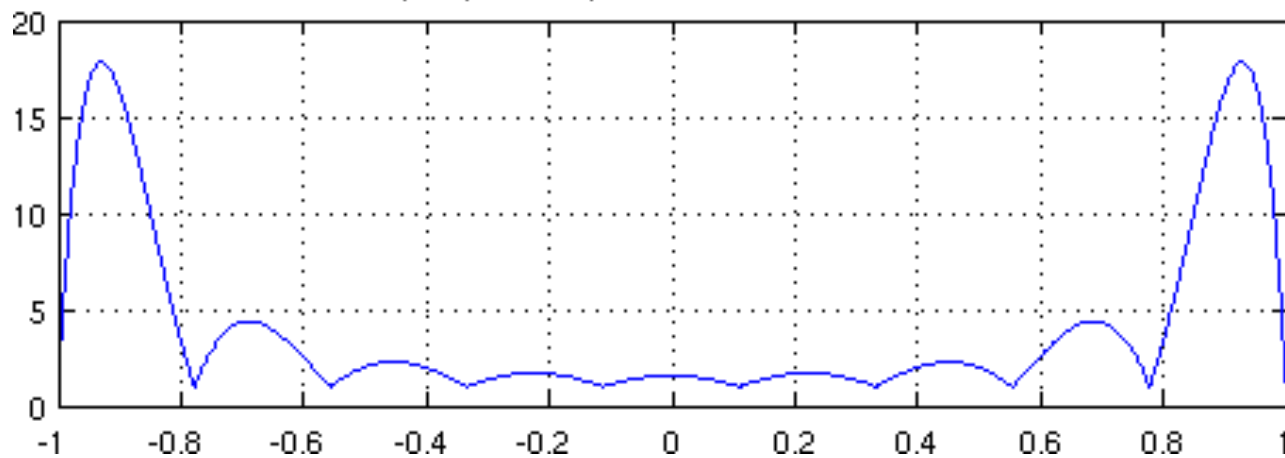
It is a very common situation that interpolation on equidistant nodes leads to high oscillations near the interval ends \rightsquigarrow **Runge's phenomenon**

For equidistant nodes, the Lebesgue constant grows dramatically near the interval ends

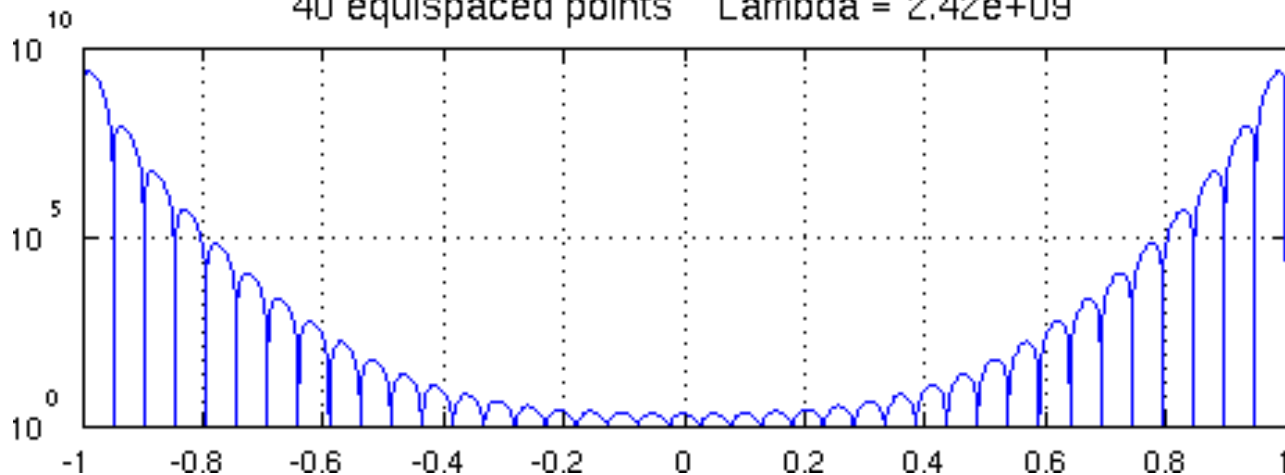


For equidistant nodes, the Lebesgue constant grows dramatically near the interval ends

10 equispaced points $\Lambda = 17.85$

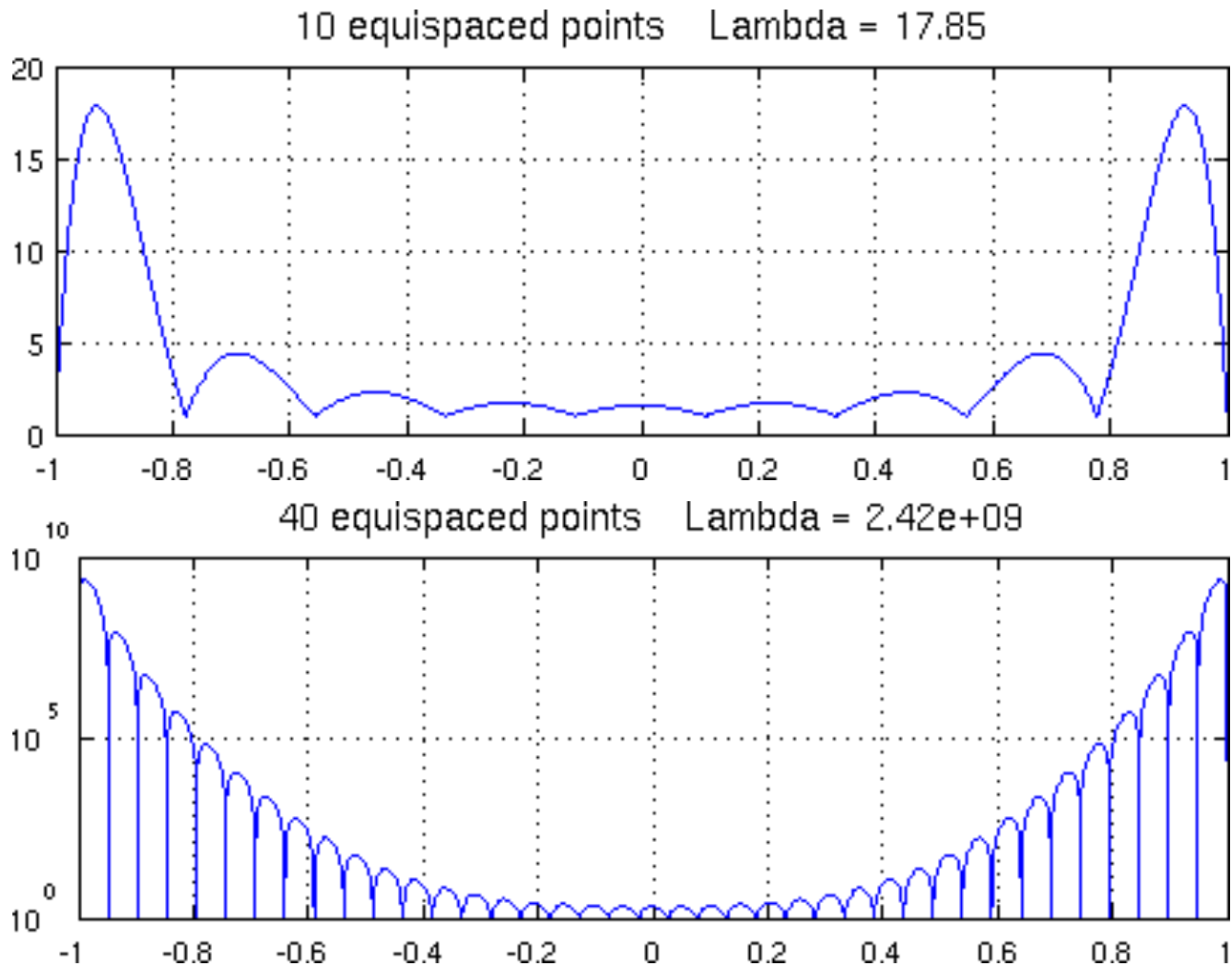


40 equispaced points $\Lambda = 2.42e+09$



What is a way out of this?

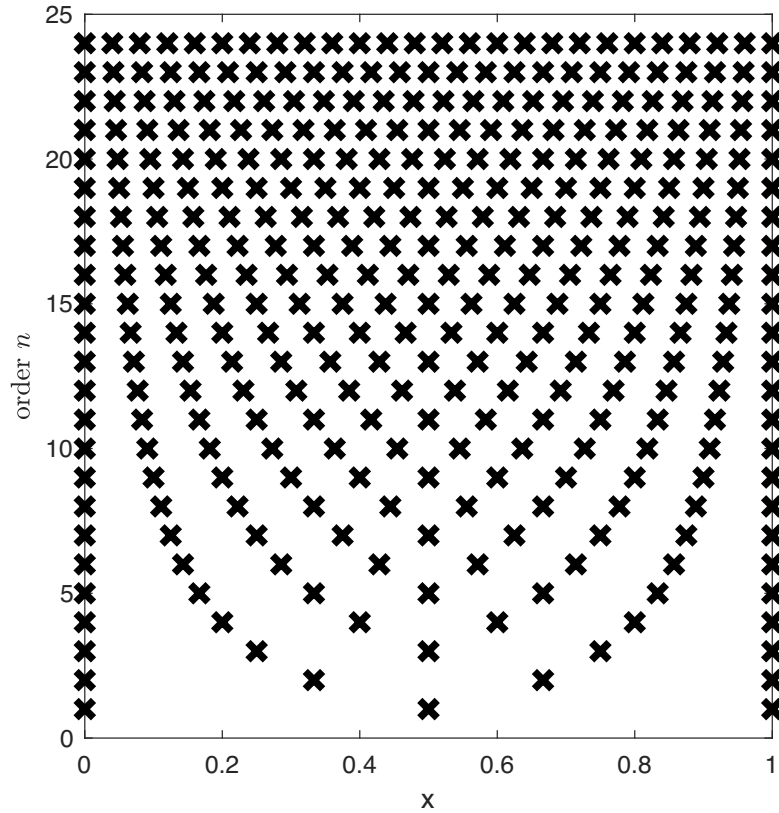
For equidistant nodes, the Lebesgue constant grows dramatically near the interval ends



What is a way out of this? \rightsquigarrow Use different nodes

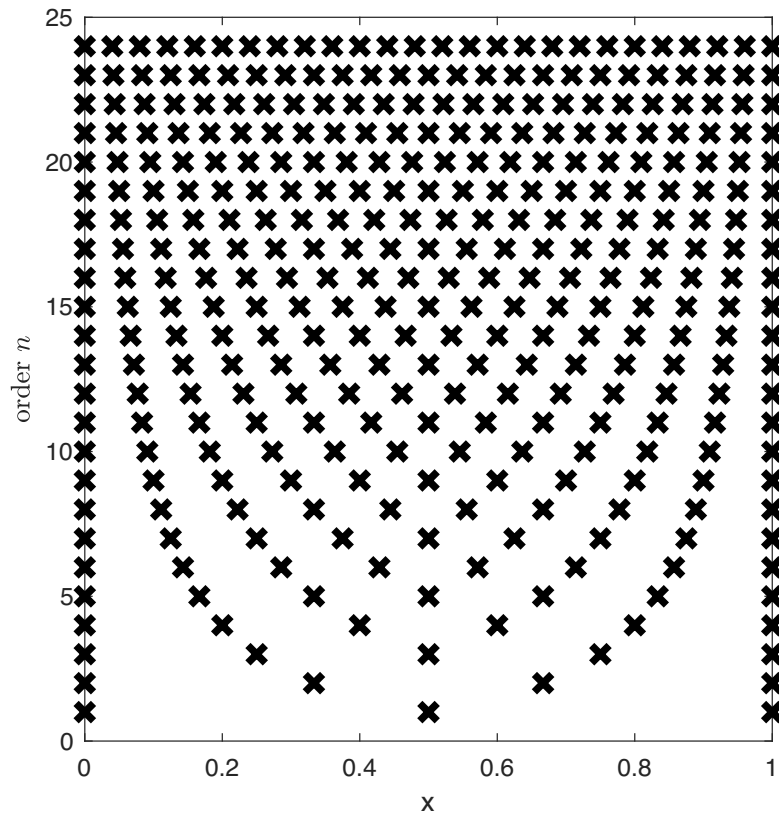
Non-equispaced points

equispaced

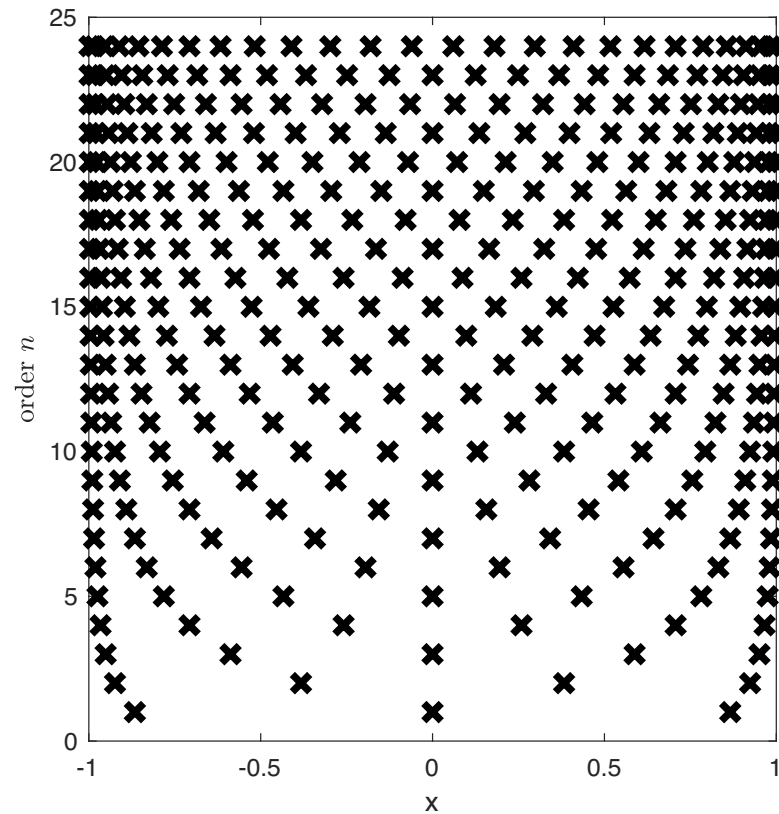


Non-equispaced points

equispaced



Chebyshev nodes



Lebesgue constants for different orders:

n	Λ_n for equidistant nodes	Λ_n for Chebyshev nodes
5	3.106292	2.104398
10	29.890695	2.489430
15	512.052451	2.727778
20	10986.533993	2.900825

Chebyshev nodes are the roots of the Chebyshev polynomials:

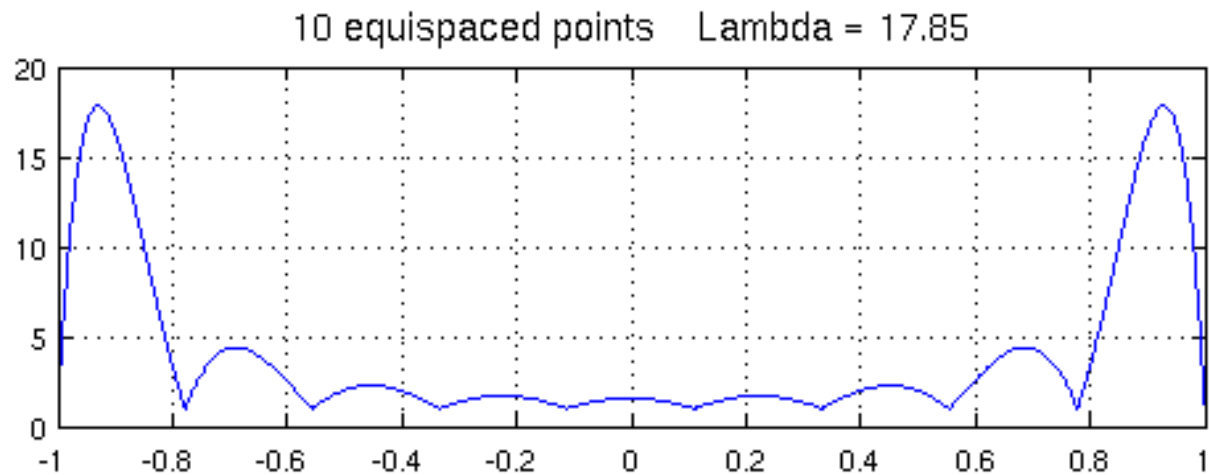
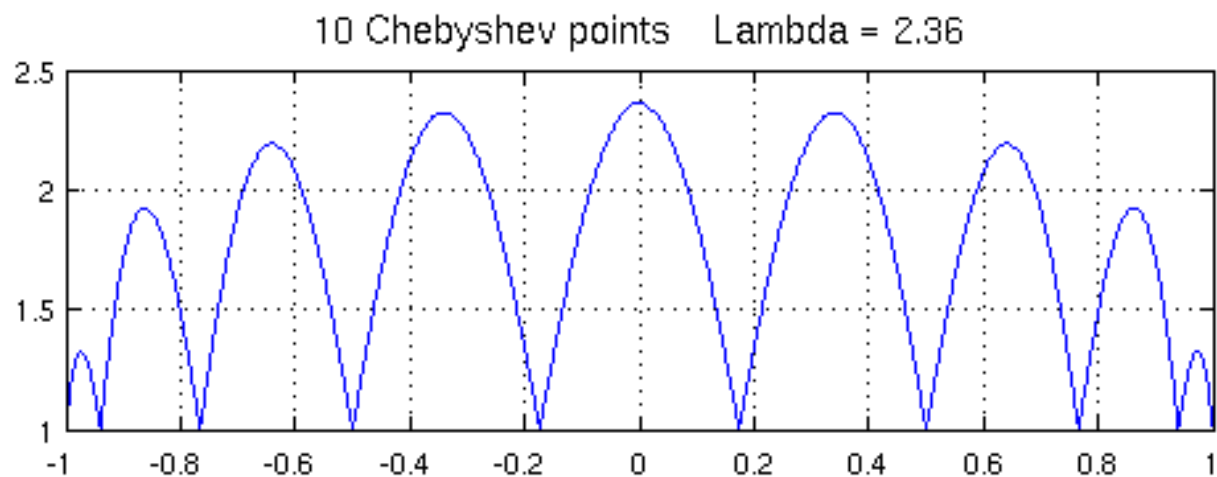
$$t_i = \cos\left(\frac{2i+1}{2n+2}\pi\right), \text{ for } i = 0, \dots, n$$

Lebesgue constant for Chebyshev nodes is bounded as

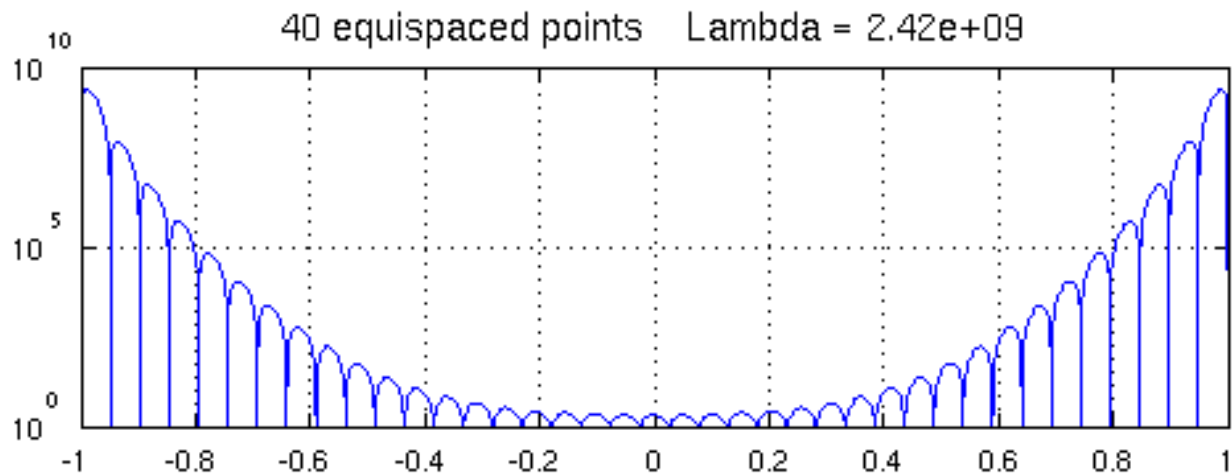
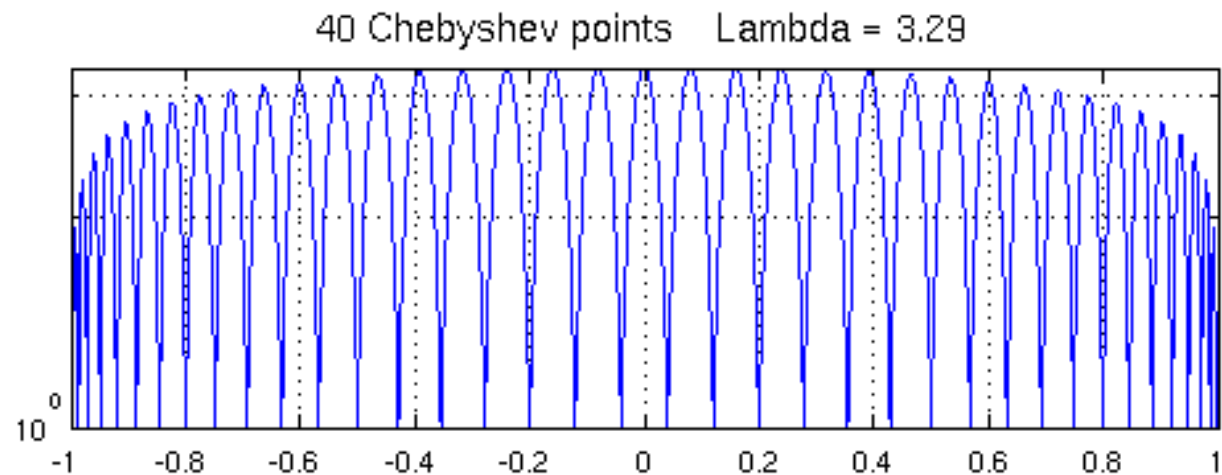
$$\Lambda_n \leq \frac{2}{\pi} \log(n+1) + 1,$$

which is close to the lower bound from previous slides

Lebesgue constant for $n = 10$, uniform vs. Chebyshev nodes:



Lebesgue constant for $n = 40$, uniform vs. Chebyshev nodes:



The Lebesgue constant for Chebyshev points is bounded as

$$\Lambda_n \leq \frac{2}{\pi} \log(n+1) + 1,$$

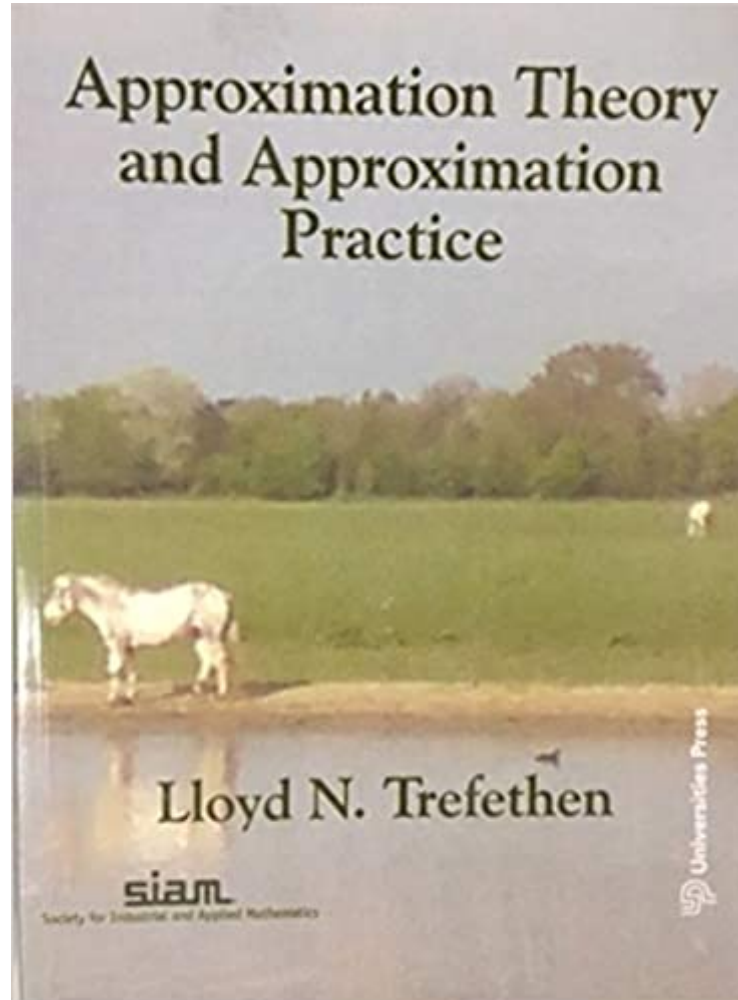
We know that for all interpolation points there exists a continuous function f such that polynomial interpolation does not converge (result by Faber, see previous slides)

This is also true for Chebyshev points! However, if f is Lipschitz continuous, then one can show uniform convergence.

The smoother the function (e.g., continuously differentiable, ν -times continuously differentiable, analytic), the faster interpolation on Chebyshev nodes converges uniformly.

This is a general principle: The more regularity there is in the function, the easier we can approximate it (faster convergence w.r.t. number of degrees of freedom of the approximation)

Beyond the material we cover here...



Pointwise evaluation of interpolating polynomials

Aitken Lemma: The interpolating polynomial $P_f(\cdot|x_0, \dots, x_n)$ satisfies the recurrence relation

$$P_f(x|x_0, \dots, x_n) = \frac{(x_0 - x)P_f(x|x_1, \dots, x_n) - (x_n - x)P_f(x|x_0, \dots, x_{n-1})}{x_0 - x_n}$$

Proof \rightsquigarrow board

Aitken's lemma

$$\varphi(x) = \frac{(x_0 - x) P_f(x | x_1, \dots, x_n) - (x_n - x) P_f(x | x_0, \dots, x_{n-1})}{x_0 - x_n}$$

We now show that φ interpolates f at x_0, \dots, x_n . Because we have uniqueness of interpolating polynomials, it follows $\varphi = P_f(\cdot | x_0, \dots, x_n)$

for $i = \underline{1}, \dots, \underline{n-1}$

$$\varphi(x_i) = \frac{(x_0 - x_i) f(x_i) - (x_n - x_i) f(x_i)}{x_0 - x_n}$$

$$= \frac{x_0 - x_i - x_n + x_i}{x_0 - x_n} f(x_i) = f(x_i)$$

$$\varphi(x_0) = \frac{0 - (x_n - x_0) f(x_0)}{x_0 - x_n} = f(x_0)$$

$$\varphi(x_n) = \frac{(x_0 - x_n) f(x_n) - 0}{x_0 - x_n} = f(x_n)$$

$$\varphi = P_f(\cdot | x_0, \dots, x_n)$$

Pointwise evaluation of interpolating polynomials

Aitken Lemma: The interpolating polynomial $P_f(\cdot|x_0, \dots, x_n)$ satisfies the recurrence relation

$$P_f(x|x_0, \dots, x_n) = \frac{(x_0 - x)P_f(x|x_1, \dots, x_n) - (x_n - x)P_f(x|x_0, \dots, x_{n-1})}{x_0 - x_n}$$

Proof \rightsquigarrow board

Introduce the notation

$$P_{ik} = P_f(x|x_{i-k}, \dots, x_i), \quad i \geq k$$

then $P_f(x|x_0, \dots, x_n) = P_{nn}$ can be computed based on the Neville scheme

$$P_{i0} = f(x_i), \quad i = 0, \dots, n$$

$$P_{ik} = P_{i,k-1} + \frac{x - x_i}{x_i - x_{i-k}} (P_{i,k-1} - P_{i-1,k-1}), \quad i \geq k$$

\rightsquigarrow a numerically stable and efficient ($\mathcal{O}(n^2)$) way to **evaluate (!)** $P_f(x|x_0, \dots, x_n)$

Hermite interpolation

Given

$$a = x_0 \leq x_1 \leq \dots \leq x_n = b$$

with possibly **duplicated** nodes. Less information than degrees of freedom?

Hermite interpolation

Given

$$a = x_0 \leq x_1 \leq \dots \leq x_n = b$$

with possibly **duplicated** nodes. Less information than degrees of freedom? Therefore, if the node x_i occurs k times, the corresponding node values correspond to $f(x_i), f'(x_i), \dots, f^{k-1}(x_i)$.

The **Hermite interpolation** polynomial $p(x)$ is a polynomial of order n , which coincides with the nodal values (and, for duplicated nodes, derivatives at nodal values) at the nodes.

Hermite interpolation

Given

$$a = x_0 \leq x_1 \leq \dots \leq x_n = b$$

with possibly **duplicated** nodes. Less information than degrees of freedom? Therefore, if the node x_i occurs k times, the corresponding node values correspond to $f(x_i), f'(x_i), \dots, f^{k-1}(x_i)$.

The **Hermite interpolation** polynomial $p(x)$ is a polynomial of order n , which coincides with the nodal values (and, for duplicated nodes, derivatives at nodal values) at the nodes.

Aitken lemma for Hermite interpolation: The (Hermite) interpolating polynomial $P = P_f(\cdot | x_0, \dots, x_n)$ satisfies, if $x_i \neq x_j$, the recurrence relation:

$$P_f(x | x_0, \dots, x_n) = \frac{(x_j - x)P_f(x | x_0, \dots, \hat{x}_i, \dots, x_n) - (x_i - x)P_f(x | x_0, \dots, \hat{x}_j, \dots, x_n)}{x_j - x_i},$$

where the hat symbol indicates that the corresponding node is omitted.

Polynomial interpolation in Newton basis

The **Newton basis** $\omega_0, \dots, \omega_n$ is given by

$$\omega_i(x) := \prod_{j=0}^{i-1} (x - x_j) \in \mathbb{P}_i.$$

Would like to find coefficients c_0, c_1, \dots, c_n of interpolating polynomial in Newton basis

$$P_f(x|x_0, \dots, x_n) = c_0\omega_0(x) + c_1\omega_1(x) + \dots + c_n\omega_n(x)$$

Newton polynomial basis

The leading coefficient a_n of the interpolation polynomial (monomial basis!)

$$P_f(x|x_0, \dots, x_n) = a_n x^n + \dots + a_0$$

is called the *n-th divided difference*, $[x_0, \dots, x_n]f := a_n$.

The divided differences are the coefficients c_0, \dots, c_n : The interpolation polynomial $P_f(\cdot|x_0, \dots, x_n)$ for $x_0 \leq x_1 \leq \dots \leq x_n$ (not necessarily distinct and thus need $f \in C^{n+1}$) is given by

$$P(x) = \sum_{i=0}^n [x_0, \dots, x_i]f \omega_i(x).$$

Furthermore,

$$f(x) = P(x) + [x_0, \dots, x_n, x]f \omega_{n+1}(x).$$

Proof \rightsquigarrow board

Newton differences

↳ for $n=0$, $p(x) = \sigma_0$

$$p(x_0) = \sigma_0 = [x_0]f = f(x_0)$$

$$w_0(x) = 1$$

↳ $n > 0$ let

$$P_{n-1} = \sum_{i=0}^{n-1} [x_0, \dots, x_i] f w_i(x)$$

interpolates x_0, \dots, x_{n-1}

$$P_n = [x_0, \dots, x_n] f x^n + \sigma_{n-1} x^{n-1} + \dots + \sigma_0$$

as

$$P_n(x) = [x_0, \dots, x_n] f w_n(x) + Q_{n-1}(x)$$

with some polynomial $Q_{n-1} \in \mathbb{P}_{n-1}$ of degree $n-1$

Notice that

$$w_n(x) = \prod_{i=0}^{n-1} (x_i - x)$$

$$w_n(x_i) = 0, \quad i = 0, \dots, n-1$$

$$\begin{aligned} Q_{n-1}(x_i) &= P_n(x_i) - [x_0, \dots, x_n] f w_n(x_i) \\ &= f(x_i) \end{aligned}$$

Q_{n-1} interpolates x_0, \dots, x_{n-1}

$$\Rightarrow Q_{n-1} = \sum_{i=0}^{n-1} [x_0, \dots, x_i] p w_i$$

$$P_n(x) = [x_0, \dots, x_n] p w_n(x) + \sum_{i=0}^{n-1} [x_0, \dots, x_i] p w_i(x)$$

Numerical Methods I

MATH-GA 2010.001/CSCI-GA 2420.001

Benjamin Peherstorfer
Courant Institute, NYU

Based on slides by G. Stadler and A. Donev

Today

Last time

- ▶ Polynomial interpolation

Today

- ▶ Interpolation

Announcements

- ▶ Homework 5 is due Mon, Nov 18 before class

Recap

Let's say the function f is continuous, i.e., $f \in C^0([a, b])$, but not even differentiable.

Define the maximum norm

$$\|f\|_\infty = \max_{x \in [a, b]} |f(x)|$$

Given points $X = \{x_0, \dots, x_n\} \subset [a, b]$, define the interpolation error

$$E_{n, \infty}(X) = \|f - P_f(\cdot | X)\|_\infty$$

and the best-approximation error with $f^* \in \mathbb{P}_n$

$$E_n^* = \|f - f^*\|_\infty \leq \|f - \tilde{f}\|_\infty, \quad \forall \tilde{f} \in \mathbb{P}_n$$

Recap

Theorem Let $f \in C^0([a, b])$ and $X = \{x_0, \dots, x_n\} \subset [a, b]$. Then

$$E_{n,\infty}(X) = E_n^*(1 + \Lambda_n(X)),$$

where $\Lambda_n(X)$ is the *Lebesgue* constant of X , defined as

$$\Lambda_n(X) = \left\| \sum_{j=0}^n |L_j^{(n)}| \right\|_{\infty},$$

where $L_j^{(n)} \in \mathbb{P}_n$ is the j -th Lagrange polynomial with

$$L_j^{(n)}(x_i) = \begin{cases} 1, & i = j, \\ 0, & \text{otherwise} \end{cases}$$

Notice the decomposition of the error into component E_n^* (independent of X but dependent on f) and $\Lambda_n(X)$ (independent of f but dependent on X)

Recap: Pointwise evaluation of interpolating polynomials

Aitken Lemma: The interpolating polynomial $P_f(\cdot|x_0, \dots, x_n)$ satisfies the recurrence relation

$$P_f(x|x_0, \dots, x_n) = \frac{(x_0 - x)P_f(x|x_1, \dots, x_n) - (x_n - x)P_f(x|x_0, \dots, x_{n-1})}{x_0 - x_n}$$

Introduce the notation

$$P_{ik} = P_f(x|x_{i-k}, \dots, x_i), \quad i \geq k$$

then $P_f(x|x_0, \dots, x_n) = P_{nn}$ can be computed based on the Neville scheme

$$P_{i0} = f(x_i), \quad i = 0, \dots, n$$

$$P_{ik} = P_{i,k-1} + \frac{x - x_j}{x_i - x_{i-k}} (P_{i,k-1} - P_{i-1,k-1}), \quad i \geq k$$

\rightsquigarrow a numerically stable and efficient ($\mathcal{O}(n^2)$) way to **evaluate (!)** $P_f(x|x_0, \dots, x_n)$

Recap: Polynomial interpolation in Newton basis

The **Newton basis** $\omega_0, \dots, \omega_n$ is given by

$$\omega_i(x) := \prod_{j=0}^{i-1} (x - x_j) \in \mathbb{P}_i.$$

Would like to find coefficients c_0, c_1, \dots, c_n of interpolating polynomial in Newton basis

$$P_f(x|x_0, \dots, x_n) = c_0\omega_0(x) + c_1\omega_1(x) + \dots + c_n\omega_n(x)$$

Recap: Newton polynomial basis

The leading coefficient a_n of the interpolation polynomial (monomial basis!)

$$P_f(x|x_0, \dots, x_n) = a_n x^n + \dots + a_0$$

is called the *n-th divided difference*, $[x_0, \dots, x_n]f := a_n$.

The divided differences are the coefficients c_0, \dots, c_n : The interpolation polynomial $P_f(\cdot|x_0, \dots, x_n)$ for $x_0 \leq x_1 \leq \dots \leq x_n$ (not necessarily distinct and thus need $f \in C^{n+1}$) is given by

$$P(x) = \sum_{i=0}^n [x_0, \dots, x_i]f \omega_i(x).$$

Furthermore,

$$f(x) = P(x) + [x_0, \dots, x_n, x]f \omega_{n+1}(x).$$

Proof \rightsquigarrow board

cont of proof:

$$P_n = [x_0, \dots, x_n] f w_n + \sum_{i=1}^{n-1} [x_0, \dots, x_i] f w_i$$

$$P_{n+1} = \sum_{i=1}^n [x_0, \dots, x_i] f w_i + [x_0, \dots, x_n, x_{n+1}] f w_{n+1}$$

this is true for any x_{n+1} and $P_{n+1}(x_{n+1}) = f(x_{n+1})$

$$\Rightarrow f(x) = P_n(x) + [x_0, \dots, x_n, \underline{x}] f w_{n+1}(x)$$

Divided differences

The divided differences $[x_0, \dots, x_n]f$ satisfy the following properties:

- ▶ $[x_0, \dots, x_n]P = 0$ for all $P \in \mathbb{P}_{n-1}$ (because $a_n = 0$ for degree $n - 1$ polynomial)
- ▶ If $x_0 = \dots = x_n$:

$$[x_0, \dots, x_n]f = \frac{f^{(n)}(x_0)}{n!}$$

- ▶ The following recurrence relation holds for $x_i \neq x_j$:

$$[x_0, \dots, x_n]f = \frac{([x_0, \dots, \hat{x}_i, \dots, x_n]f - [x_0, \dots, \hat{x}_j, \dots, x_n]f)}{x_j - x_i}$$

- ▶ $[x_0, \dots, x_n]f = \frac{1}{n!}f^{(n)}(\tau)$ with $a \leq \tau \leq b$, if f in C^{n+1}

Divided differences

Let us use divided differences to compute the coefficients for the Newton basis for the cubic interpolation polynomial p that satisfies $p(0) = 1$, $p(0.5) = 2$, $p(1) = 0$, $p(2) = 3$.

$$p(\underline{0}) = 1, \quad p(\underline{0.5}) = 2, \quad p(\underline{1}) = 0, \quad p(\underline{2}) = 3$$

x_i	
0	$[x_0]f = 1$
$\frac{1}{2}$	$[x_1]f = 2$ $[x_0, x_1]f = \frac{[x_1]f - [x_0]f}{x_1 - x_0} = \frac{2 - 1}{\frac{1}{2} - 0} = \underline{2}$
1	$[x_2]f = 0$ $[x_1, x_2]f = \frac{[x_2]f - [x_1]f}{x_2 - x_1} = \frac{0 - 2}{1 - \frac{1}{2}} = \underline{-4}$ →
2	$[x_3]f = 3$

$$[x_0, x_1, x_2]f = \frac{[x_1, x_2]f - [x_0, x_1]f}{x_2 - x_0} = \frac{-4 - 2}{1 - 0} = \underline{-6}$$

$$[x_0, x_1, x_2, x_3]f = \underline{\underline{\frac{16}{3}}}$$

$$p(x) = \sum_{i=0}^3 c_i \omega_i(x)$$

$$\omega_i(x) = \prod_{j=0}^{i-1} (x - x_j)$$

$$= c_0 \cdot 1 + c_1 (x - 0) + c_2 (x - 0)(x - \frac{1}{2}) + c_3 (x - 0)(x - \frac{1}{2})(x - 1)$$

$$= 1 + 2x - 6x(x - \frac{1}{2}) + \frac{16}{3} x(x - \frac{1}{2})(x - 1)$$

Divided differences

Let us use divided differences to compute the coefficients for the Newton basis for the cubic interpolation polynomial p that satisfies $p(0) = 1$, $p(0.5) = 2$, $p(1) = 0$, $p(2) = 3$.

x_j			
0	$[x_0]f = 1$		
0.5	$[x_1]f = 2$	$[x_0x_1]f = \frac{[x_1]f - [x_0]f}{x_1 - x_0} = 2$	
1	$[x_2]f = 0$	$[x_1x_2]f = \frac{[x_2]f - [x_1]f}{x_2 - x_1} = -4$	$[x_0x_1x_2]f = -6$
2	$[x_3]f = 3$	$[x_2x_3]f = \frac{[x_3]f - [x_2]f}{x_3 - x_2} = 3$	$[x_1x_2x_3]f = \frac{14}{3} \quad \frac{16}{3}$

Thus, the interpolating polynomial is

$$p(x) = 1 + 2x + (-6)x(x - 0.5) + \frac{16}{3}x(x - 0.5)(x - 1).$$

Divided differences

Let us now use divided differences to compute the coefficients for the Newton basis for the cubic interpolation polynomial p that satisfies $p(0) = 1$, $p'(0) = 2$, $p''(0) = 1$, $p(1) = 3$.

x_i				
0	$[x_0]f = 1$			
0	$[x_0]f = 1$	$[x_0x_1]f = p'(0) = 2$		
0	$[x_0]f = 1$	$[x_1x_2]f = p'(0) = 2$	$[x_0x_1x_2]f = \frac{p''(0)}{2!} = \frac{1}{2}$	
1	$[x_3]f = 3$	$[x_2x_3]f = \frac{[x_3]f - [x_0]f}{x_3 - x_0} = 2$	$[x_1x_2x_3]f = 0$	$[x_0x_1x_2x_3]f = -\frac{1}{2}$

Thus, the interpolating polynomial is

$$p(t) = 1 + 2t + \frac{1}{2}t^2 + \left(-\frac{1}{2}\right)t^3$$

Polynomial interpolation

- ▶ Polynomial interpolation
- ▶ Hermite interpolation
- ▶ (Least squares with polynomials)
- ▶ What else?

Polynomial interpolation

- ▶ Polynomial interpolation
- ▶ Hermite interpolation
- ▶ (Least squares with polynomials)
- ▶ **What else?** Splines, i.e., piecewise polynomial interpolation

Splines

Assume $(l + 2)$ pairwise disjoint nodes:

$$a = x_0 < x_1 < \dots < x_{l+1} = b.$$

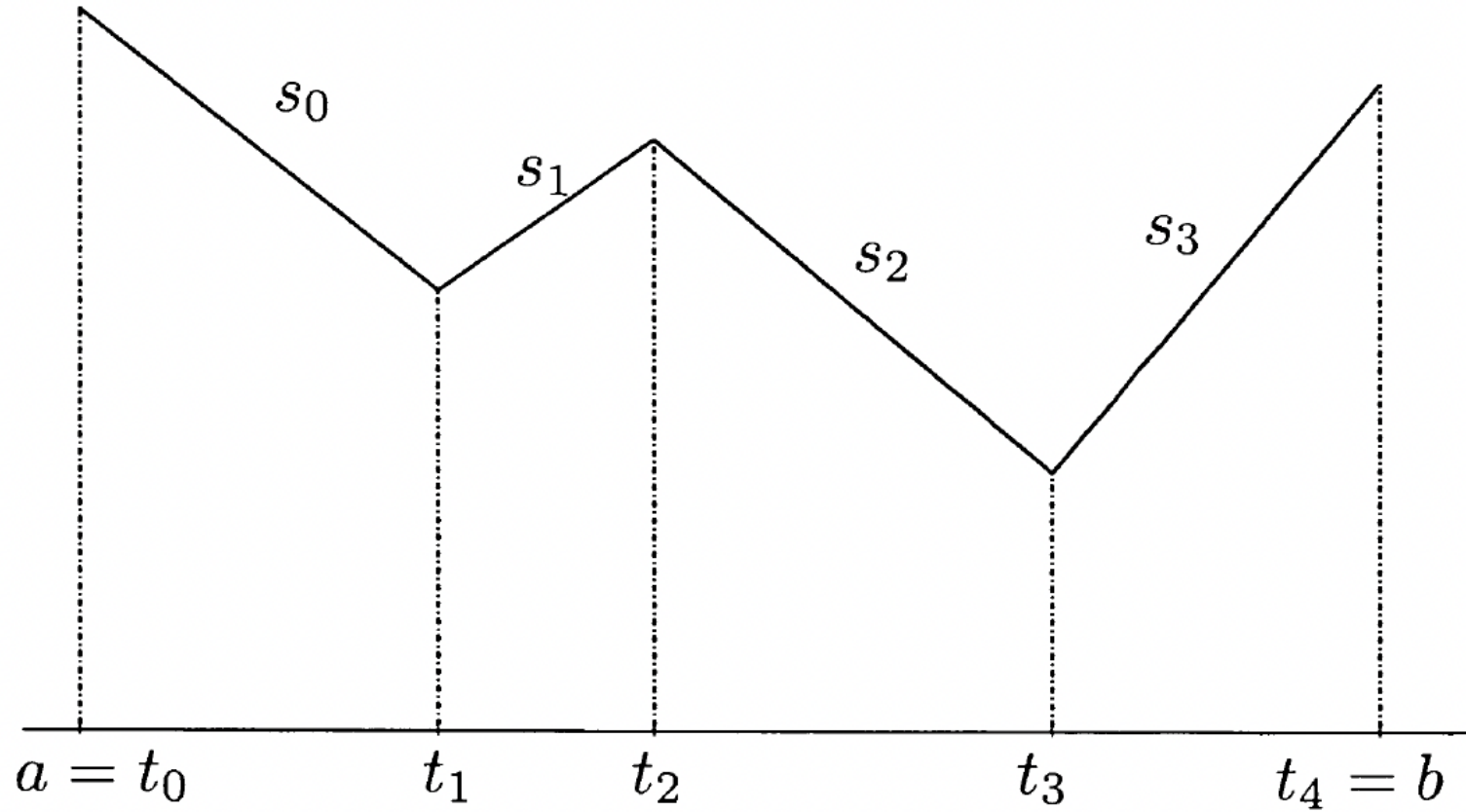
A **spline** of degree $k - 1$ (order k) is a function in C^{k-2} which on each interval $[x_i, x_{i+1}]$ coincides with a polynomial in \mathbb{P}_{k-1} .

Most important examples:

- ▶ linear splines, $k = 2$
- ▶ cubic splines, $k = 4$

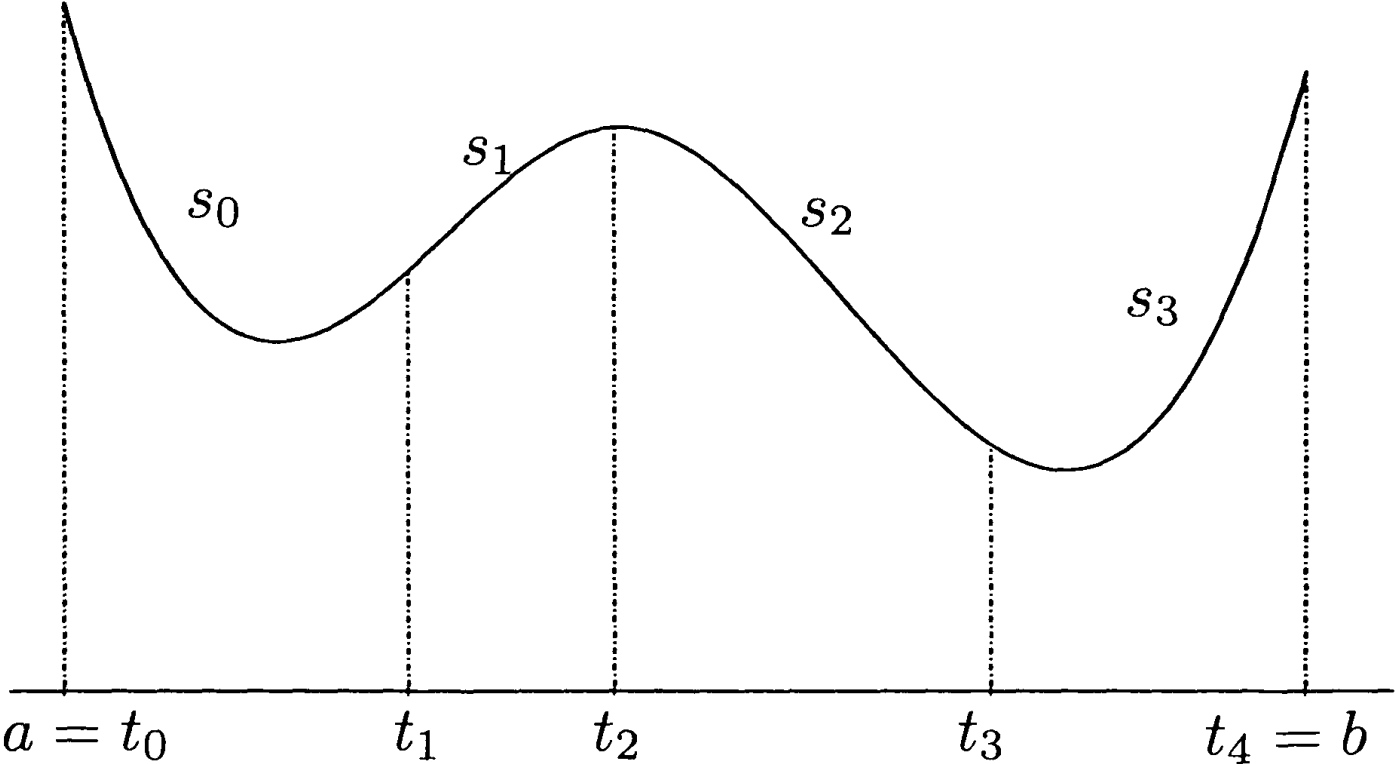
Splines

Linear spline (piecewise linear polynomial)



Splines

Cubic splines look smooth:



Trigonometric Interpolation for periodic functions

Instead of polynomials, use $\sin(jt)$, $\cos(jt)$ for different $j \in \mathbb{N}$.

For $N \geq 1$, we define the set of complex trigonometric polynomials of degree $\leq N - 1$ as

$$\mathcal{T}_{N-1} := \left\{ \sum_{j=0}^{N-1} c_j e^{ijt}, c_j \in \mathbb{C} \right\},$$

where $i = \sqrt{-1}$.

Complex interpolation problem: Given pairwise distinct nodes $t_0, \dots, t_{N-1} \in [0, 2\pi)$ and corresponding nodal values $f_0, \dots, f_{N-1} \in \mathbb{C}$, find a trigonometric polynomial $p \in \mathcal{T}_{N-1}$ such that $p(t_i) = f_i$, for $i = 0, \dots, N - 1$.

- ▶ There exists exactly one $p \in \mathcal{T}_{N-1}$, which solves this interpolation problem.
- ▶ Choose the equidistant nodes $t_k := \frac{2\pi k}{N}$ for $k = 0, \dots, N-1$. The trigonometric polynomial that satisfies $p(t_i) = f_i$ for $i = 0, \dots, N-1$ has the coefficients

$$c_j = \frac{1}{N} \sum_{k=0}^{N-1} e^{-\frac{2\pi ijk}{N}} f_k.$$

- ▶ For equidistant nodes, the linear map from $\mathbb{C}^N \rightarrow \mathbb{C}^N$ defined by $(f_0, \dots, f_{N-1}) \mapsto (c_0, \dots, c_{N-1})$ is called the discrete Fourier transformation (DFT). **The Fast Fourier Transform (FFT) is a (very famous!) algorithm that computes the DFT and its inverse in $O(N \log N)$ flops. \rightsquigarrow Numerical Methods II**

Conclusions

- ▶ Interpolation means approximating function values in the interior of a domain when there are known samples of the function at a set of interior and boundary nodes
- ▶ Given a set of basis functions, interpolation amounts to solving a linear system of equations for the coefficients; there are clever choices to avoid explicitly computing the solution of the system of equations
- ▶ Interpolation on equidistant nodes is a bad idea: Not accurate and not stable! Instead, use Chebyshev nodes, then also higher-order polynomials are safe!
- ▶ Another option is to use piecewise polynomial interpolation such as splines, which is very common when solving PDEs numerically
- ▶ What about higher dimensions?

Solving nonlinear equations

Solving nonlinear equations (“root finding”)

We want to solve the nonlinear equation

$$f(x) = 0, \quad x \in \mathbb{R}.$$

We could also have $n < \infty$ equations in n unknowns with $f : \mathbb{R}^n \rightarrow \mathbb{R}^n$

$$f(\mathbf{x}) = \mathbf{0}$$

In general, we will need an iterative approach that constructs x_1, x_2, x_3, \dots such that

$$\lim_{k \rightarrow \infty} x_k = x^*,$$

with $f(x^*) = 0$.

Solving nonlinear equations (“root finding”)

We want to solve the nonlinear equation

$$f(x) = 0, \quad x \in \mathbb{R}.$$

We could also have $n < \infty$ equations in n unknowns with $f : \mathbb{R}^n \rightarrow \mathbb{R}^n$

$$f(\mathbf{x}) = \mathbf{0}$$

In general, we will need an iterative approach that constructs x_1, x_2, x_3, \dots such that

$$\lim_{k \rightarrow \infty} x_k = x^*,$$

with $f(x^*) = 0$.

What are important properties of a method for solving nonlinear equations?

Solving nonlinear equations (“root finding”)

We want to solve the nonlinear equation

$$f(x) = 0, \quad x \in \mathbb{R}.$$

We could also have $n < \infty$ equations in n unknowns with $f : \mathbb{R}^n \rightarrow \mathbb{R}^n$

$$f(\mathbf{x}) = \mathbf{0}$$

In general, we will need an iterative approach that constructs x_1, x_2, x_3, \dots such that

$$\lim_{k \rightarrow \infty} x_k = x^*,$$

with $f(x^*) = 0$.

What are important properties of a method for solving nonlinear equations?

- ▶ Does it converge? From which starting point x_0 ?

Solving nonlinear equations (“root finding”)

We want to solve the nonlinear equation

$$f(x) = 0, \quad x \in \mathbb{R}.$$

We could also have $n < \infty$ equations in n unknowns with $f : \mathbb{R}^n \rightarrow \mathbb{R}^n$

$$f(\mathbf{x}) = \mathbf{0}$$

In general, we will need an iterative approach that constructs x_1, x_2, x_3, \dots such that

$$\lim_{k \rightarrow \infty} x_k = x^*,$$

with $f(x^*) = 0$.

What are important properties of a method for solving nonlinear equations?

- ▶ Does it converge? From which starting point x_0 ?
- ▶ How quickly does it converge?

Solving nonlinear equations (“root finding”)

We want to solve the nonlinear equation

$$f(x) = 0, \quad x \in \mathbb{R}.$$

We could also have $n < \infty$ equations in n unknowns with $f : \mathbb{R}^n \rightarrow \mathbb{R}^n$

$$f(\mathbf{x}) = \mathbf{0}$$

In general, we will need an iterative approach that constructs x_1, x_2, x_3, \dots such that

$$\lim_{k \rightarrow \infty} x_k = x^*,$$

with $f(x^*) = 0$.

What are important properties of a method for solving nonlinear equations?

- ▶ Does it converge? From which starting point x_0 ?
- ▶ How quickly does it converge?
- ▶ How expensive is each step?

Bisection method

The bisection method exploits that given a continuous function $f : [a, b] \rightarrow \mathbb{R}$, such that $f(a)f(b) < 0$, there exists $x^* \in (a, b)$ with $f(x^*) = 0$

- ▶ Assumption: f is continuous over $[a, b]$ (very weak assumption!)
- ▶ We have chosen a reasonable interval $[a, b]$ so that there exists a solution $x^* \in (a, b)$ with $f(x^*) = 0$

Bisection method

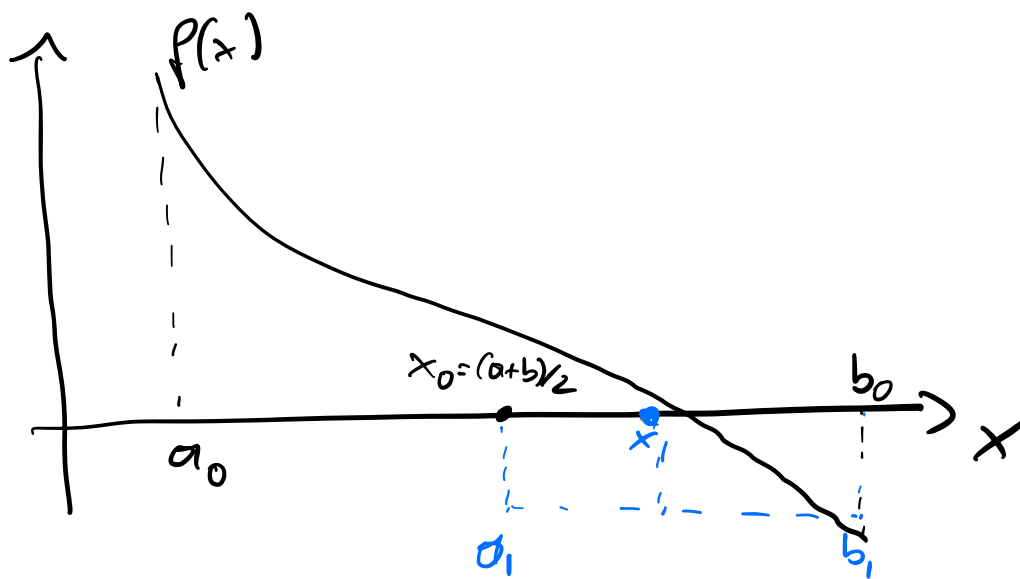
The bisection method exploits that given a continuous function $f : [a, b] \rightarrow \mathbb{R}$, such that $f(a)f(b) < 0$, there exists $x^* \in (a, b)$ with $f(x^*) = 0$

- ▶ Assumption: f is continuous over $[a, b]$ (very weak assumption!)
- ▶ We have chosen a reasonable interval $[a, b]$ so that there exists a solution $x^* \in (a, b)$ with $f(x^*) = 0$

Set $a_0 = a, b_0 = b, x_0 = (a + b)/2$ and iterate for $k = 0, 1, 2, 3, \dots$ as follows:

1. Set $a_{k+1} = a_k, b_{k+1} = x_k$ if $f(x_k)f(a_k) < 0$
2. Set $a_{k+1} = x_k, b_{k+1} = b_k$ if $f(x_k)f(b_k) < 0$
3. Set $x_{k+1} = (a_{k+1} + b_{k+1})/2$
4. Terminate if $|b_{k+1} - a_{k+1}| \leq \epsilon$

Visualization \rightsquigarrow board



Convergence: set $I_k = [a_k, b_k]$

$$|I_k| = b_k - a_k = \frac{b-a}{2^k}$$

Denote the error $e_k = x_k - x^*$ of step k

$$|e_k| = |x_k - x^*| \leq \frac{1}{2} |I_k| = \frac{b-a}{2^{k+1}}$$

$$\lim_{k \rightarrow \infty} |e_k| = 0$$

\Rightarrow global convergence

\Rightarrow non-monotone convergence (see plot)

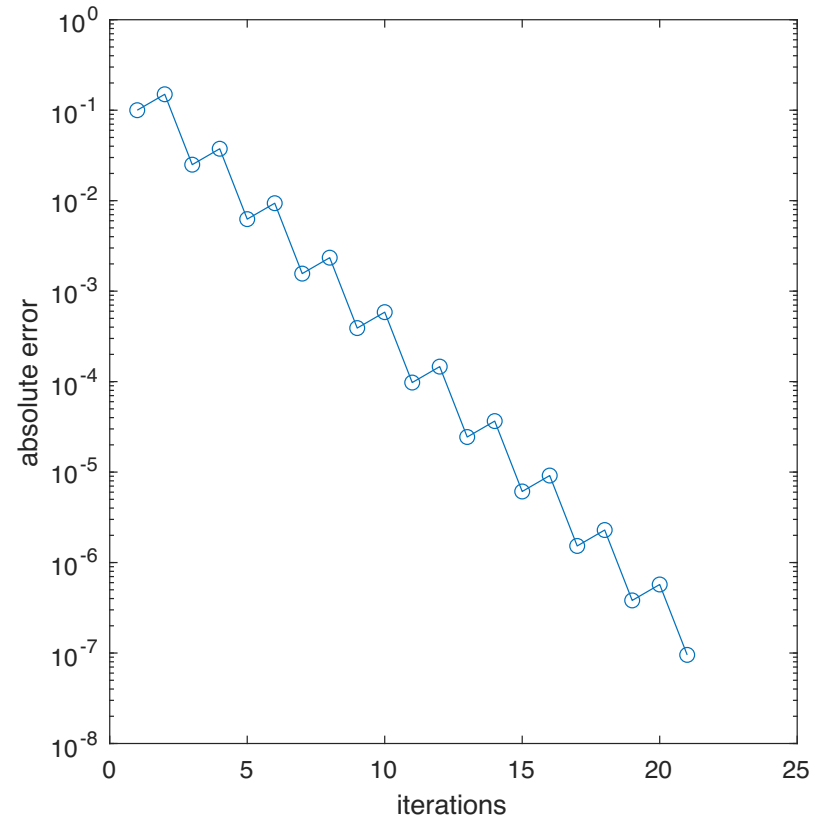
\Rightarrow not linear convergence

Numerical example

Experiment: Solve $f(x) = x^2 - c = 0$ over $[0.5, 1.5]$ with $c = 0.81$ and $x_0 = 1$

```
1: a = 0.5; b = 1.5; c = 0.81; xStar = sqrt(c);
2: f = @(x)x^2 - c;
3: x = (a + b)/2;
4:
5: res = [x, xStar];
6: for k=1:20
7:     if(f(a)*f(x) < 0)
8:         b = x;
9:     else
10:        a = x;
11:    end
12:    x = (a + b)/2;
13:    res(end + 1, :) = [x, xStar];
14: end
```

1:	1.0000e+00	9.0000e-01
2:	7.5000e-01	9.0000e-01
3:	8.7500e-01	9.0000e-01
4:	9.3750e-01	9.0000e-01
5:	9.0625e-01	9.0000e-01
6:	8.9062e-01	9.0000e-01
7:	8.9844e-01	9.0000e-01
8:	9.0234e-01	9.0000e-01
9:	9.0039e-01	9.0000e-01
10:	8.9941e-01	9.0000e-01
11:	8.9990e-01	9.0000e-01
12:	9.0015e-01	9.0000e-01
13:	9.0002e-01	9.0000e-01
14:	8.9996e-01	9.0000e-01
15:	8.9999e-01	9.0000e-01
16:	9.0001e-01	9.0000e-01
17:	9.0000e-01	9.0000e-01



- ▶ Bisection is a slow but sure method.
- ▶ It uses no information about the value of the function or its derivatives - only the sign
- ▶ There are variants that achieve faster convergence \rightsquigarrow textbook by Quarterioni
- ▶ How can we achieve faster convergence in general?

- ▶ Bisection is a slow but sure method.
- ▶ It uses no information about the value of the function or its derivatives - only the sign
- ▶ There are variants that achieve faster convergence \rightsquigarrow textbook by Quarterioni
- ▶ **How can we achieve faster convergence in general?** \rightsquigarrow need to use additional information, at least the function value instead of just the sign

More general formulation via fixed point iterations

Reformulation as fixed point method so that x^* is fixed point

$$x^* = \Phi(x^*)$$

Corresponding iteration: Choose x_0 (initialization) and compute x_1, x_2, \dots from

$$x_{k+1} = \Phi(x_k)$$

We now want to study when this iteration converges to x^* with $f(x^*) = 0$

Convergence of fixed point methods

A mapping $\Phi : [a, b] \rightarrow \mathbb{R}$ is called **contractive** on $[a, b]$ if there is a $0 \leq \Theta < 1$ such that

$$|\Phi(x) - \Phi(y)| \leq \Theta|x - y| \text{ for all } x, y \in [a, b].$$

If Φ is continuously differentiable on $[a, b]$, then

$$\Theta = \sup_{x, y \in [a, b]} \frac{|\Phi(x) - \Phi(y)|}{|x - y|} = \sup_{z \in [a, b]} |\Phi'(z)|$$

Convergence of fixed point methods

Let $\Phi : [a, b] \rightarrow [a, b]$ be contractive with constant $\Theta < 1$. Then:

- ▶ There exists a unique fixed point \bar{x} with $\bar{x} = \Phi(\bar{x})$
- ▶ For any starting guess x_0 in $[a, b]$, the fixed point iteration converges to \bar{x} and

$$|x_{k+1} - x_k| \leq \Theta |x_k - x_{k-1}|$$

and

$$|\bar{x} - x_k| \leq \frac{\Theta^k}{1 - \Theta} |x_1 - x_0|.$$

The second expression allows to estimate the required number of iterations.

↪ board

Convergence of fixed point methods:

For all $x_0 \in I$, it holds for x_0, x_1, x_2, \dots

$$|x_{k+1} - x_k| = |\phi(x_k) - \phi(x_{k-1})| \leq \Theta |x_k - x_{k-1}|$$

because ϕ is contracting with constant Θ .

By induction, we obtain

$$|x_{k+1} - x_k| \leq \Theta^k |x_1 - x_0|$$

Thus



$$|x_{k+m} - x_k| \leq \overbrace{|x_{k+m} - x_{k+m-1}| + \dots + |x_{k+1} - x_k|}$$

$$\leq (\Theta^{k+m-1} + \dots + \Theta^k) |x_1 - x_0|$$

$$\leq \frac{\Theta^k}{1-\Theta} |x_1 - x_0|$$

because $0 \leq \Theta < 1$

$$(1 + \Theta + \dots + \Theta^{m-1}) \leq \sum_{j=0}^{\infty} \Theta^j = \frac{1}{1-\Theta}$$

\Rightarrow Cauchy sequence \Rightarrow convergence

Show $\lim_{n \rightarrow \infty} x_n = x^*$ that is fixed point of ϕ :

$$\begin{aligned} |x^* - \phi(x^*)| &= |x^* - x_{n+1} + x_{n+1} - \phi(x^*)| \\ &= |x^* - x_{n+1} + \phi(x_n) - \phi(x^*)| \\ &\leq |x^* - x_{n+1}| + |\phi(x_n) - \phi(x^*)| \\ &\leq \underbrace{|x^* - x_{n+1}|}_{\xrightarrow{n \rightarrow \infty} 0} + \theta \underbrace{|x^* - x_n|}_{\xrightarrow{n \rightarrow \infty} 0} \end{aligned}$$

$$|x^* - \phi(x^*)| \leq 0$$

$$|x^* - \phi(x^*)| = 0$$

\Rightarrow fixed point

Uniqueness:

If there are two fixed points x^* and y^* then

$$0 \leq |x^* - y^*| = |\phi(x^*) - \phi(y^*)|$$

$$\leq \theta |x^* - y^*|$$

Because $\theta < 1$, this is possible only if

$$|x^* - y^*| = 0$$

\Rightarrow uniqueness

Newton's method

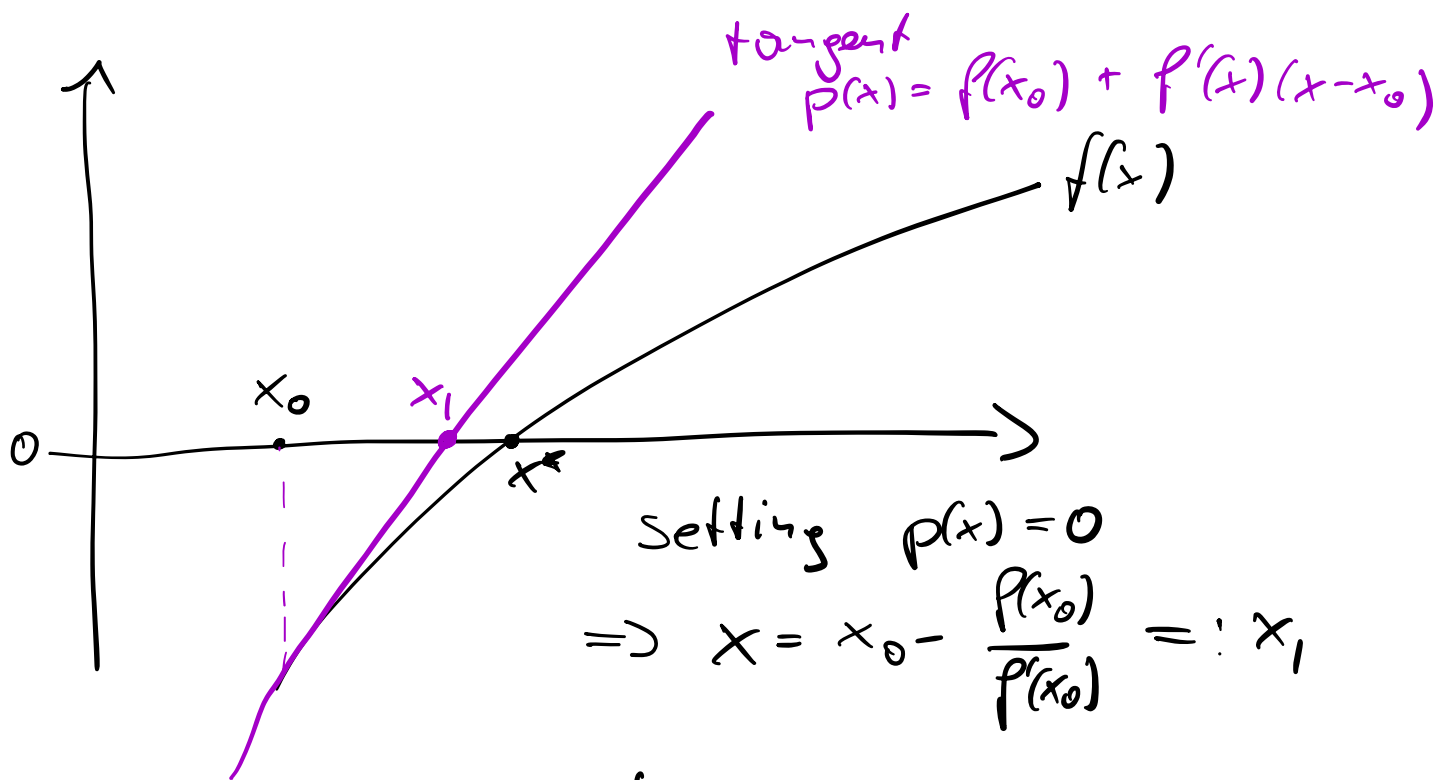
What is the standard approach in numerics when we encounter a nonlinear problem?

Newton's method

What is the standard approach in numerics when we encounter a nonlinear problem?

↪ we linearize

↪ board



fixed point iteration

$$\phi(x) = x - \frac{f(x)}{f'(x)}$$

\Rightarrow need that f is differentiable

$\Rightarrow f'$ does not vanish

Numerical Methods I

MATH-GA 2010.001/CSCI-GA 2420.001

Benjamin Peherstorfer
Courant Institute, NYU

Based on slides by G. Stadler and A. Donev

Today

Last time

- ▶ Interpolation
- ▶ Solving systems of nonlinear equations
- ▶ Bisection methods

Today

- ▶ Newton method

Announcements

- ▶ Homework 5 posted and is due Mon, Nov 18 before class

Recap: Solving nonlinear equations (“root finding”)

We want to solve the nonlinear equation

$$f(x) = 0, \quad x \in \mathbb{R}.$$

We could also have $n < \infty$ equations in n unknowns with $f : \mathbb{R}^n \rightarrow \mathbb{R}^n$

$$f(\mathbf{x}) = \mathbf{0}$$

In general, we will need an iterative approach that constructs x_1, x_2, x_3, \dots such that

$$\lim_{k \rightarrow \infty} x_k = x^*,$$

with $f(x^*) = 0$.

What are important properties of a method for solving nonlinear equations?

- ▶ Does it converge? From which starting point x_0 ?
- ▶ How quickly does it converge?
- ▶ How expensive is each step?

Recap: Bisection method

The bisection method exploits that given a continuous function $f : [a, b] \rightarrow \mathbb{R}$, such that $f(a)f(b) < 0$, there exists $x^* \in (a, b)$ with $f(x^*) = 0$

- ▶ Assumption: f is continuous over $[a, b]$ (very weak assumption!)
- ▶ We have chosen a reasonable interval $[a, b]$ so that there exists a solution $x^* \in (a, b)$ with $f(x^*) = 0$

Recap: Bisection method

The bisection method exploits that given a continuous function $f : [a, b] \rightarrow \mathbb{R}$, such that $f(a)f(b) < 0$, there exists $x^* \in (a, b)$ with $f(x^*) = 0$

- ▶ Assumption: f is continuous over $[a, b]$ (very weak assumption!)
- ▶ We have chosen a reasonable interval $[a, b]$ so that there exists a solution $x^* \in (a, b)$ with $f(x^*) = 0$

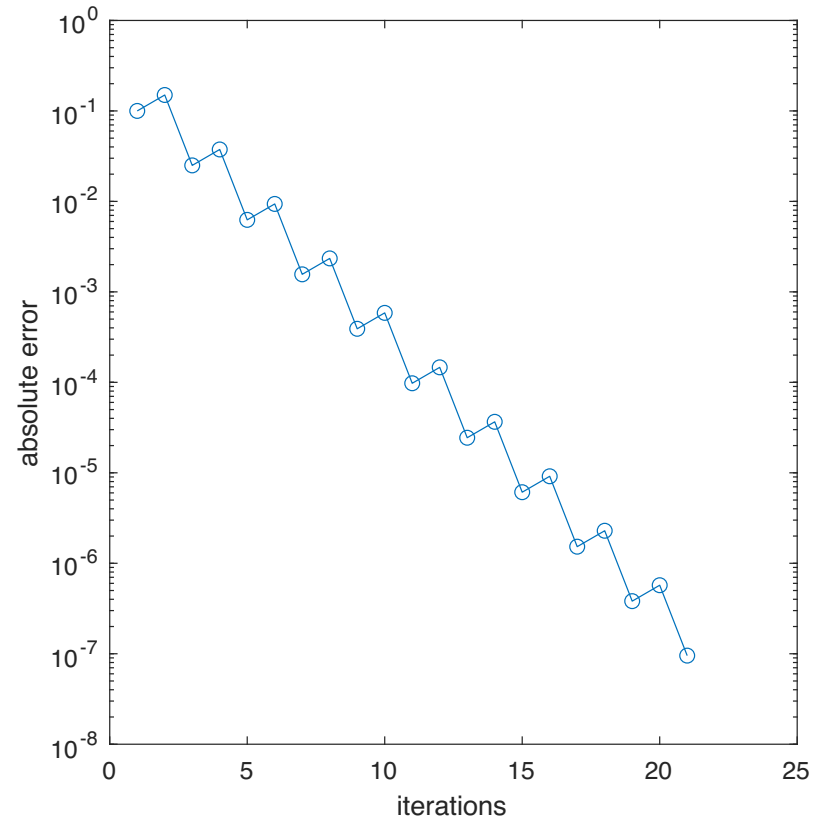
Set $a_0 = a, b_0 = b, x_0 = (a + b)/2$ and iterate for $k = 0, 1, 2, 3, \dots$ as follows:

1. Set $a_{k+1} = a_k, b_{k+1} = x_k$ if $f(x_k)f(a_k) < 0$
2. Set $a_{k+1} = x_k, b_{k+1} = b_k$ if $f(x_k)f(b_k) < 0$
3. Set $x_{k+1} = (a_{k+1} + b_{k+1})/2$
4. Terminate if $|b_{k+1} - a_{k+1}| \leq \epsilon$

Visualization \rightsquigarrow board

Recap

1:	1.0000e+00	9.0000e-01
2:	7.5000e-01	9.0000e-01
3:	8.7500e-01	9.0000e-01
4:	9.3750e-01	9.0000e-01
5:	9.0625e-01	9.0000e-01
6:	8.9062e-01	9.0000e-01
7:	8.9844e-01	9.0000e-01
8:	9.0234e-01	9.0000e-01
9:	9.0039e-01	9.0000e-01
10:	8.9941e-01	9.0000e-01
11:	8.9990e-01	9.0000e-01
12:	9.0015e-01	9.0000e-01
13:	9.0002e-01	9.0000e-01
14:	8.9996e-01	9.0000e-01
15:	8.9999e-01	9.0000e-01
16:	9.0001e-01	9.0000e-01
17:	9.0000e-01	9.0000e-01

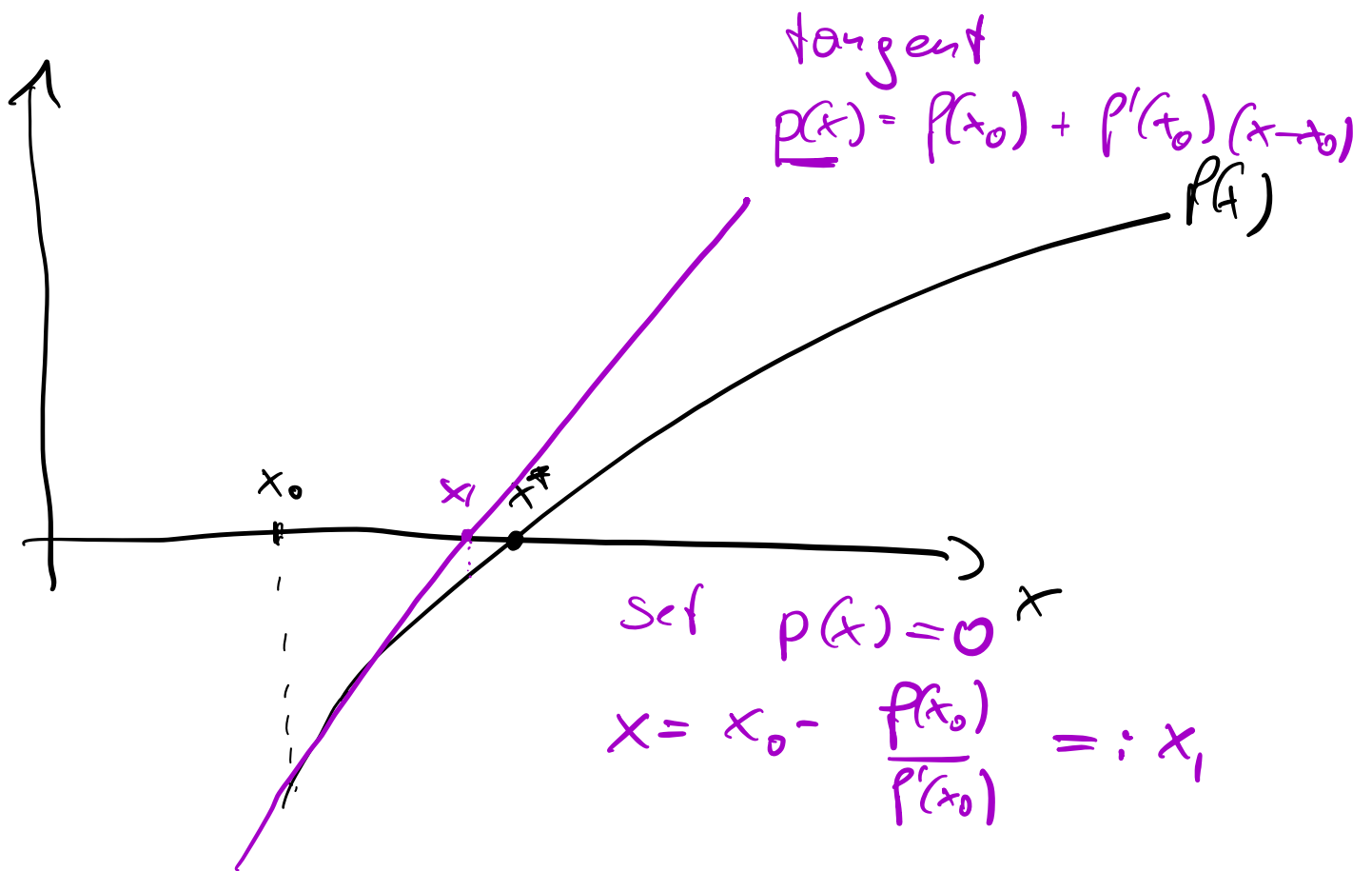


Recap: Newton's method

What is the standard approach in numerics when we encounter a nonlinear problem?

↪ we linearize

↪ board



fixed point iteration

$$\phi(x) = x - \frac{f(x)}{f'(x)}$$

- > need differentiable f
- > not vanishing $f'(x)$

In one dimension, solve $f(x) = 0$:

Start with x_0 , and compute x_1, x_2, \dots from

$$x_{k+1} = x_k - \frac{f(x_k)}{f'(x_k)}, \quad k = 0, 1, \dots$$

Requires $f'(x_k) \neq 0$ to be well-defined (i.e., tangent has nonzero slope).

In one dimension, solve $f(x) = 0$:

Start with x_0 , and compute x_1, x_2, \dots from

$$x_{k+1} = x_k - \frac{f(x_k)}{f'(x_k)}, \quad k = 0, 1, \dots$$

Requires $f'(x_k) \neq 0$ to be well-defined (i.e., tangent has nonzero slope).

Experiment: Solve $f(x) = x^2 - c = 0$ with $c = 0.81$ and $x_0 = 1$

$$\phi(x) = x - \frac{f(x)}{f'(x)} = x - \frac{x^2 - c}{2x} = x - \frac{x}{2} + \frac{c}{2x} = \frac{1}{2} \left(x + \frac{c}{x} \right)$$

Iterations

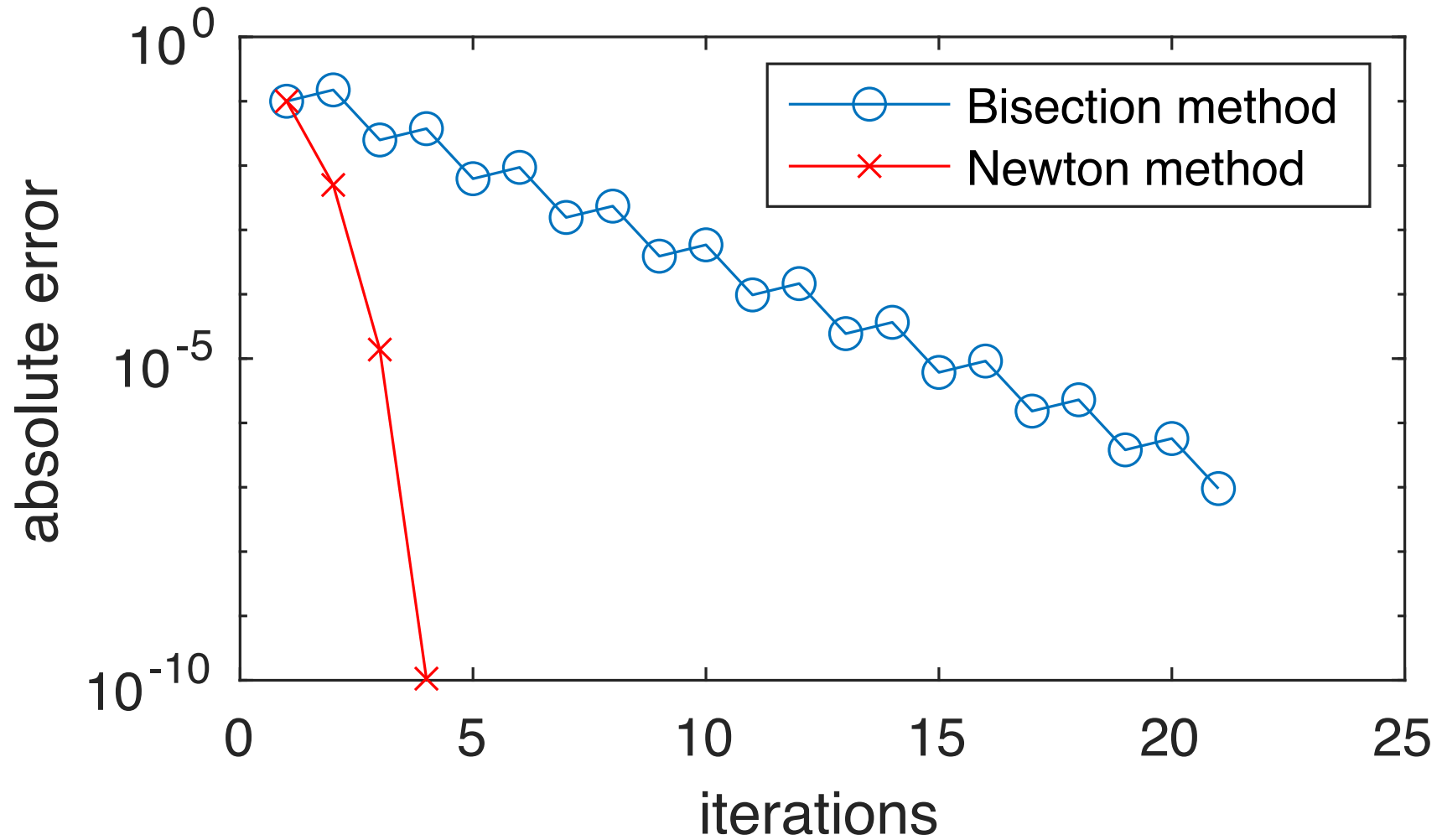
$$x_{k+1} = \frac{1}{2} \left(x_k + \frac{c}{x_k} \right)$$

```
1: format longE
2: c = 0.81;
3: xStar = sqrt(c);
4: x = 1;
5: res = [x, xStar];
6: for i=1:4
7:     x = 0.5*(x + c/x);
8:     res(end + 1, :) = [x, xStar];
9: end
10: res
```

```
1: format longE
2: c = 0.81;
3: xStar = sqrt(c);
4: x = 1;
5: res = [x, xStar];
6: for i=1:4
7:     x = 0.5*(x + c/x);
8:     res(end + 1, :) = [x, xStar];
9: end
10: res
```

```
1: res =
2:     1.0000000000000000e+00     9.0000000000000000e-01
3:     9.0500000000000000e-01     9.0000000000000000e-01
4:     9.000138121546961e-01     9.0000000000000000e-01
5:     9.000000001059849e-01     9.0000000000000000e-01
6:     9.0000000000000000e-01     9.0000000000000000e-01
```

↪ very quick convergence; certainly faster than linear convergence



Newton's method

Let $F : \mathbb{R}^n \rightarrow \mathbb{R}^n$, $n \geq 1$ and solve

$$F(\mathbf{x}) = 0.$$

Truncated Taylor expansion of F about starting point \mathbf{x}^0 :

$$F(\mathbf{x}) \approx F(\mathbf{x}^0) + F'(\mathbf{x}^0)(\mathbf{x} - \mathbf{x}^0).$$

Hence:

$$\mathbf{x}^1 = \mathbf{x}^0 - F'(\mathbf{x}^0)^{-1}F(\mathbf{x}^0)$$

Newton iteration: Start with $\mathbf{x}^0 \in \mathbb{R}^n$, and for $k = 0, 1, \dots$ compute

$$F'(\mathbf{x}^k)\Delta\mathbf{x}^k = -F(\mathbf{x}^k), \quad \mathbf{x}^{k+1} = \mathbf{x}^k + \Delta\mathbf{x}^k$$

Requires that $F'(\mathbf{x}^k) \in \mathbb{R}^{n \times n}$ is invertible.

Newton's method

Newton iteration: Start with $\mathbf{x}^0 \in \mathbb{R}^n$, and for $k = 0, 1, \dots$ compute

$$F'(\mathbf{x}^k)\Delta\mathbf{x}^k = -F(\mathbf{x}^k), \quad \mathbf{x}^{k+1} = \mathbf{x}^k + \Delta\mathbf{x}^k$$

Equivalently:

$$\mathbf{x}^{k+1} = \mathbf{x}^k - F'(\mathbf{x}^k)^{-1}F(\mathbf{x}^k)$$

Newton's method is **affine invariant**, that is, the sequence is invariant to affine transformations \rightsquigarrow **board**

Newton is affine invariant:

Solving $F(x) = 0$ is equivalent to solving

$$AF(x) = 0$$

for regular $A \in \mathbb{R}^{n \times n}$

Set $G(x) = AF(x)$, then apply Newton to G

$$\begin{aligned} y_{k+1} &= y_k - G'(y_k)^{-1} G(y_k) \\ &= y_k - (AF'(y_k))^{-1} (AF(y_k)) \\ &= y_k - F'(y_k)^{-1} A^{-1} A F(y_k) \\ &= y_k - F'(y_k)^{-1} F(y_k) \end{aligned}$$

\Rightarrow iterates are the same even if we transform $F \rightarrow G$ via regular matrix

\Rightarrow not affected by bad scaling



vs



Convergence of Newton's method

Assumptions on F : $D \subset \mathbb{R}^n$ open and convex, $F : D \rightarrow \mathbb{R}^n$ continuously differentiable with $F'(\mathbf{x})$ invertible for all \mathbf{x} , and there exists $\omega \geq 0$ such that

$$\|F'(\mathbf{x})^{-1}(F'(\mathbf{x} + s\mathbf{v}) - F'(\mathbf{x}))\mathbf{v}\| \leq s\omega\|\mathbf{v}\|^2$$

for all $s \in [0, 1]$, $\mathbf{x} \in D$, $\mathbf{v} \in \mathbb{R}^n$ with $\mathbf{x} + \mathbf{v} \in D$.

Assumptions on \mathbf{x}^* and \mathbf{x}^0 : There exists a solution $\mathbf{x}^* \in D$ and a starting point $\mathbf{x}^0 \in D$ such that

$$\rho := \|\mathbf{x}^* - \mathbf{x}^0\| \leq \frac{2}{\omega} \text{ and } B_\rho(\mathbf{x}^*) \subset D$$

where

$$B_\rho(\mathbf{x}^*) = \{\mathbf{y} \in \mathbb{R}^n \mid \|\mathbf{y} - \mathbf{x}^*\| < \rho\}$$

Q: Meaning of ω ?

Theorem: Under the assumptions of the previous slide, the Newton sequence \mathbf{x}^k stays in $B_\rho(\mathbf{x}^*)$ and $\lim_{k \rightarrow \infty} \mathbf{x}^k = \mathbf{x}^*$, and

$$\|\mathbf{x}^{k+1} - \mathbf{x}^*\| \leq \frac{\omega}{2} \|\mathbf{x}^k - \mathbf{x}^*\|^2$$

Moreover, the solution x^* is unique in $B_{2/\omega}(x^*)$.

Proof: \rightsquigarrow board

Consider the Lipschitz condition

$$\bullet \|F'(x)^{-1} (F'(x+sv) - F'(x))v\|_2 \leq sw \|v\|_2^2$$
$$\forall s \in [0,1], x \in D, v \in \mathbb{R}^n \quad x+v \in D$$

and

$$\bullet \int_0^1 [F'(x + s(\gamma-x)) - F'(x)] (\gamma-x) ds = \dots =$$
$$= F(\gamma) - F(x) - F'(x)(\gamma-x)$$

Now take $v = \gamma - x$ with $\gamma \in D$, then

$$\int_0^1 \|F'(x)^{-1} [F'(x + s(\gamma-x)) - F'(x)] (\gamma-x)\| ds$$
$$\leq \int_0^1 sw \|\gamma-x\|^2 ds$$

$$\|F'(x)^{-1} [F(\gamma) - F(x) - F'(x)(\gamma-x)]\|$$
$$= \left\| \int_0^1 F'(x)^{-1} [F'(x + s(\gamma-x)) - F'(x)] (\gamma-x) ds \right\|$$
$$\leq \int_0^1 \|F'(x)^{-1} [F'(x + s(\gamma-x)) - F'(x)] (\gamma-x)\| ds$$
$$\leq \int_0^1 sw \|\gamma-x\|^2 ds$$
$$= w \frac{s^2}{2} \|\gamma-x\|^2 \Big|_0^1 = \frac{w}{2} \|\gamma-x\|^2 \quad (*)$$

Now convergence of Newton

$$\begin{aligned}x^{n+1} - x^* &= x^n - F'(x^n)^{-1} (F(x^n) - F(x^*)) - x^* \\ &= x^n - x^* - F'(x^n)^{-1} (F(x^n) - F(x^*))\end{aligned}$$

$$= F'(x^n)^{-1} [-F(x^n) + F(x^*) - F'(x^n)(x^* - x^n)]$$

with (*) follows

$$\|x^{n+1} - x^*\| \leq \frac{L}{2} \|x^* - x^n\|^2 \Rightarrow \text{quadratic rate}$$

Now show that sequence (x^n) remains in ball $B_p(x^*)$

Have:

$$0 < \|x^n - x^*\| \leq \|x^* - x^0\| =: p$$

$$\begin{aligned}\|x^{n+1} - x^*\| &\leq \underbrace{\frac{L}{2} \|x^* - x^n\|}_{\leq p} \|x^* - x^n\| \\ &\leq \underline{p \frac{L}{2}}\end{aligned}$$

Now we need assumption about starting point
We have that $x^0 \in D$ satisfies

$$p = \|x^* - x^0\| < \frac{2}{L} \quad \text{and} \quad B_p(x^*) \subseteq D$$

$$\Rightarrow p \frac{L}{2} < 1$$

$$\|x^{n+1} - x^*\| < \|x^n - x^*\| \leq \|x^* - x^0\| = p$$

$\Rightarrow (x^n)$ remains in $B_p(x^*)$

Thus $\{x^*\} \subseteq B_p(x^*) \subseteq D$

important for F' to remain invertible

Uniqueness: the solution x^* is unique

$$\text{in } B_{\frac{\rho}{L}}(x^*)$$

$\hookrightarrow x^{**} \in B_{\frac{\rho}{L}}(x^*)$ is another solution

$$\text{so that } F(x^{**}) = 0 \text{ and } \|x^* - x^{**}\| < \frac{\rho}{L}$$

With (*) we obtain

$$\|x^{**} - x^*\| = \|0 - 0 - F'(x^*)^{-1} F'(x^*) (x^{**} - x^*)\|$$
$$= \|F'(x^*)^{-1} (F(x^{**}) - F(x^*) - F'(x^*) (x^{**} - x^*))\|$$

$$\stackrel{\text{use (*)}}{\leq} \frac{L}{2} \|x^{**} - x^*\| \|x^{**} - x^*\|$$

$$< 1$$

thus only $x^{**} = x^*$

Theorem: Under the assumptions of the previous slide, the Newton sequence \mathbf{x}^k stays in $B_\rho(\mathbf{x}^*)$ and $\lim_{k \rightarrow \infty} \mathbf{x}^k = \mathbf{x}^*$, and

$$\|\mathbf{x}^{k+1} - \mathbf{x}^*\| \leq \frac{\omega}{2} \|\mathbf{x}^k - \mathbf{x}^*\|^2$$

Moreover, the solution x^* is unique in $B_{2/\omega}(x^*)$.

Proof: \rightsquigarrow board

Summary: The Newton method converges locally and quadratically.

Role of initialization

Choice of **initialization** x^0 is critical. Depending on the initialization, the Newton iteration might

- ▶ not converge (it could “blow up” or “oscillate” between two points)
- ▶ converge to different solutions
- ▶ fail because it hits a point where the Jacobian is not invertible (this cannot happen if the conditions of the convergence theorem are satisfied)
- ▶ ...

Sometimes, **continuation ideas** must be used to find good initializations: Solve simpler problems first and use solution as starting point for harder problems.

General comments

The “more nonlinear” a problem, the harder it is to solve.

$$\|F'(\mathbf{x})^{-1}(F'(\mathbf{x} + s\mathbf{v}) - F'(\mathbf{x}))\mathbf{v}\| \leq s\omega\|\mathbf{v}\|^2$$

Very nonlinear $\rightsquigarrow F'(x)$ changes a lot $\rightsquigarrow \omega$ large (need x_0 closer to x^* required)

General comments

The “more nonlinear” a problem, the harder it is to solve.

$$\|F'(\mathbf{x})^{-1}(F'(\mathbf{x} + s\mathbf{v}) - F'(\mathbf{x}))\mathbf{v}\| \leq s\omega\|\mathbf{v}\|^2$$

Very nonlinear $\rightsquigarrow F'(x)$ changes a lot $\rightsquigarrow \omega$ large (need x_0 closer to x^* required)

Computation of Jacobian can be costly/complicated

\rightsquigarrow sometimes approximate $F'(x^k)$

General comments

The “more nonlinear” a problem, the harder it is to solve.

$$\|F'(\mathbf{x})^{-1}(F'(\mathbf{x} + s\mathbf{v}) - F'(\mathbf{x}))\mathbf{v}\| \leq s\omega\|\mathbf{v}\|^2$$

Very nonlinear $\rightsquigarrow F'(x)$ changes a lot $\rightsquigarrow \omega$ large (need x_0 closer to x^* required)

Computation of Jacobian can be costly/complicated

\rightsquigarrow sometimes approximate $F'(x^k)$

There's **no reliable black-box solver** for nonlinear problems; at least for higher-dimensional problems, the structure of the problem must be taken into account.

General comments

The “more nonlinear” a problem, the harder it is to solve.

$$\|F'(\mathbf{x})^{-1}(F'(\mathbf{x} + s\mathbf{v}) - F'(\mathbf{x}))\mathbf{v}\| \leq s\omega\|\mathbf{v}\|^2$$

Very nonlinear $\rightsquigarrow F'(x)$ changes a lot $\rightsquigarrow \omega$ large (need x_0 closer to x^* required)

Computation of Jacobian can be costly/complicated

\rightsquigarrow sometimes approximate $F'(x^k)$

There's **no reliable black-box solver** for nonlinear problems; at least for higher-dimensional problems, the structure of the problem must be taken into account.

“Classification of mathematical problems as linear and nonlinear is like classification of the Universe as bananas and non-bananas.”

Nonlinear least squares—Gauss-Newton

Nonlinear least-squares problems

Assume a least squares problem, where the parameters \mathbf{x} do *not* enter linearly into the model.

Instead of

$$\min_{\mathbf{x} \in \mathbb{R}^n} \|A\mathbf{x} - \mathbf{b}\|^2,$$

we have with $F : D \rightarrow \mathbb{R}^m$, $m \geq n$, $D \subset \mathbb{R}^n$:

$$\min_{\mathbf{x} \in \mathbb{R}^n} g(\mathbf{x}) := \frac{1}{2} \|F(\mathbf{x})\|^2, \quad \text{where } F(\mathbf{x})_i = \varphi(t_i, \mathbf{x}) - b_i, 1 \leq i \leq m$$

A (local) minimum \mathbf{x}^* of this optimization problem satisfies:

$$g'(\mathbf{x}) = 0, \quad g''(\mathbf{x}) \text{ is positive definite.}$$

Nonlinear least-squares problems

The **derivative** of $g(\cdot)$ is

$$G(\mathbf{x}) := g'(\mathbf{x}) = F'(\mathbf{x})F(\mathbf{x})$$

Setting $G(\mathbf{x}) = 0$ gives a nonlinear system in \mathbf{x} , $G : D \rightarrow \mathbb{R}^m$.

Let's try to solve it $G(\mathbf{x}) = 0$ using Newton's method:

$$G'(\mathbf{x}^k)\Delta\mathbf{x}^k = -G(\mathbf{x}^k), \quad \mathbf{x}^{k+1} = \mathbf{x}^k + \Delta\mathbf{x}^k$$

where

$$G'(\mathbf{x}) = F'(\mathbf{x})^T F'(\mathbf{x}) + F''(\mathbf{x})^T F(\mathbf{x}) \quad \text{Hessian of } g \text{ (objective)}$$

\rightsquigarrow second-order information of F enters through $F''(X)^T$

Q: What are you observing about how second-order information enters?

Nonlinear least-squares problems

The **derivative** of $g(\cdot)$ is

$$G(\mathbf{x}) := g'(\mathbf{x}) = F'(\mathbf{x})F(\mathbf{x})$$

Setting $G(\mathbf{x}) = 0$ gives a nonlinear system in \mathbf{x} , $G : D \rightarrow \mathbb{R}^m$.

Let's try to solve it $G(\mathbf{x}) = 0$ using Newton's method:

$$G'(\mathbf{x}^k)\Delta\mathbf{x}^k = -G(\mathbf{x}^k), \quad \mathbf{x}^{k+1} = \mathbf{x}^k + \Delta\mathbf{x}^k$$

where

$$G'(\mathbf{x}) = F'(\mathbf{x})^T F'(\mathbf{x}) + F''(\mathbf{x})^T F(\mathbf{x}) \quad \text{Hessian of } g \text{ (objective)}$$

\rightsquigarrow second-order information of F enters through $F''(X)^T$

Q: What are you observing about how second-order information enters? \rightsquigarrow multiplied with $F(x)$

Nonlinear least-squares problems

If the data is compatible with the model, which means that the model can perfectly fit the data with zero training error, then $F(\mathbf{x}^*) = 0$

Then, term involving $F''(\mathbf{x}^*)$ drops out in $G'(\mathbf{x}^*)$ anyway as we move towards \mathbf{x}^* .

If $\|F(\mathbf{x}^*)\|$ is small, and thus data *almost* compatible with model, then neglecting that term might not make the convergence much slower.

Also, it's expensive to compute $F''(\mathbf{x})$

Nonlinear least-squares problems—Gauss-Newton

The resulting Newton method for the nonlinear least squares problem is called **Gauss-Newton method**: Initialize \mathbf{x}^0 and for $k = 0, 1, \dots$ solve

$$F'(\mathbf{x}^k)^T F'(\mathbf{x}^k) \Delta \mathbf{x}^k = -F'(\mathbf{x}^k)^T F(\mathbf{x}^k) \quad (\text{solve}) \quad (1)$$

$$\mathbf{x}^{k+1} = \mathbf{x}^k + \Delta \mathbf{x}^k. \quad (\text{update step})$$

Nonlinear least-squares problems—Gauss-Newton

The resulting Newton method for the nonlinear least squares problem is called **Gauss-Newton method**: Initialize \mathbf{x}^0 and for $k = 0, 1, \dots$ solve

$$F'(\mathbf{x}^k)^T F'(\mathbf{x}^k) \Delta \mathbf{x}^k = -F'(\mathbf{x}^k)^T F(\mathbf{x}^k) \quad (\text{solve}) \quad (1)$$

$$\mathbf{x}^{k+1} = \mathbf{x}^k + \Delta \mathbf{x}^k. \quad (\text{update step})$$

Q: Should we implement GN like this?

Nonlinear least-squares problems—Gauss-Newton

The resulting Newton method for the nonlinear least squares problem is called **Gauss-Newton method**: Initialize \mathbf{x}^0 and for $k = 0, 1, \dots$ solve

$$F'(\mathbf{x}^k)^T F'(\mathbf{x}^k) \Delta \mathbf{x}^k = -F'(\mathbf{x}^k)^T F(\mathbf{x}^k) \quad (\text{solve}) \quad (1)$$

$$\mathbf{x}^{k+1} = \mathbf{x}^k + \Delta \mathbf{x}^k. \quad (\text{update step})$$

Q: Should we implement GN like this?

The solve step is the normal equation for the linear least squares problem

$$\min_{\Delta \mathbf{x}} \|F'(\mathbf{x}^k) \Delta \mathbf{x}^k + F(\mathbf{x}^k)\|. \quad (2)$$

so we better solve (2) rather than directly (1)

Convergence of Gauss-Newton method

Assumptions on F : $D \subset \mathbb{R}^n$ open and convex, $F : D \rightarrow \mathbb{R}^m$, $m \geq n$ continuously differentiable with $F'(\mathbf{x})$ has full rank for all \mathbf{x} , and let $\omega \geq 0$, $0 \leq \kappa^* < 1$ such that

$$\|F'(\mathbf{x})^+(F'(\mathbf{x} + s\mathbf{v}) - F'(\mathbf{x}))\mathbf{v}\| \leq s\omega\|\mathbf{v}\|^2$$

for all $s \in [0, 1]$, $\mathbf{x} \in D$, $\mathbf{v} \in \mathbb{R}^n$ with $\mathbf{x} + \mathbf{v} \in D$.

Assumptions on \mathbf{x}^* and \mathbf{x}^0 : Assume there exists a solution $\mathbf{x}^* \in D$ of the least squares problem and a starting point $\mathbf{x}^0 \in D$ such that

$$\|F'(\mathbf{x})^+F(\mathbf{x}^*)\| \leq \kappa^*\|\mathbf{x} - \mathbf{x}^*\|$$

$$\rho := \|\mathbf{x}^* - \mathbf{x}^0\| \leq \frac{2(1 - \kappa^*)}{\omega} := \sigma$$

Theorem: Then, the sequence \mathbf{x}^k stays in $B_\rho(\mathbf{x}^*)$ and $\lim_{k \rightarrow \infty} \mathbf{x}^k = \mathbf{x}^*$, and

$$\|\mathbf{x}^{k+1} - \mathbf{x}^*\| \leq \frac{\omega}{2}\|\mathbf{x}^k - \mathbf{x}^*\|^2 + \underbrace{\kappa^*\|\mathbf{x}^k - \mathbf{x}^*\|}_{\rightsquigarrow \text{linear convergence if } \kappa^* > 0!}$$

\rightsquigarrow linear convergence if $\kappa^* > 0!$

\rightsquigarrow we usually want to choose models that are “almost compatible” which means κ^* is often very small

Conclusions

- ▶ Solving nonlinear systems of equations (“root finding”) is iterative in nature in general
- ▶ The order of convergence matters; quadratic is good enough but mind costs per step
- ▶ Newton’s method is second order but requires derivatives/Jacobian evaluations.
- ▶ In higher dimensions, a good initial guess is critical for Newton’s method
- ▶ There are many variants of Newton’s method (e.g., Quasi-Newton methods) that avoid the computational costs of computing the Hessian
- ▶ (Machine learning is using first-order methods only anyway...)

Numerical Methods I

MATH-GA 2010.001/CSCI-GA 2420.001

Benjamin Peherstorfer
Courant Institute, NYU

Based on slides by G. Stadler and A. Donev

Today

Two weeks ago

- ▶ Function approximation
- ▶ Interpolation with polynomials
- ▶ Interpolation beyond polynomials

Today

- ▶ Numerical integration
- ▶ Newton-Cotes

Announcements

- ▶ Homework 5 is posted and due Mon, Nov 18 before class

Numerical quadrature

Numerical integration

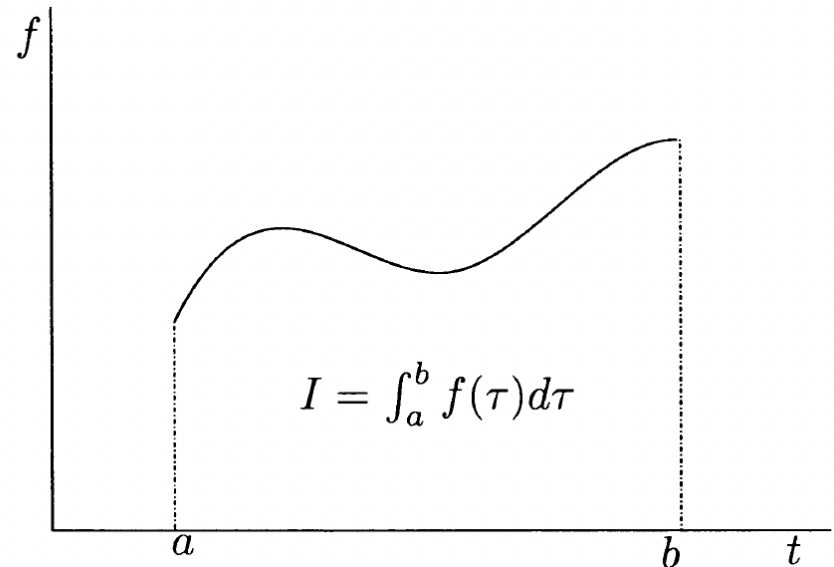
We want to approximate the definite integral

$$I(f) = I_a^b(f) = \int_a^b f(t) dt$$

numerically.

Properties of the integral:

- ▶ I is linear
- ▶ positive, i.e., if f is nonnegative, then $I(f)$ is nonnegative
- ▶ additive w.r.t. the interval bounds:
 $I_a^c = I_a^b + I_b^c$



Condition of numerical integration

Lets study the map

$$([a, b], f) \rightarrow \int_a^b f(t) dt,$$

where we use the L^1 -norm for f

$$\|f\|_1 = \int_a^b |f(t)| dt = I(|f|)$$

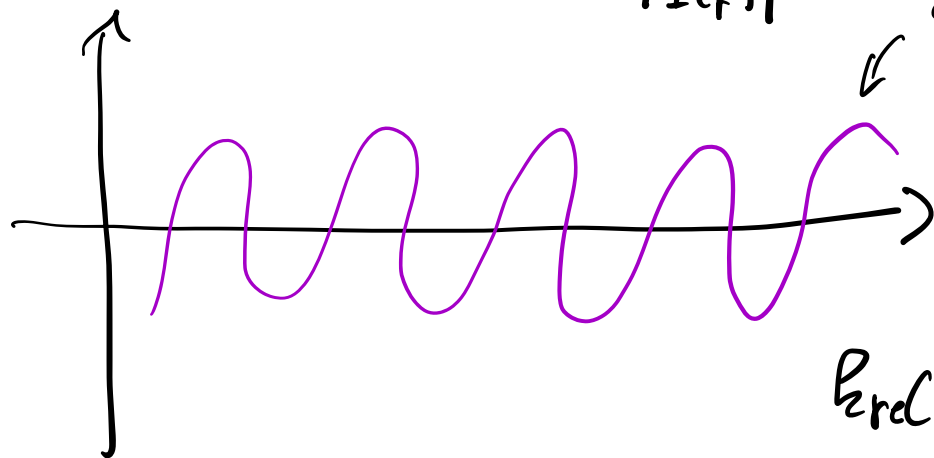
The absolute and relative condition numbers of integration are:

$$\begin{aligned}\kappa_{\text{abs}} &= 1, \\ \kappa_{\text{rel}} &= \frac{I(|f|)}{|I(f)|}.\end{aligned}$$

What does this mean?

Condition number

$$k_{\text{rel}} = \frac{I(f)}{|I(f)|}$$



changes sign after

$$I(f) \approx 0$$

$$|I(f)| \gg 0$$

$$k_{\text{rel}} = \frac{\text{large}}{\text{small}} \gg 1$$

comparable to addition

$$(a, b) \mapsto a + b$$

$$a, b \in \mathbb{R}$$

$$k_{\text{rel}} = \frac{|a| + |b|}{|a + b|}$$

if $a \approx -b$, then poorly conditioned

Condition of numerical integration

Lets study the map

$$([a, b], f) \rightarrow \int_a^b f(t) dt,$$

where we use the L^1 -norm for f

$$\|f\|_1 = \int_a^b |f(t)| dt = I(|f|)$$

The absolute and relative condition numbers of integration are:

$$\begin{aligned}\kappa_{\text{abs}} &= 1, \\ \kappa_{\text{rel}} &= \frac{I(|f|)}{|I(f)|}.\end{aligned}$$

What does this mean? Integration is harmless w.r.t. the absolute condition number, and problematic w.r.t. the relative condition number if $I(f)$ is small and f changes sign
 \rightsquigarrow **board**

Numerical integration

We are looking for a map

$$\hat{I} : \begin{cases} C([a, b]) & \rightarrow \mathbb{R} \\ f & \mapsto \hat{I}(f) \end{cases}$$

such that the integration error $|I(f) - \hat{I}(f)|$ is small.

What ideas come to your mind?

Numerical integration

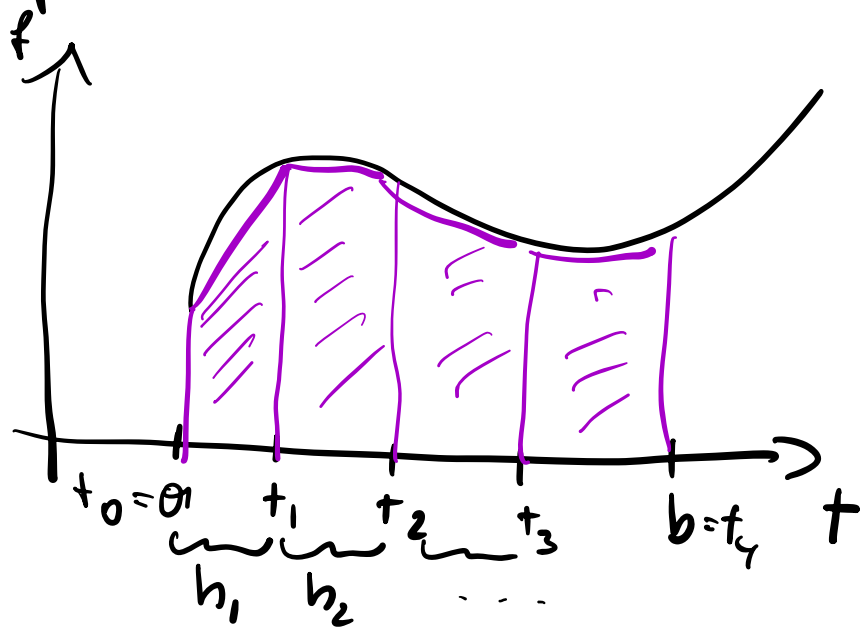
We are looking for a map

$$\hat{I} : \begin{cases} C([a, b]) & \rightarrow \mathbb{R} \\ f & \mapsto \hat{I}(f) \end{cases}$$

such that the integration error $|I(f) - \hat{I}(f)|$ is small.

What ideas come to your mind? Example: Trapezoidal sum \rightsquigarrow board

Trapezoidal sum



$$T^{(n)} = \sum_{i=1}^n T_i, \quad T_i = \frac{h_i}{2} (f(t_{i-1}) + f(t_i))$$

Lower and upper bound sum:

$$R_{\min}^{(n)} = \sum_{i=1}^n h_i \min_{t \in [t_{i-1}, t_i]} f(t)$$

$$R_{\max}^{(n)} = \sum_{i=1}^n h_i \max_{t \in [t_{i-1}, t_i]} f(t)$$

$$R_{\min}^{(n)} \leq T^{(n)} \leq R_{\max}^{(n)}$$

$$\downarrow$$

$$I(f)$$

$$\downarrow$$

$$I(f)$$

for continuous f
 $n \rightarrow \infty$

$$T^{(n)} \xrightarrow{n \rightarrow \infty} I(f)$$

Numerical integration

We are looking for a map

$$\hat{I} : \begin{cases} C([a, b]) & \rightarrow \mathbb{R} \\ f & \mapsto \hat{I}(f) \end{cases}$$

such that the integration error $|I(f) - \hat{I}(f)|$ is small.

What ideas come to your mind? Example: Trapezoidal sum \rightsquigarrow board

General quadrature formula:

$$\hat{I}(f) = \sum_{i=0}^n \lambda_i f(t_i),$$

with **weights** λ_i and **nodal points** t_i , $i = 0, 1, \dots, n$.

Newton-Cotes formulas

Trapezoidal rule replaces f by easy-to-integrate piecewise linear approximation \hat{f} .
What other easy-to-integrate approximations could we use?

Newton-Cotes formulas

Trapezoidal rule replaces f by easy-to-integrate piecewise linear approximation \hat{f} .
What other easy-to-integrate approximations could we use?

Polynomials!

Quick recap on polynomial interpolation

Recap: Interpolation

Consider n pairs of data samples $(x_i, y_i), i = 1, \dots, n$ with

$$y_i = f(x_i)$$

Based on $\{(x_i, y_i)\}_{i=1}^n$, we now would like to find an approximation $\tilde{f} \in \mathcal{V}_n$ that is “close” to f .

For example, we could enforce the interpolation condition, namely that it holds

$$\tilde{f}(x_i) = f(x_i), \quad i = 1, \dots, n$$

We could also use regression ($m > n$) and minimize, e.g.,

$$\frac{1}{m} \sum_{i=1}^m |y_i - \tilde{f}(x_i)|^2$$

Recap: Orthogonal polynomials

Define an **inner product between functions**:

$$(f, g) = \int_a^b \omega(x) f(x) g(x) dx,$$

where $\omega(x) > 0$ for $a \leq x \leq b$ is a **weight function**. The **induced norm** is $\|f\| := \sqrt{(f, f)}$.

Let $P_0, P_1, P_2, \dots, P_K$ be polynomials of $0, 1, 2, \dots, K$ order, respectively. They are called orthogonal polynomials on $[a, b]$ with respect to the weight function $\omega(x)$ if it holds

$$(P_i, P_j) = \int_a^b \omega(x) P_i(x) P_j(x) dx = \delta_{ij} \gamma_i, \quad i, j = 0, \dots, K,$$

with $\gamma_i = \|P_i\|^2 > 0$.

Recap

To define orthogonal polynomials uniquely, we normalize them so that the leading coefficient is one, i.e.,

$$P_k(x) = x^k + \dots$$

Theorem: There exist uniquely determined orthogonal polynomials $P_k \in \mathbb{P}_k$ with leading coefficient 1. These polynomials satisfy the 3-term recurrence relation:

$$P_k(x) = (x + a_k)P_{k-1}(x) + b_k P_{k-2}(x), \quad k = 2, 3, \dots$$

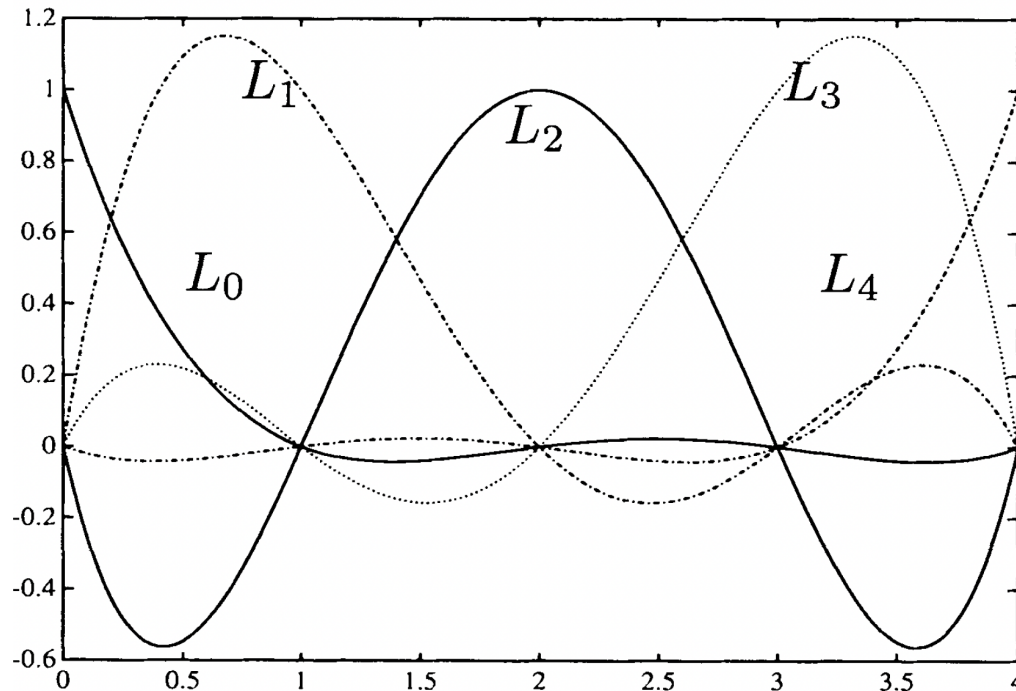
with starting values $P_0 = 1$, $P_1 = x + a_1$, where

$$a_k = -\frac{(xP_{k-1}, P_{k-1})}{(P_{k-1}, P_{k-1})}, \quad b_k = -\frac{(P_{k-1}, P_{k-1})}{(P_{k-2}, P_{k-2})}$$

Recap: Lagrange basis

The Lagrange polynomials $L_0, \dots, L_n \in \mathbb{P}_n$ are uniquely defined for distinct x_0, \dots, x_n

$$L_i(x_j) = \delta_{ij}, \quad L_i \in \mathbb{P}_n.$$



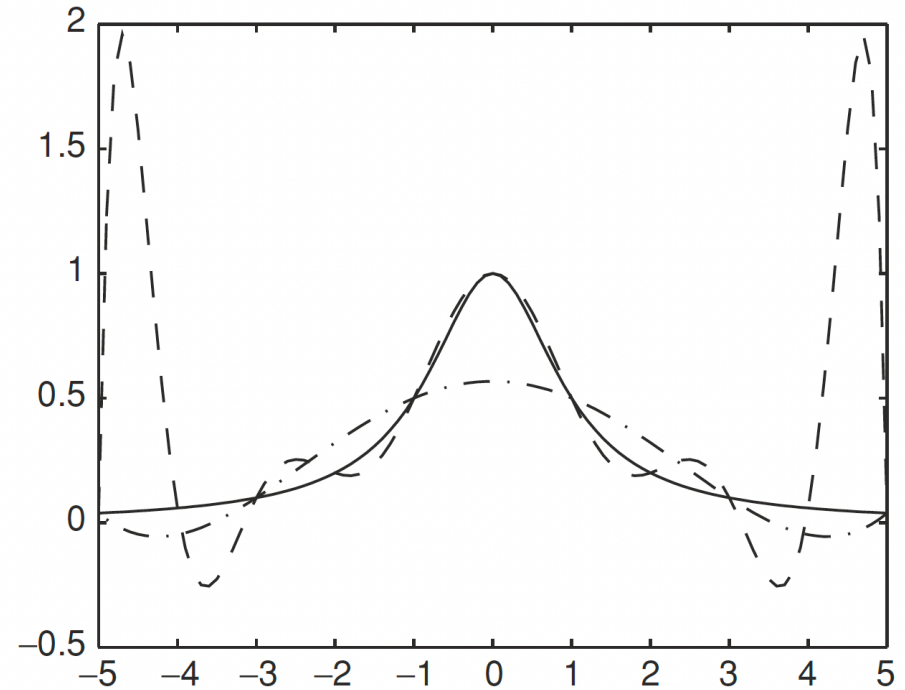
Lagrange polynomials up to order $n = 4$ for equidistant x_0, \dots, x_4 . [Figure: Deuffhard]

Recap: Runge's phenomenon

Let's approximate

$$f(x) = \frac{1}{1+x^2}, \quad -5 \leq x \leq 5,$$

using polynomial interpolation on equally spaced nodes in $[-5, 5]$.



[Figure: Quarteroni]

Interpolants of degree $n = 5$ and $n = 10$ of $f(x) = 1/(1+x^2)$ on equidistant nodes. It can be shown that polynomial interpolation does not converge for $|x| > 4$ for this f . It is a very common situation that interpolation on equidistant nodes leads to high oscillations near the interval ends \rightsquigarrow **Runge's phenomenon**

Recap: Newton polynomial basis

The leading coefficient a_n of the interpolation polynomial (monomial basis!)

$$P_f(x|x_0, \dots, x_n) = a_n x^n + \dots + a_0$$

is called the n -th divided difference, $[x_0, \dots, x_n]f := a_n$.

The divided differences are the coefficients c_0, \dots, c_n : The interpolation polynomial $P_f(\cdot|x_0, \dots, x_n)$ for $x_0 \leq x_1 \leq \dots \leq x_n$ (not necessarily distinct and thus need $f \in C^{n+1}$) is given by

$$P(x) = \sum_{i=0}^n [x_0, \dots, x_i]f \omega_i(x).$$

Furthermore,

$$f(x) = P(x) + [x_0, \dots, x_n, x]f \omega_{n+1}(x).$$

Back to quadrature

Newton-Cotes formulas continued

Given **fixed** nodes t_0, \dots, t_n , use polynomial approximation

$$\hat{f} = P_f(t|t_0, \dots, t_n) = \sum_{i=0}^n f(t_i)L_{in}(t)$$

with Lagrange polynomials L_{0n}, \dots, L_{nn}

Thus:

$$\hat{I}(f) = (b - a) \sum_{i=0}^n \lambda_{in} f(t_i),$$

where

$$\lambda_{in} = \frac{1}{b - a} \int_a^b L_{in}(t) dt$$

Newton-Cotes formulas (cont'd)

Quadrature formulas defined in this way are exact for polynomials $P \in \mathbb{P}_n$ of degree less than or equal to n

$$\hat{I}(P) = I(P_n(P)) = I(P), \quad \text{for all } P \in \mathbb{P}_n$$

Theorem: For $(n + 1)$ pairwise distinct nodes t_0, \dots, t_n , there exists exactly one quadrature formula (i.e., unique weights $\lambda_0, \dots, \lambda_n$)

$$\hat{I}(f) = (b - a) \sum_{i=0}^n \lambda_i f(t_i),$$

that is exact for all $p \in \mathbb{P}_n$.

Proof \rightsquigarrow **board**

$n+1$ pairwise distinct nodes t_0, \dots, t_n , then there exists one and only one quadrature rule

$$\hat{I}(f) = (b-a) \sum_{i=0}^n \lambda_i f(t_i)$$

which is exact for all polynomials $P \in \mathbb{P}_n$.

↳ have \hat{I} that is exact for $P \in \mathbb{P}_n$

insert Lagrange polynomials for t_0, \dots, t_n

$$L_{ij} \in \mathbb{P}_n$$

$$\begin{aligned} I(L_{ij}) &= \hat{I}(L_{ij}) = (b-a) \sum_{j=0}^n \lambda_j \underbrace{L_{ij}(t_j)}_{\delta_{ij}} \\ \text{exact!} \nearrow &= (b-a) \lambda_j \end{aligned}$$

$$\lambda_j = \frac{1}{(b-a)} I(L_{ij})$$

As long as $I(L_{ij}) = \hat{I}(L_{ij})$ we get the same weights.

\Rightarrow unique

Equidistantly spaced nodes

$$h_i = h = \frac{b - a}{n}, \quad t_i = a + ih, \quad i = 0, \dots, n$$

then quadrature formulas are called the *Newton-Cotes formulas* with weights

$$\lambda_{in} = \frac{1}{b - a} \int_a^b \prod_{i \neq j} \frac{t - t_i}{t_i - t_j} dt = \frac{1}{n} \int_0^n \prod_{i \neq j} \frac{s - j}{i - j} ds$$

These weights are independent of the interval boundaries a and b and can be pre-computed once and for all:

Table 9.1. Newton-Cotes weights λ_{in} for $n = 1, \dots, 4$.

n	$\lambda_{0n}, \dots, \lambda_{nn}$	Error	Name
1	$\frac{1}{2} \quad \frac{1}{2}$	$\frac{h^3}{12} f''(\tau)$	Trapezoidal rule
2	$\frac{1}{6} \quad \frac{4}{6} \quad \frac{1}{6}$	$\frac{h^5}{90} f^{(4)}(\tau)$	Simpson's rule, Kepler's barrel rule
3	$\frac{1}{8} \quad \frac{3}{8} \quad \frac{3}{8} \quad \frac{1}{8}$	$\frac{3h^5}{80} f^{(4)}(\tau)$	Newton's 3/8-rule
4	$\frac{7}{90} \quad \frac{32}{90} \quad \frac{12}{90} \quad \frac{32}{90} \quad \frac{7}{90}$	$\frac{8h^7}{945} f^{(6)}(\tau)$	Milne's rule

Lemma Let $f \in C^2([a, b])$ be a twice continuously differentiable function. Then the integration error of the trapezoidal rule

$$T = \frac{b-a}{2}(f(a) + f(b))$$

with step size $h = b - a$ can be expressed by

$$T - \int_a^b f = \frac{(b-a)^3}{12} f''(\tau),$$

for some $\tau \in [a, b]$

Proof \rightsquigarrow board

Integration errors of other Newton-Cotes formulas are listed in the table on the previous slide.

Helper mean value theorem:

Let $g, h \in C([a, b])$ be cont. functions on $[a, b]$,
where g has only one sign,

either $g(t) \geq 0$ or $g(t) < 0$ on $[a, b]$

then there exists $\tau \in [a, b]$ such that

$$\int_a^b h(t) g(t) dt = h(\tau) \int_a^b g(t) dt$$

Let us assume $g(t) \geq 0$.

$$\begin{aligned} \min_{t \in [a, b]} h(t) \int_a^b g(s) ds &\leq \int_a^b h(s) g(s) ds \\ &\leq \max_{t \in [a, b]} h(t) \int_a^b g(s) ds \end{aligned}$$

Therefore, for cont. function

$$F(t) = \int_a^b h(s) g(s) ds - h(t) \int_a^b g(s) ds$$

there exists $t_0, t_1 \in [a, b]$ with

$$F(t_0) \geq 0$$

$$F(t_1) \leq 0$$

Because cont. function F , we find $\tau \in [a, b]$

$$F(\tau) = 0.$$

$$\Rightarrow \int_a^b h(s) g(s) ds = h(\tau) \int_a^b g(s) ds$$

Error of trapezoidal rule:

Let $f \in C^2([a, b])$ be twice cont. diff. Then the approximation error of the trap. rule

$$T = \frac{b-a}{2} (f(a) + f(b))$$

with step size $h = b-a$ can be expressed as

$$T - \int_a^b f = \frac{h^3}{12} f''(\tau), \quad \tau \in [a, b].$$

Linear interpolation $P \in \mathcal{P}_1$ satisfies

$$f(t) = P(t) + [\alpha, b, t] f (t-a)(t-b)$$

(\rightarrow Newton interpolation)

Also know that for $f \in C^2([a, b])$

$$[\alpha, b, t] f = \frac{f''(\tau(t))}{2} \quad \tau(t) \in [a, b]$$

$$\int_a^b f = \int_a^b P(t) dt + \int_a^b [\alpha, b, t] f \underbrace{(t-a)(t-b)}_{\leq 0, a \leq t \leq b} dt$$

Because $g(t) = (t-a)(t-b)$ has one sign over $[a, b]$,

$\exists \bar{t} \in [a, b]$ such that

$$\int [\alpha, b, t] f g(t) = \underbrace{[\alpha, b, \bar{t}] f}_{\text{constant}} \int g(t)$$

$$\int_a^b f = \underbrace{T}_{\substack{\text{trapezoidal} \\ \text{rule integral}}} + \frac{f''(\tau(\bar{t}))}{2} \underbrace{\int_a^b (t-a)(t-b) dt}_{= -\frac{(b-a)^3}{6}}$$

$$T - \int_a^b f = \frac{(b-a)^3}{12} f''(\tau) \quad \tau \in (a, b)$$

We use equidistantly spaced nodes?

*almost equal because the f does not change too much (smoothness) from t_i to t_{i+1}

We use equidistantly spaced nodes? We know that this leads to poorly conditioned interpolation problems!

*almost equal because the f does not change too much (smoothness) from t_i to t_{i+1}

We use equidistantly spaced nodes? We know that this leads to poorly conditioned interpolation problems!

Recall: Polynomial interpolation can lead to Runge's phenomenon \rightsquigarrow high frequency oscillations \rightsquigarrow we saw that the relative condition number of numerical integration can increase for oscillatory functions

Another point of view: The weights $\lambda_{n0}, \dots, \lambda_{nn}$ of Newton-Cotes formulas can become *negative* for larger n \rightsquigarrow cancellation because we subtract "almost equal numbers*" in

$$\hat{I}(f) = (b - a) \sum_{i=0} \lambda_i f(t_i)$$

Weights are positive up to order 7, then some start to become negative.

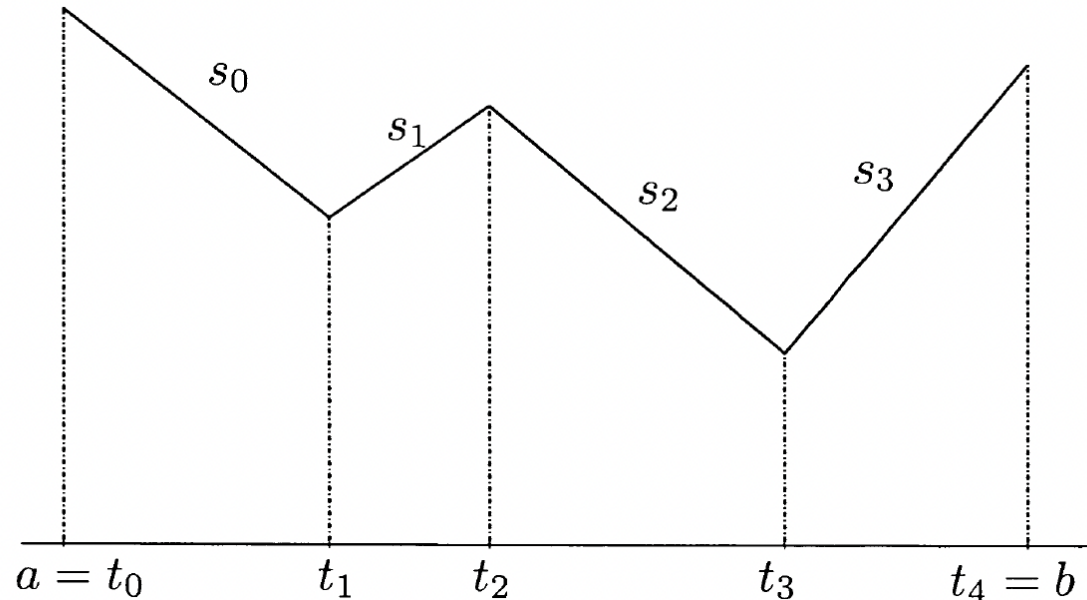
*almost equal because the f does not change too much (smoothness) from t_i to t_{i+1}

Trapezoidal sums

To avoid poorly conditioned problems, let us split the integration interval $[a, b]$ into n sub-intervals $[t_{i-1}, t_i]$, $i = 1, \dots, n$. Then consider the rule

$$\hat{I}(f) = \sum_{i=1}^n \hat{I}_{t_{i-1}}^{t_i}(f),$$

where $\hat{I}_{t_{i-1}}^{t_i}$ is a quadrature formula on the interval $[t_{i-1}, t_i]$.



Same assumptions on f as before. We have seen already the trapezoidal sum with $h = (b - a)/n$

$$T(h) = \sum_{i=1}^n T_i = h \left(\frac{1}{2}(f(a) + f(b)) + \sum_{i=1}^{n-1} f(a + ih) \right)$$

It has error

$$T(h) - \int_a^b f = \frac{(b-a)h^2}{12} f''(\tau), \quad \tau \in [a, b]$$

Proof \rightsquigarrow board

Trapezoidal sum

$\hookrightarrow \exists \tau_i \in [t_{i-1}, t_i]$ such that

$$T_i - \int_{t_{i-1}}^{t_i} f = \frac{h^3}{12} f''(\tau_i) \quad \left. \vphantom{\int} \right\} \text{just showed this}$$

$$\begin{aligned} T(h) - \int_0^b f &= \sum_{i=1}^n (T_i - \int_{t_{i-1}}^{t_i} f) \\ &= \sum_{i=1}^n \frac{h^3}{12} f''(\tau_i) \\ &= \sum_{i=1}^n \frac{h^2}{12} \frac{(b-a)}{n} f''(\tau_i) \\ &= \frac{(b-a)h^2}{12} \frac{1}{n} \sum_{i=1}^n f''(\tau_i) \end{aligned}$$

Now get rid of dependence of τ on intervals

$$\min_{t \in [a,b]} f''(t) \leq \frac{1}{n} \sum_{i=1}^n f''(\tau_i) \leq \max_{t \in [a,b]} f''(t)$$

Mean value theorem: $\exists \tau \in [a,b]$

$$\frac{1}{n} \sum_{i=1}^n f''(\tau_i) = f''(\tau)$$

$$\Rightarrow T(h) - \int_0^b f = \frac{(b-a)h^2}{12} f''(\tau)$$

Same assumptions on f as before. We have seen already the trapezoidal sum with $h = (b - a)/n$

$$T(h) = \sum_{i=1}^n T_i = h \left(\frac{1}{2}(f(a) + f(b)) + \sum_{i=1}^{n-1} f(a + ih) \right)$$

It has error

$$T(h) - \int_a^b f = \frac{(b-a)h^2}{12} f''(\tau), \quad \tau \in [a, b]$$

Proof \rightsquigarrow **board**

What did we achieve with this?

Same assumptions on f as before. We have seen already the trapezoidal sum with $h = (b - a)/n$

$$T(h) = \sum_{i=1}^n T_i = h \left(\frac{1}{2}(f(a) + f(b)) + \sum_{i=1}^{n-1} f(a + ih) \right)$$

It has error

$$T(h) - \int_a^b f = \frac{(b-a)h^2}{12} f''(\tau), \quad \tau \in [a, b]$$

Proof \rightsquigarrow **board**

What did we achieve with this? \rightsquigarrow We can increase n (and thus decrease h) to reduce the error without increasing the degree of the underlying polynomial



Numerical Methods I

MATH-GA 2010.001/CSCI-GA 2420.001

Benjamin Peherstorfer
Courant Institute, NYU

Based on slides by G. Stadler and A. Donev

Today

Last week

- ▶ Numerical integration
- ▶ Newton-Cotes formulas

Today

- ▶ Gauss quadrature
- ▶ Integration in multiple dimensions

Announcements

- ▶ Homework 6 is posted and due Mon, Dec 2 before class

Recap: Condition of numerical integration

Lets study the map

$$([a, b], f) \rightarrow \int_a^b f(t) dt,$$

where we use the L^1 -norm for f

$$\|f\|_1 = \int_a^b |f(t)| dt = I(|f|)$$

The absolute and relative condition numbers of integration are:

$$\begin{aligned}\kappa_{\text{abs}} &= 1, \\ \kappa_{\text{rel}} &= \frac{I(|f|)}{|I(f)|}.\end{aligned}$$

What does this mean?

Recap: Condition of numerical integration

Lets study the map

$$([a, b], f) \rightarrow \int_a^b f(t) dt,$$

where we use the L^1 -norm for f

$$\|f\|_1 = \int_a^b |f(t)| dt = I(|f|)$$

The absolute and relative condition numbers of integration are:

$$\begin{aligned}\kappa_{\text{abs}} &= 1, \\ \kappa_{\text{rel}} &= \frac{I(|f|)}{|I(f)|}.\end{aligned}$$

What does this mean? Integration is harmless w.r.t. the absolute condition number, and problematic w.r.t. the relative condition number if $I(f)$ is small and f changes sign

Recap: Newton-Cotes formulas

Given **fixed** nodes t_0, \dots, t_n , use polynomial approximation

$$\hat{f} = P_f(t|t_0, \dots, t_n) = \sum_{i=0}^n f(t_i)L_{in}(t)$$

with Lagrange polynomials L_{0n}, \dots, L_{nn}

Thus:

$$\hat{I}(f) = (b - a) \sum_{i=0}^n \lambda_{in} f(t_i),$$

where

$$\lambda_{in} = \frac{1}{b - a} \int_a^b L_{in}(t) dt$$

Recap: Newton-Cotes formulas (cont'd)

Quadrature formulas defined in this way are exact for polynomials $P \in \mathbb{P}_n$ of degree less than or equal to n

$$\hat{I}(P) = I(P_n(P)) = I(P), \quad \text{for all } P \in \mathbb{P}_n$$

Theorem: For $(n + 1)$ pairwise distinct nodes t_0, \dots, t_n , there exists exactly one quadrature formula (i.e., unique weights $\lambda_0, \dots, \lambda_n$)

$$\hat{I}(f) = (b - a) \sum_{i=0}^n \lambda_i f(t_i),$$

that is exact for all $p \in \mathbb{P}_n$.

Recap

Equidistantly spaced nodes

$$h_i = h = \frac{b - a}{n}, \quad t_i = a + ih, \quad i = 0, \dots, n$$

then quadrature formulas are called the *Newton-Cotes formulas* with weights

$$\lambda_{in} = \frac{1}{b - a} \int_a^b \prod_{i \neq j} \frac{t - t_i}{t_i - t_j} dt = \frac{1}{n} \int_0^n \prod_{i \neq j} \frac{s - j}{i - j} ds$$

These weights are independent of the interval boundaries a and b and can be pre-computed once and for all:

Table 9.1. Newton-Cotes weights λ_{in} for $n = 1, \dots, 4$.

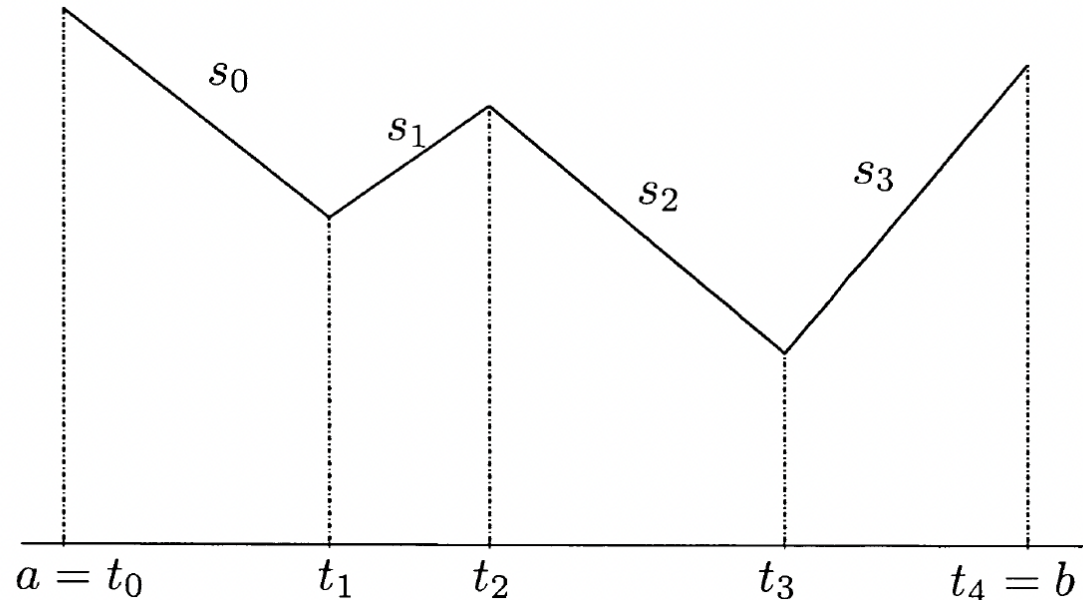
n	$\lambda_{0n}, \dots, \lambda_{nn}$	Error	Name
1	$\frac{1}{2} \quad \frac{1}{2}$	$\frac{h^3}{12} f''(\tau)$	Trapezoidal rule
2	$\frac{1}{6} \quad \frac{4}{6} \quad \frac{1}{6}$	$\frac{h^5}{90} f^{(4)}(\tau)$	Simpson's rule, Kepler's barrel rule
3	$\frac{1}{8} \quad \frac{3}{8} \quad \frac{3}{8} \quad \frac{1}{8}$	$\frac{3h^5}{80} f^{(4)}(\tau)$	Newton's 3/8-rule
4	$\frac{7}{90} \quad \frac{32}{90} \quad \frac{12}{90} \quad \frac{32}{90} \quad \frac{7}{90}$	$\frac{8h^7}{945} f^{(6)}(\tau)$	Milne's rule

Recap: Trapezoidal sums

To avoid poorly conditioned problems, let us split the integration interval $[a, b]$ into n sub-intervals $[t_{i-1}, t_i], i = 1, \dots, n$. Then consider the rule

$$\hat{I}(f) = \sum_{i=1}^n \hat{I}_{t_{i-1}}^{t_i}(f),$$

where $\hat{I}_{t_{i-1}}^{t_i}$ is a quadrature formula on the interval $[t_{i-1}, t_i]$.



Recap: Error of Trapezoidal sum

Trapezoidal sum with $h = (b - a)/n$

$$T(h) = \sum_{i=1}^n T_i = h \left(\frac{1}{2}(f(a) + f(b)) + \sum_{i=1}^{n-1} f(a + ih) \right)$$

It has error

$$T(h) - \int_a^b f = \frac{(b-a)h^2}{12} f''(\tau), \quad \tau \in [a, b]$$

What did we achieve with this? \rightsquigarrow

Recap: Error of Trapezoidal sum

Trapezoidal sum with $h = (b - a)/n$

$$T(h) = \sum_{i=1}^n T_i = h \left(\frac{1}{2}(f(a) + f(b)) + \sum_{i=1}^{n-1} f(a + ih) \right)$$

It has error

$$T(h) - \int_a^b f = \frac{(b-a)h^2}{12} f''(\tau), \quad \tau \in [a, b]$$

What did we achieve with this? \rightsquigarrow We can increase n (and thus decrease h) to reduce the error without increasing the degree of the underlying polynomial

Gauss-Christoffel quadrature

Besides piecewise approximations, what else could we do to avoid issues with high-degree polynomials?

Gauss-Christoffel quadrature

Besides piecewise approximations, what else could we do to avoid issues with high-degree polynomials? **Change the nodes**

So far, we allowed to choose weights but the nodes were given (e.g., equidistant). Now, let's allow changing the nodes t_0, \dots, t_n too.

What highest degree can we hope for?

Gauss-Christoffel quadrature

Besides piecewise approximations, what else could we do to avoid issues with high-degree polynomials? **Change the nodes**

So far, we allowed to choose weights but the nodes were given (e.g., equidistant). Now, let's allow changing the nodes t_0, \dots, t_n too.

What highest degree can we hope for? The best we can hope for is exact interpolation up to polynomials of degree $2n + 1$ ($2(n + 1)$ coefficients) based on a non-rigorous counting argument of having $n + 1$ DoFs because of nodes t_0, \dots, t_n and $n + 1$ DoFs because of weights $\lambda_0, \dots, \lambda_n$.

Also, for generalization, we consider quadrature of weighted integrals, with a positive weight function $\omega(t)$:

$$I(f) = \int_a^b \omega(t)f(t) dt$$

with weight functions $\omega(t) = 1, \omega(t) = 1/\sqrt{1-t^2}, \dots$

Our goal is the construction of quadrature formulas of the form

$$\hat{I}_n(f) = \sum_{i=0}^n \lambda_{in} f(\tau_{in}),$$

to approximate $I(f)$. Thus, for a given n , we seek $n + 1$ nodes $\tau_{0n}, \dots, \tau_{nn}$ and $n + 1$ weights $\lambda_{0n}, \dots, \lambda_{nn}$ so that polynomials up to as high degree N as possible can be integrated exactly

$$\hat{I}_n(P) = I(P), \quad \text{for all } P \in \mathbb{P}_N$$

Gauss-Christoffel quadrature (cont'd)

Can we say something about the nodes $\tau_{0n}, \dots, \tau_{nn}$?

If \tilde{I}_n is exact for all polynomials $P \in \mathbb{P}_{2n+1}$, then the polynomials $\{P_j\}$ that are defined by their root representation

$$P_{j+1}(t) = (t - \tau_{0j}) \dots (t - \tau_{jj}) \in \mathbb{P}_{n+1}$$

are orthogonal with respect to the scalar product that is induced by w

$$(f, g) = \int_a^b w(t) f(t) g(t) dt$$

Proof: \tilde{I}_n exact for polynomials in \mathbb{P}_{2n+1} :

for $j < n+1$, have $P_j P_{n+1} \in \mathbb{P}_{2n+1}$

$$\begin{aligned} (P_j, P_{n+1}) &= \int w P_j P_{n+1} = \tilde{I}_n(P_j P_{n+1}) \\ &= \sum_{i=0}^n \lambda_{in} P_j(\tau_{in}) \underbrace{P_{n+1}(\tau_{in})}_{=0} = 0 \end{aligned}$$

because $P_{n+1}(t) = \prod_{i=0}^n (t - \tau_{in})$.

Gauss-Christoffel quadrature (cont'd)

Can we say something about the nodes $\tau_{0n}, \dots, \tau_{nn}$?

Theorem: For $n \in \mathbb{N}$, consider nodes $\tau_{0n}, \dots, \tau_{nn}$. For an n , define \hat{I} as

$$\hat{I}_n(f) = \sum_{i=0}^n \lambda_{in} f(\tau_{in})$$

and let it be exact for polynomials $p \in \mathbb{P}_{2n+1}$

$$\hat{I}_n(p) = \int_a^b \omega(t) p(t) dt.$$

Then, the polynomials $\{P_n\}$ given by

$$P_{n+1}(t) = (t - \tau_{0n}) \cdot \dots \cdot (t - \tau_{nn}) \in \mathbb{P}_{n+1}$$

are orthogonal with respect to the scalar product induced by $\omega(t)$ \rightsquigarrow **proof**

Gauss-Christoffel quadrature (cont'd)

Therefore, the nodes τ_{in} have to be roots of pairwise orthogonal polynomials $\{P_j\}$ of degree $\deg(P_j) = j$

For a given ω , we already know that the set of orthogonal polynomials $\{P_j\}$ is unique (if leading coefficient is 1)

We also know that the roots of these polynomials are real and have to lie in $(a, b) \rightsquigarrow$
proof board

For each ω and $n \in \mathbb{N}$, this gives us a unique set of nodes, namely the roots of the corresponding orthogonal polynomial P_{n+1}

We thus have fixed $n + 1$ degrees of freedom so far...

The orthogonal polynomial $P_k \in \mathbb{P}_k$ possess exactly k simple roots in (a, b) .

Let t_1, \dots, t_m be m distinct points in (a, b) of which P_k changes sign. We now show that $m = k$.

Consider the polynomial

$$Q(t) = (t-t_1)(t-t_2) \dots (t-t_m)$$

The polynomial changes sign at the same points in (a, b) as P_k . Thus: $w(t)Q(t)P_k(t)$ does not change sign

$$(Q, P_k) = \int_a^b w(t) Q(t) P_k(t) dt \neq 0$$

But P_k orthogonal to all $P \in \mathbb{P}_{k-1}$, thus

$$\deg(Q) = m \geq k.$$

What do we know about the weights for fixed t_0, \dots, t_n when we want to integrate exactly polynomials up to degree n ?

What do we know about the weights for fixed t_0, \dots, t_n when we want to integrate exactly polynomials up to degree n ?

Because we want to be able to exactly integrate polynomials up to degree $2n + 1$, we already know that to even achieve exactness up to degree n , the weights are fixed for given nodes t_0, \dots, t_n :

$$\lambda_{in} = \frac{1}{b-a} \int_a^b L_{in}(t) dt, \quad \text{Lagrange poly } L_{in}(\tau_{jn}) = \delta_{ij}$$

(“For $n + 1$ pairwise distinct nodes, there exists only one quadrature formula that exactly integrates polynomials up to degree n .”)

Theorem: Let $\tau_{0n}, \dots, \tau_{nn}$ be the roots of the $(n + 1)$ st orthogonal polynomial for the weight ω . Then any quadrature formula \hat{I} is exact for polynomials up to order n if and only if it is exact up to order $2n + 1$.

Proof

Let $\tau_{0n}, \dots, \tau_{nn}$ be the roots of the $(n+1)$ -st orthogonal polynomial P_{n+1} with weight w . Then any quadrature rule

$$\hat{I}_n(f) = \sum_{i=0}^n \lambda_i f(\tau_{in})$$

satisfies

$$\hat{I}_n \text{ exact on } \mathbb{P}_n \iff \hat{I}_n \text{ exact on } \mathbb{P}_{2n+1}$$

" \implies ": Suppose that \hat{I}_n is exact on \mathbb{P}_n .

$P \in \mathbb{P}_{2n+1}$, then there exist polynomials

$$Q, R \in \mathbb{P}_n$$

such that

$$P = Q P_{n+1} + R$$

P_{n+1} is orthogonal to \mathbb{P}_n w.r.t. weight w

$$\int w P = \underbrace{\int w Q P_{n+1}}_{=0} + \int w R = \int w R = \hat{I}_n(R)$$

Additionally, with $P_{n+1}(\tau_{in}) = 0$ obtain

$$\hat{I}_n(R) = \sum_{i=0}^n \lambda_{in} R(\tau_{in})$$

$$= \sum_{i=0}^n \lambda_{in} (Q(\tau_{in}) \underbrace{P_{n+1}(\tau_{in})}_{=0} + R(\tau_{in})) = \hat{I}_n(P)$$

$q \leftarrow i: \dots$

Weight functions for Gauss-Christoffel quadrature

$\omega(t)$	Interval $I = [a, b]$	Orthogonal polynomials
$\frac{1}{\sqrt{1-t^2}}$	$[-1, 1]$	Chebyshev polynomials T_n
e^{-t}	$[0, \infty]$	Laguerre polynomials L_n
e^{-t^2}	$[-\infty, \infty]$	Hermite polynomials H_n
1	$[-1, 1]$	Legendre polynomials P_n

Corresponding quadrature rules are usually prefixed with “Gauss-”, i.e., “Gauss-Legendre quadrature”, or “Gauss-Chebyshev quadrature”.

Summary of Gauss quadrature

There exist uniquely determined nodes $\tau_{0n}, \dots, \tau_{nn}$ and weights $\lambda_{0n}, \dots, \lambda_{nn}$ such that the quadrature formula

$$\hat{I}_n(f) = \sum_{i=0}^n \lambda_{in} f(\tau_{in})$$

integrates exactly all polynomials of degree less than or equal to $2n + 1$, i.e.,

$$\hat{I}_n(P) = \int_a^b \omega P, \quad P \in \mathbb{P}_{2n+1}.$$

The nodes τ_{in} are the roots of the $n + 1$ -st orthogonal polynomial P_{n+1} with respect to the weight function ω and the weights are

$$\lambda_{in} = \frac{1}{b-a} \int_a^b L_{in}(t) dt,$$

with the Lagrange polynomials $L_{in}(\tau_{jn}) = \delta_{ij}$. Furthermore, the weights are all positive $\lambda_{in} > 0$. \rightsquigarrow proof

Unique nodes τ_0, \dots, τ_n and weights $\lambda_0, \dots, \lambda_n$ such that

$$\hat{I}_n(f) = \sum_{i=0}^n \lambda_i f(\tau_i)$$

integrates exactly $P \in \mathbb{P}_{2n+1}$. The weights are positive.

Let $Q \in \mathbb{P}_{2n+1}$ be a polynomial so that τ_{2n} is the only node at which it does not vanish:

$$Q(\tau_{2n}) \neq 0, \quad Q(\tau_i) = 0 \quad i \neq 2n$$

Then

$$\int_a^b w Q = \hat{I}_n(Q) = \lambda_{2n} Q(\tau_{2n})$$

and so $\lambda_{2n} = \frac{1}{Q(\tau_{2n})} \int_a^b w Q$.

if we set

$$Q(t) = \left(\frac{P_{n+1}(t)}{(t - \tau_{2n})} \right)^2$$

then $Q \in \mathbb{P}_{2n}$ has these properties

Weights

$$\lambda_{2n} = \frac{1}{Q(\tau_{2n})} \int_a^b w Q = \int_a^b w \left(\frac{P_{n+1}(t)}{P_{n+1}'(\tau_{2n})(t - \tau_{2n})} \right)^2 dt > 0$$

$$Q(\tau_{2n}) = [P_{n+1}'(\tau_{2n})]^2$$

Approximation error

For any function $f \in C^{2n+2}$, the approximation error of Gauss quadrature can be expressed in the form

$$\int_a^b \omega f - \hat{I}_n(f) = \frac{f^{(2n+2)}(\tau)}{(2n+2)!} (P_{n+1}, P_{n+1}),$$

with some $\tau \in [a, b]$.

↪ proof

For any function $f \in C^{2n+2}$, the approximation error of Gauss quadrature can be expressed as

$$\int_a^b \omega f - \hat{I}_n(f) = \frac{f^{(2n+2)}(\tau)}{(2n+2)!} (P_{n+1}, P_{n+1})$$

with some $\tau \in [a, b]$, P_{n+1} ortho. poly.

We employ the Newton remainder of the Hermite interpolant $P \in \underline{P}_{2n+1}$ of f for the nodes

$$\tau_{0n}, \tau_{0n}, \dots, \tau_{nn}, \tau_{nn}$$

so that

$$f(t) = P(t) + [t, \tau_{0n}, \tau_{0n}, \dots, \tau_{nn}, \tau_{nn}] f (t - \tau_{0n})^2 \dots (t - \tau_{nn})^2$$

Notice that

$$(t - \tau_{0n})^2 \dots (t - \tau_{nn})^2 = P_{n+1}^2$$

Because \hat{I}_n integrates P exactly, we have

$$\begin{aligned} \int \omega f &= \underbrace{\int \omega P}_{\hat{I}_n(f)} + \frac{f^{(2n+2)}(\tau)}{(2n+2)!} \int \omega P_{n+1}^2 \\ &= \underbrace{\sum_{i=0}^n \lambda_{in} \underbrace{P(\tau_{in})}_{=f(\tau_{in})}}_{\hat{I}_n(f)} + \frac{f^{(2n+2)}(\tau)}{(2n+2)!} (P_{n+1}, P_{n+1}) \end{aligned}$$

Summary of Gauss quadrature (cont'd)

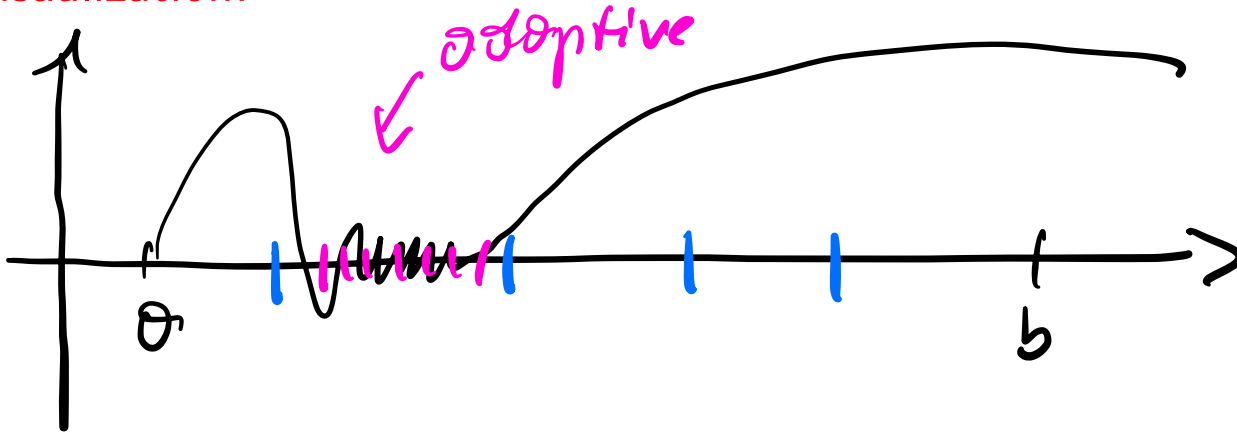
- ▶ Gauss quadrature gives the highest degree $2n + 1$ of polynomials that can be exactly integrated with $n + 1$ function evaluations in general
- ▶ Accuracy in Gauss-(Chebyshev, Laguerre, Hermite, Legendre,...) w.r.t. to polynomial degree can only be improved by increasing number of points; not by better weights
- ▶ Of particular interest are quadrature points for infinite intervals (Laguerre, Hermite)
- ▶ Interval partitioning superior, but only possible for $\omega \equiv 1$ (Gauss-Legendre); otherwise weight function is different in each sub-interval
- ▶ Gauss quadrature is typically used in finite-element approximation to integrate over a local element. However, in many other cases in scientific computing, interval partitioning is superior.

Adaptive interval partitioning

Idea: On each sub-interval, estimate the quadrature error by either:

- ▶ Using a higher-order quadrature (e.g., Simpson rule), or
- ▶ Comparing the error on a subinterval with the error on a refinement

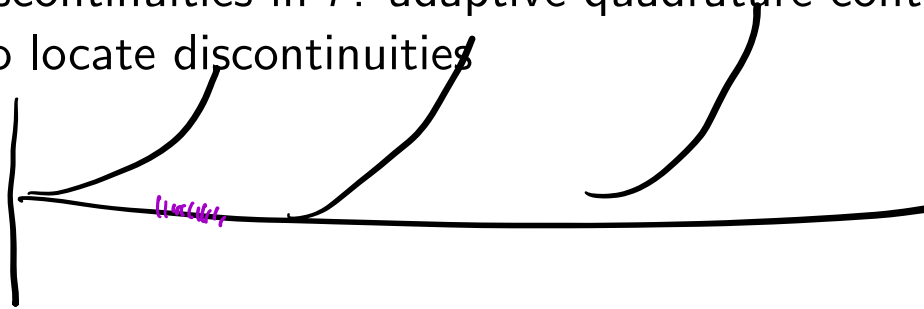
Then, subdivide the interval depending on the error estimation, and repeat \rightsquigarrow
visualization!



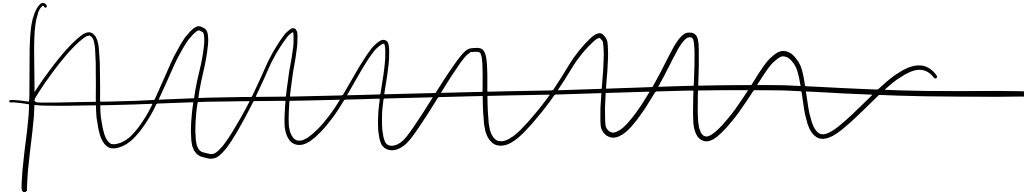
Such a method refines the nodes **a posteriori** (after having seen the function to integrate).

Difficult cases for numerical integration

- ▶ (Unknown) discontinuities in f : adaptive quadrature continues to refine, which can be used to locate discontinuities



- ▶ Highly oscillating integrals



- ▶ (Weakly) singular integrals



Numerical Methods I

MATH-GA 2010.001 / CSCI-GA 2420.001

Benjamin Peherstorfer
Courant Institute, NYU

Today

Last time

- ▶ Interpolation and quadrature in one dimension

Today

- ▶ Function approximation and interpolation in higher dimensions
- ▶ Quadrature in higher dimensions

Announcements

- ▶ Homework 6 is due Mon, Dec 2 before class

Recap: Newton-Cotes formulas

Given **fixed** nodes t_0, \dots, t_n , use polynomial approximation

$$\hat{f} = P_f(t|t_0, \dots, t_n) = \sum_{i=0}^n f(t_i) L_{in}(t)$$

with Lagrange polynomials L_{0n}, \dots, L_{nn}

Thus:

$$\hat{I}(f) = (b - a) \sum_{i=0}^n \lambda_{in} f(t_i),$$

where

$$\lambda_{in} = \frac{1}{b - a} \int_a^b L_{in}(t) dt$$

Recap

Equidistantly spaced nodes

$$h_i = h = \frac{b - a}{n}, \quad t_i = a + ih, \quad i = 0, \dots, n$$

then quadrature formulas are called the *Newton-Cotes formulas* with weights

$$\lambda_{in} = \frac{1}{b - a} \int_a^b \prod_{i \neq j} \frac{t - t_i}{t_i - t_j} dt = \frac{1}{n} \int_0^n \prod_{i \neq j} \frac{s - j}{i - j} ds$$

These weights are independent of the interval boundaries a and b and can be pre-computed once and for all:

Table 9.1. Newton-Cotes weights λ_{in} for $n = 1, \dots, 4$.

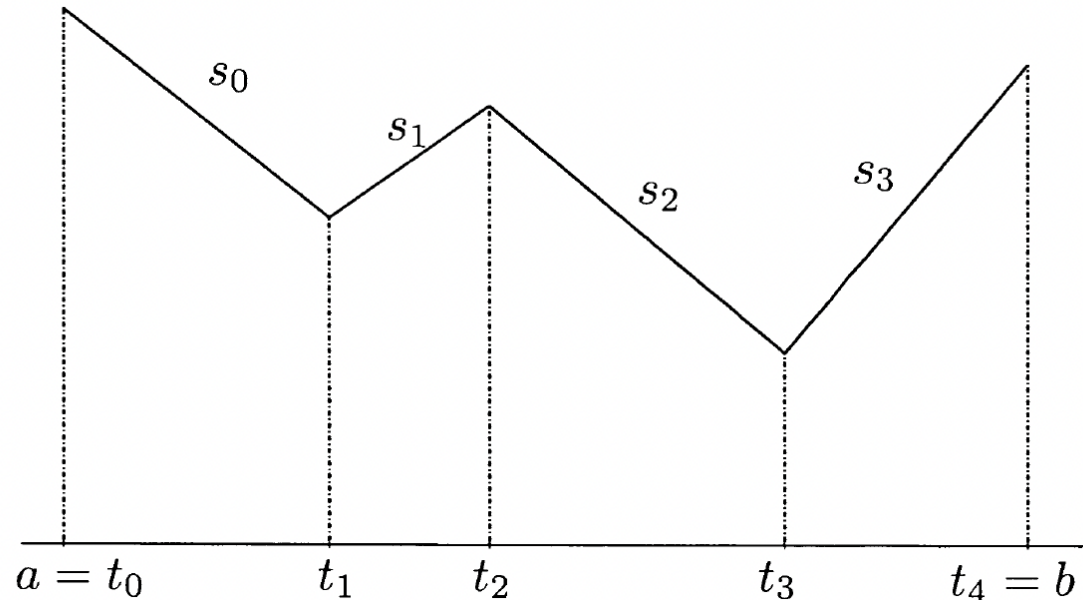
n	$\lambda_{0n}, \dots, \lambda_{nn}$	Error	Name
1	$\frac{1}{2} \quad \frac{1}{2}$	$\frac{h^3}{12} f''(\tau)$	Trapezoidal rule
2	$\frac{1}{6} \quad \frac{4}{6} \quad \frac{1}{6}$	$\frac{h^5}{90} f^{(4)}(\tau)$	Simpson's rule, Kepler's barrel rule
3	$\frac{1}{8} \quad \frac{3}{8} \quad \frac{3}{8} \quad \frac{1}{8}$	$\frac{3h^5}{80} f^{(4)}(\tau)$	Newton's 3/8-rule
4	$\frac{7}{90} \quad \frac{32}{90} \quad \frac{12}{90} \quad \frac{32}{90} \quad \frac{7}{90}$	$\frac{8h^7}{945} f^{(6)}(\tau)$	Milne's rule

Recap: Trapezoidal sums

To avoid poorly conditioned problems, let us split the integration interval $[a, b]$ into n sub-intervals $[t_{i-1}, t_i]$, $i = 1, \dots, n$. Then consider the rule

$$\hat{I}(f) = \sum_{i=1}^n \hat{I}_{t_{i-1}}^{t_i}(f),$$

where $\hat{I}_{t_{i-1}}^{t_i}$ is a quadrature formula on the interval $[t_{i-1}, t_i]$.



Recap: Error of Trapezoidal sum

Trapezoidal sum with $h = (b - a)/n$

$$T(h) = \sum_{i=1}^n T_i = h \left(\frac{1}{2}(f(a) + f(b)) + \sum_{i=1}^{n-1} f(a + ih) \right)$$

It has error

$$T(h) - \int_a^b f = \frac{(b-a)h^2}{12} f''(\tau), \quad \tau \in [a, b]$$

What did we achieve with this? \rightsquigarrow

Recap: Error of Trapezoidal sum

Trapezoidal sum with $h = (b - a)/n$

$$T(h) = \sum_{i=1}^n T_i = h \left(\frac{1}{2}(f(a) + f(b)) + \sum_{i=1}^{n-1} f(a + ih) \right)$$

It has error

$$T(h) - \int_a^b f = \frac{(b - a)h^2}{12} f''(\tau), \quad \tau \in [a, b]$$

What did we achieve with this? \rightsquigarrow We can increase n (and thus decrease h) to reduce the error without increasing the degree of the underlying polynomial

Summary of Gauss quadrature

There exist uniquely determined nodes $\tau_{0n}, \dots, \tau_{nn}$ and weights $\lambda_{0n}, \dots, \lambda_{nn}$ such that the quadrature formula

$$\hat{I}_n(f) = \sum_{i=0}^n \lambda_{in} f(\tau_{in})$$

integrates exactly all polynomials of degree less than or equal to $2n + 1$, i.e.,

$$\hat{I}_n(P) = \int_a^b \omega P, \quad P \in \mathbb{P}_{2n+1}.$$

The nodes τ_{in} are the roots of the $n + 1$ -st orthogonal polynomial P_{n+1} with respect to the weight function ω and the weights are

$$\lambda_{in} = \frac{1}{b-a} \int_a^b L_{in}(t) dt,$$

with the Lagrange polynomials $L_{in}(\tau_{jn}) = \delta_{ij}$. Furthermore, the weights are all positive $\lambda_{in} > 0$.

Integration in higher dimensions

Integration in higher dimensions

A separable integral can be integrated dimension-wise:

$$I_L = \int_a^b \int_a^b \phi(x, y) dx dy = \int_a^b \phi^{(x)}(x) dx \int_a^b \phi^{(y)}(y) dy,$$

where

$$\phi(x, y) = \phi^{(x)}(x)\phi^{(y)}(y) \quad (1)$$

Recall that one idea of numerical quadrature is to replace f with an easy-to-integrate \hat{f}

- ▶ Choose a basis $\phi_1(x, y), \dots, \phi_n(x, y)$ and approximate

$$f(x, y) \approx \sum_{i=1}^n c_i \phi_i(x, y),$$

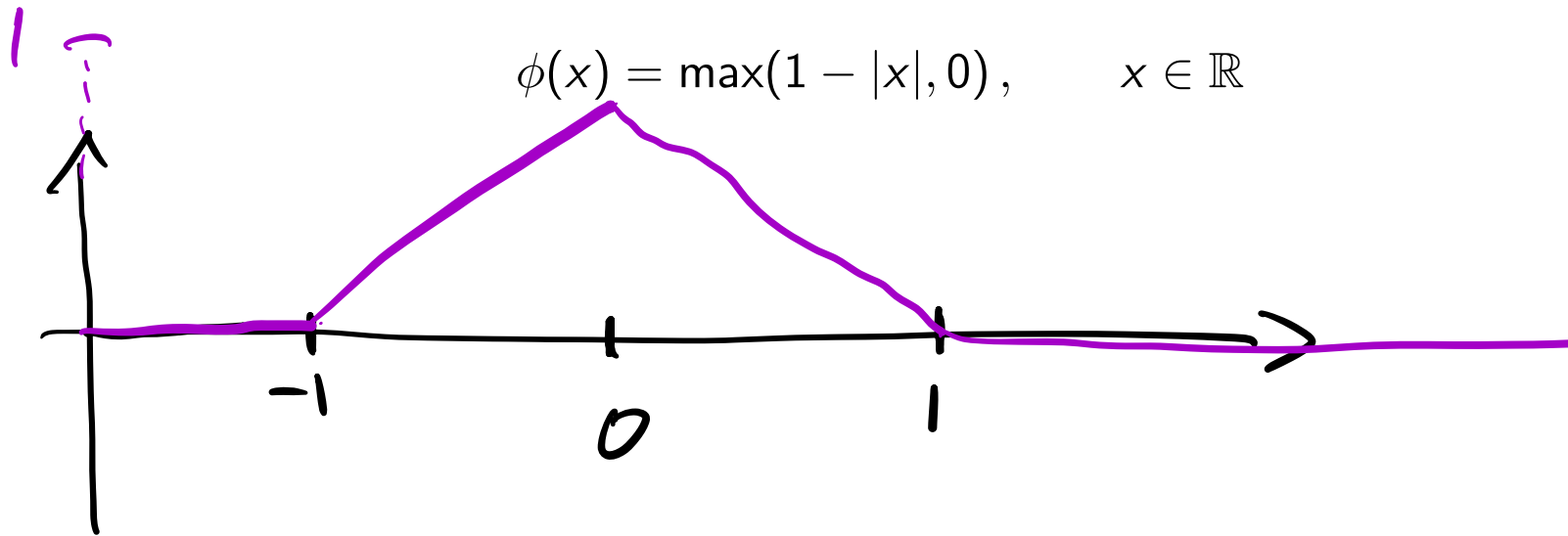
and choose $\phi_i(x, y)$ such that (1) holds.

- ▶ Then integrate

$$I(f) \approx \sum_{i=1}^n c_i \hat{I}(\phi_i^{(x)}) \hat{I}(\phi_i^{(y)})$$

One way to build such a basis is via tensor products of linear functions
(multi-dimensional analog to piecewise linear quadrature)

Define the “mother hat” function

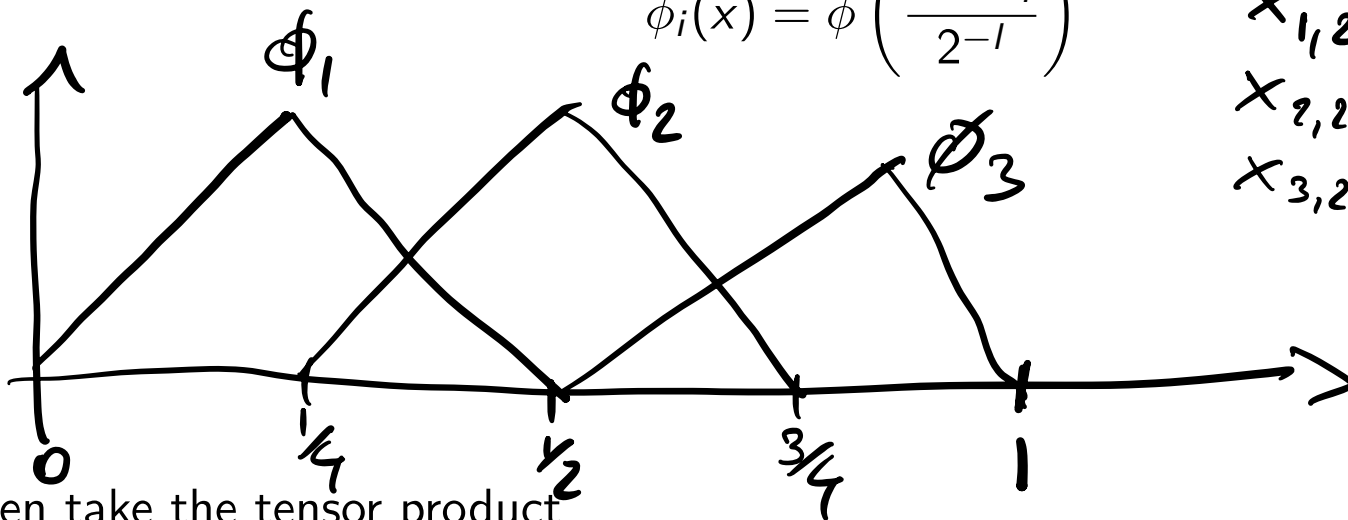


If we want $n = 2^l + 1$ basis functions, then translate and dilate it to center it on the grid points

$$x_i = i2^{-l}, \quad i = 0, \dots, n-1$$

$$\phi_i(x) = \phi\left(\frac{x - x_i}{2^{-l}}\right)$$

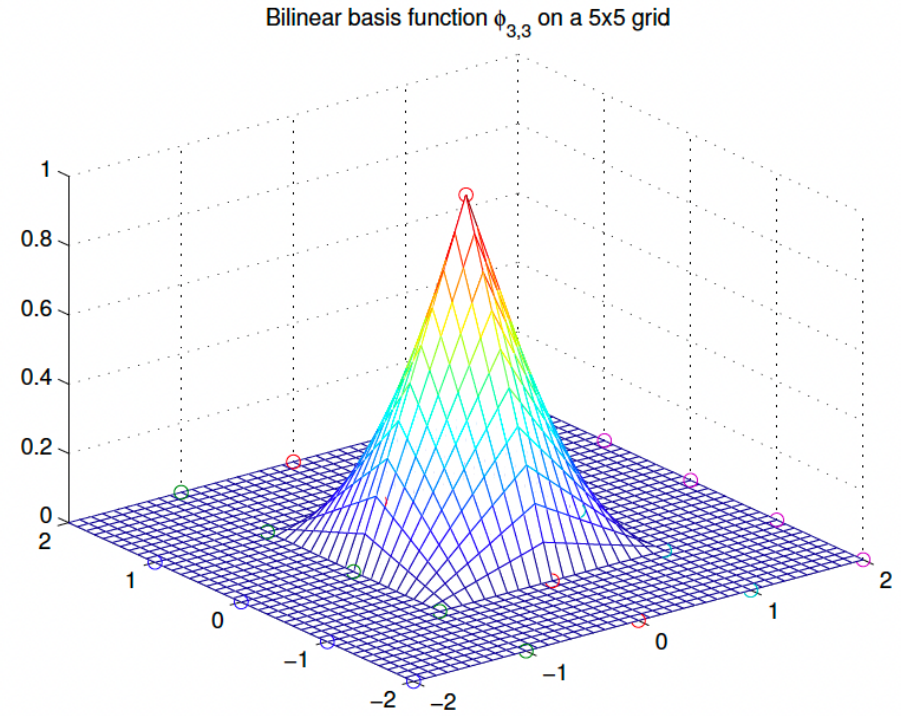
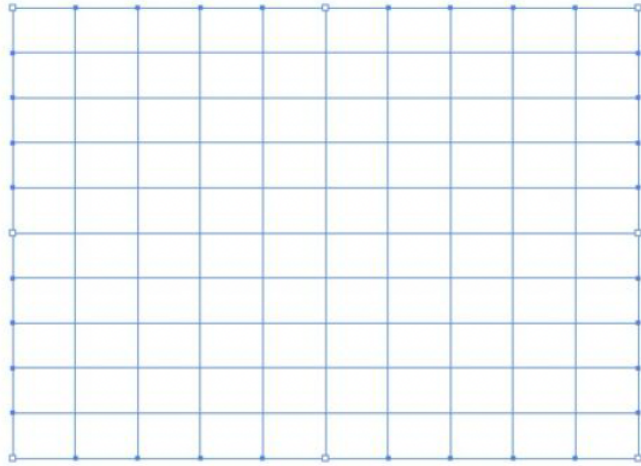
$$\begin{aligned} x_{1,2} &= 1 \times 2^{-2} = \frac{1}{4} \\ x_{2,2} &= 2 \times 2^{-2} = \frac{2}{4} \\ x_{3,2} &= 3 \times 2^{-2} = \frac{3}{4} \end{aligned}$$



Then take the tensor product

$$\phi_{i,j}(x, y) = \phi_i(x)\phi_j(y), \quad i, j = 0, \dots, n-1$$

lead to a basis of the piecewise bilinear functions in two dimensions



[Figure: A. Donev]

What type of basis is this?

What type of basis is this? The basis $\{\phi_{ij}\}$ is a nodal point basis for piecewise bilinear functions in $[0, 1]^2 \rightsquigarrow$

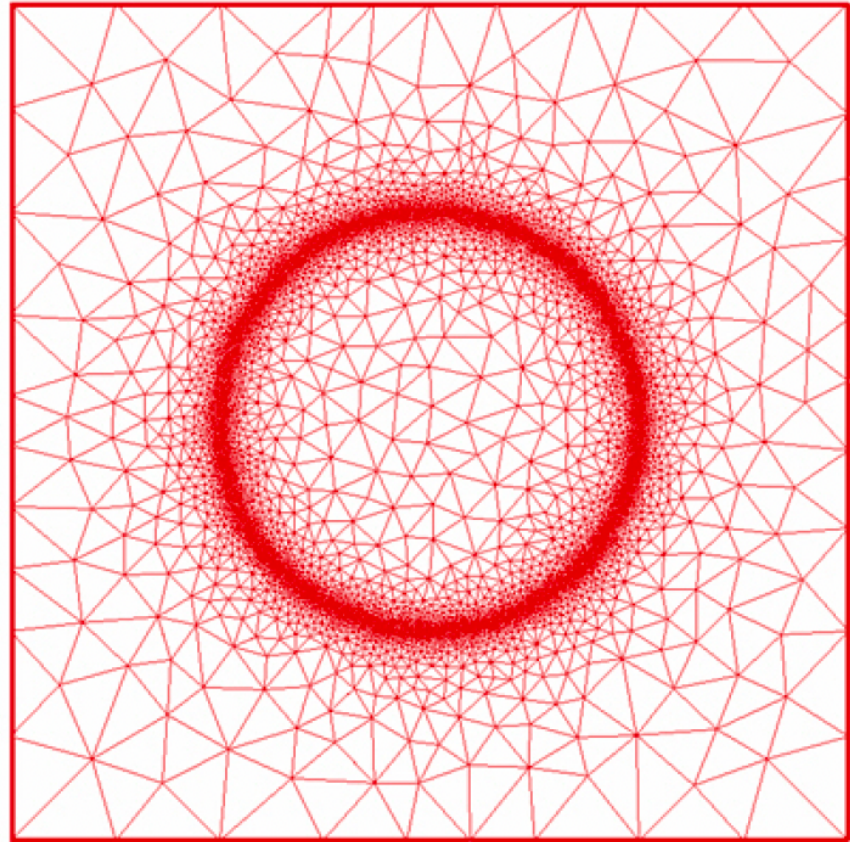
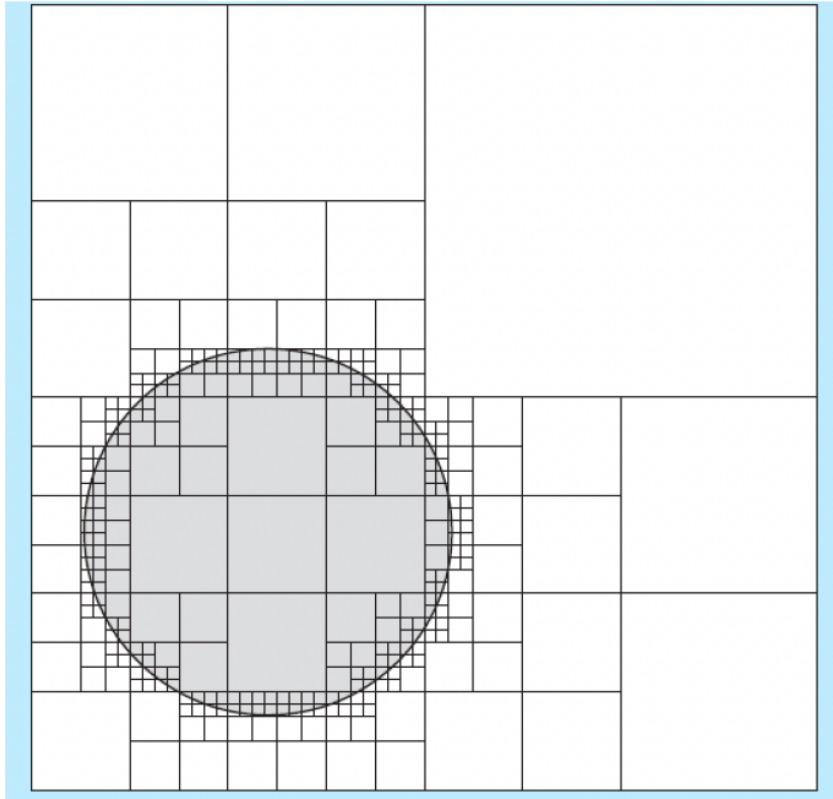
$$\hat{f}(x, y) = \sum_{i,j} f(x_i, y_j) \phi_{ij}(x, y)$$

interpolates f at the nodes $\{(x_i, y_j)\}_{i,j=0}^n$

Integrate the bilinear basis functions $\{\phi_{ij}\}$ to obtain the weights λ_{ij} and then approximate

$$\hat{I}(f) = \sum_{i,j} f(x_i, y_i) \lambda_{ij},$$

In multiple dimensions, adaptivity is essential to keep the computational costs manageable



[Figure: A. Donev]

Sparse grids

Curse of dimensionality

- ▶ The curse of dimensionality is a term coined by Bellmann (1961) that refers to an exponential increase of costs with the dimension of a problem
- ▶ For example, consider an approximation with a prescribed accuracy $\epsilon > 0$, let the costs of achieving this approximation scale as $\mathcal{O}(\epsilon^{-d})$ in d dimensions \rightsquigarrow exponential increase of the costs as we increase d
- ▶ Consider a simple uniform grid over the domain $\Omega = [0, 1]^d$. To have a mesh with mesh width $h = 1/9$ in $d = 1$, we need $N = 10$ grid points. In $d = 2$ dimensions, need $N^2 = 100$ grid points. In $d = 5$ dimensions, need $N^5 = 10^5 \rightsquigarrow$ exponential growth of cost and storage requirements as we increase dimension d while keeping mesh width h (“accuracy”) fixed

Curse of dimensionality

- ▶ The curse of dimensionality is a term coined by Bellmann (1961) that refers to an exponential increase of costs with the dimension of a problem
- ▶ For example, consider an approximation with a prescribed accuracy $\epsilon > 0$, let the costs of achieving this approximation scale as $\mathcal{O}(\epsilon^{-d})$ in d dimensions \rightsquigarrow exponential increase of the costs as we increase d
- ▶ Consider a simple uniform grid over the domain $\Omega = [0, 1]^d$. To have a mesh with mesh width $h = 1/9$ in $d = 1$, we need $N = 10$ grid points. In $d = 2$ dimensions, need $N^2 = 100$ grid points. In $d = 5$ dimensions, need $N^5 = 10^5 \rightsquigarrow$ exponential growth of cost and storage requirements as we increase dimension d while keeping mesh width h (“accuracy”) fixed
- ▶ The curse can be circumvented (to some extent) **if?**

Curse of dimensionality

- ▶ The curse of dimensionality is a term coined by Bellmann (1961) that refers to an exponential increase of costs with the dimension of a problem
- ▶ For example, consider an approximation with a prescribed accuracy $\epsilon > 0$, let the costs of achieving this approximation scale as $\mathcal{O}(\epsilon^{-d})$ in d dimensions \rightsquigarrow exponential increase of the costs as we increase d
- ▶ Consider a simple uniform grid over the domain $\Omega = [0, 1]^d$. To have a mesh with mesh width $h = 1/9$ in $d = 1$, we need $N = 10$ grid points. In $d = 2$ dimensions, need $N^2 = 100$ grid points. In $d = 5$ dimensions, need $N^5 = 10^5 \rightsquigarrow$ exponential growth of cost and storage requirements as we increase dimension d while keeping mesh width h (“accuracy”) fixed
- ▶ The curse can be circumvented (to some extent) **if?** if we make stronger assumptions on the functions to approximate \rightsquigarrow topic of today

Sparse grids*

*Follows lecture by H.-J. Bungartz. See also, Bungartz, Griebel, *Sparse grids*, Acta Numerica, 2004

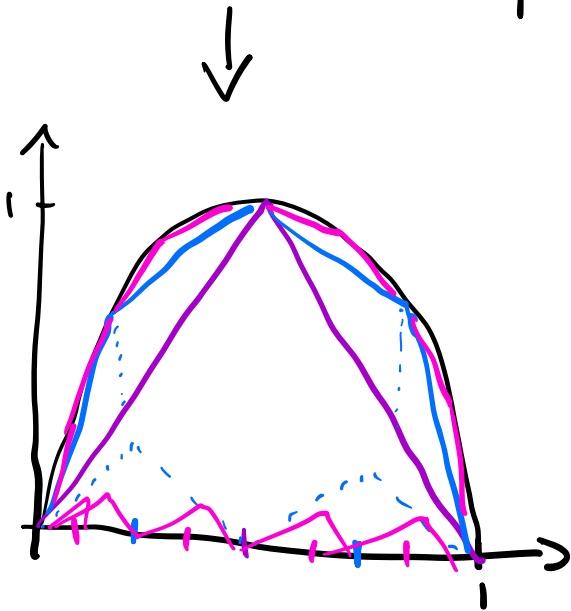
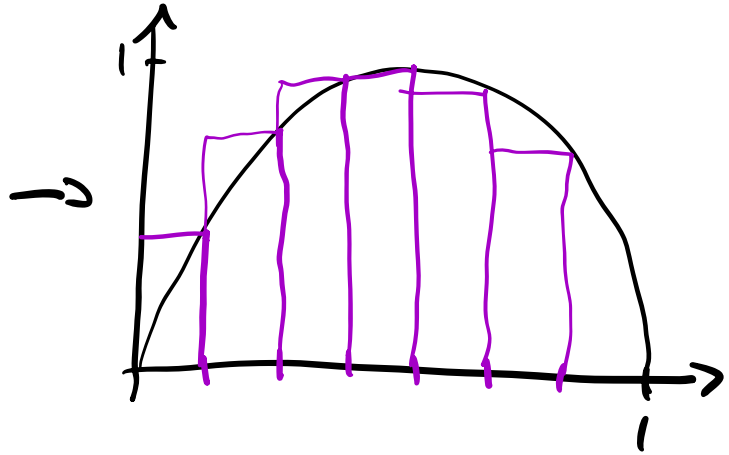
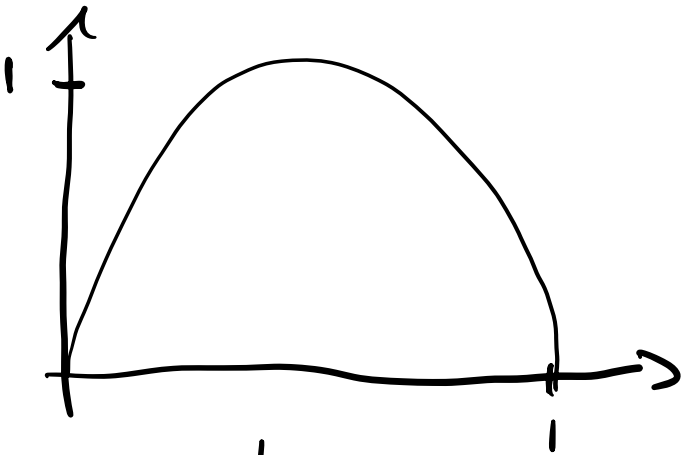
SG: Motivating example

Approximate the integral

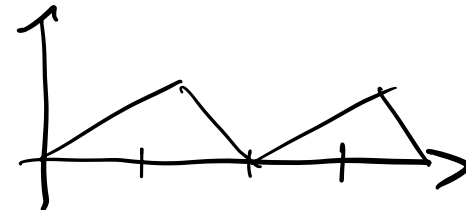
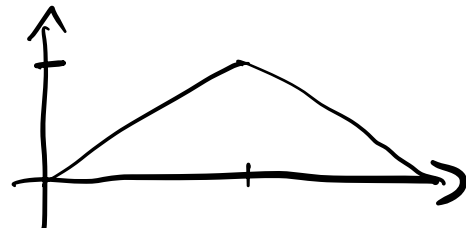
$$\int_0^1 4x(1-x)dx = \frac{2}{3}$$

Board!

$$\int_0^1 4x(1-x) dx = \frac{2}{3}$$



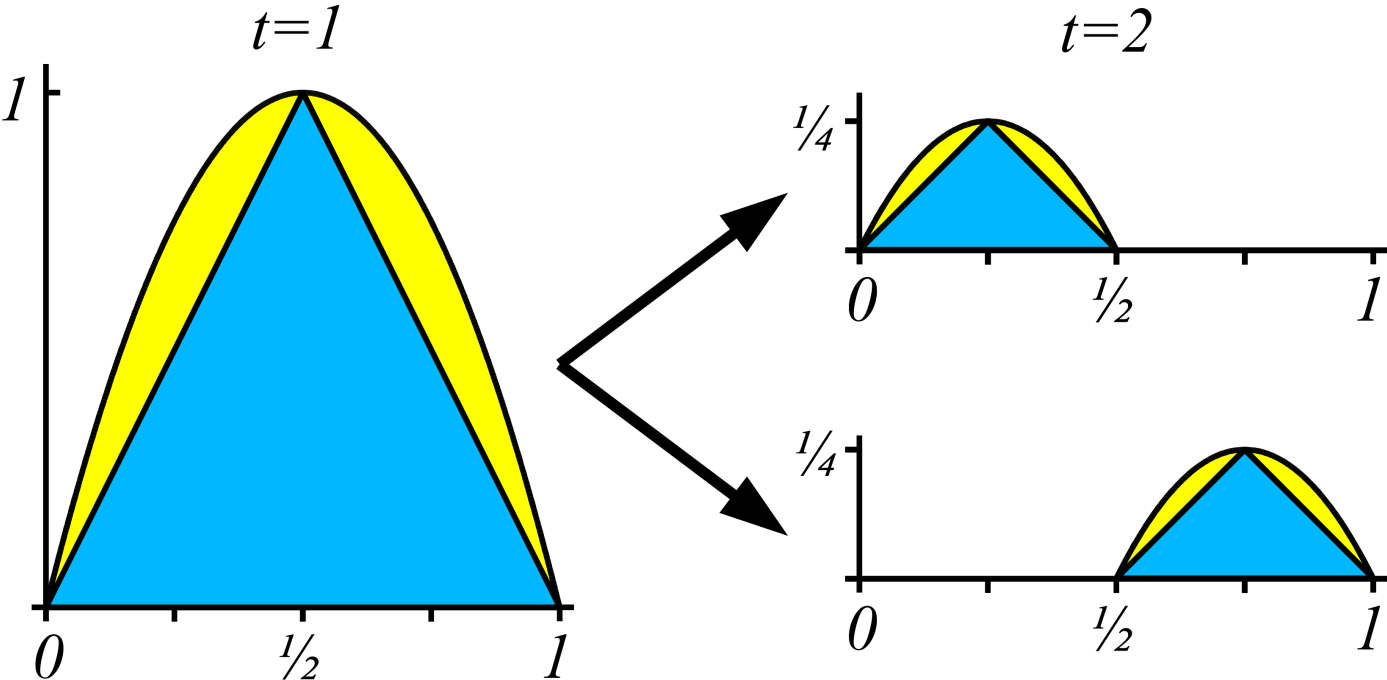
multi-level



SG: Hierarchical decomposition 1D

Archimedes Quadrature

$$\int_0^1 4x(1-x)dx = \frac{2}{3}$$



depth	1	2	3	4	...	t
interval length h	1/2	1/4	1/8	1/16	...	2^{-t}
number of triangles	1	2	4	8	...	$\frac{1}{2}2^t$
surplus v	1	1/4	1/16	1/64	...	$4 \cdot 2^{-2t}$
triangle area D_1	1/2	1/16	1/128	1/1024	...	$4 \cdot 2^{-3t}$
sum (of this t)	1/2	1/8	1/32	1/128	...	$2 \cdot 2^{-2t}$
overall sum ($\leq t$)	1/2	5/8	21/32	85/128	...	$\frac{2}{3} (1 - 2^{-2t})$
error	1/6	1/24	1/96	1/384	...	$\frac{2}{3}2^{-2t}$

What do we observe with respect to depth t ?

depth	1	2	3	4	...	t
interval length h	1/2	1/4	1/8	1/16	...	2^{-t}
number of triangles	1	2	4	8	...	$\frac{1}{2}2^t$
surplus v	1	1/4	1/16	1/64	...	$4 \cdot 2^{-2t}$
triangle area D_1	1/2	1/16	1/128	1/1024	...	$4 \cdot 2^{-3t}$
sum (of this t)	1/2	1/8	1/32	1/128	...	$2 \cdot 2^{-2t}$
overall sum ($\leq t$)	1/2	5/8	21/32	85/128	...	$\frac{2}{3} (1 - 2^{-2t})$
error	1/6	1/24	1/96	1/384	...	$\frac{2}{3}2^{-2t}$

What do we observe with respect to depth t ? Contribution (e.g., “sum (of this t)”) goes down with rate 2^{-2t}

SG: Approximation of functions

Analyze this approach in more general context.

Let ϕ_1, \dots, ϕ_n be basis functions and represent

$$u(x) = \sum_{i=1}^n \alpha_i \phi_i(x)$$

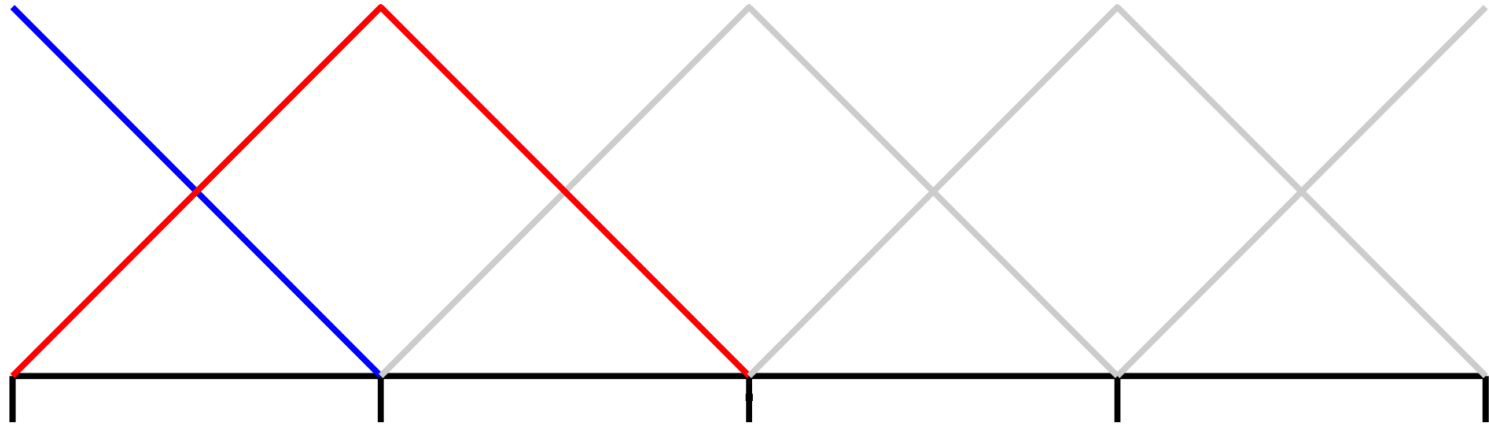
Obtain

$$\int_a^b u(x) dx = \sum_{i=1}^n \alpha_i \int_a^b \phi_i(x) dx,$$

i.e., a quadrature rule as a weighted sum of the coefficients $\alpha_1, \dots, \alpha_n$

SG: Approximation of functions (cond't)

Represent a continuous, piecewise linear u in nodal point basis



The coefficient $\alpha_1, \dots, \alpha_n$ are the function values of u at the nodal points

$$u(x) = \sum_{i=1}^n \alpha_i \phi_i(x)$$

\Rightarrow instead of nodal point basis, consider a **hierarchical** basis

SG: Piecewise linear functions

Consider only functions $u : [0, 1] \rightarrow \mathbb{R}$ with $u(0) = u(1) = 0$ in the following.

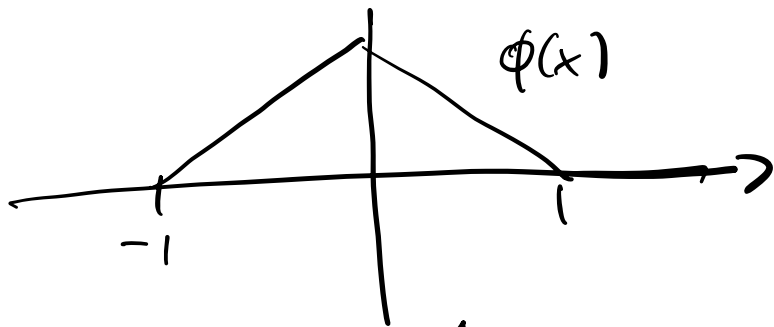
Need the following quantities

- ▶ mesh width $h_l = 2^{-l}$
- ▶ grid points $x_{l,i} = ih_l = i2^{-l}$
- ▶ Basis function

$$\phi_{l,i}(x) = \phi\left(\frac{x - x_{l,i}}{h_l}\right), \quad \phi(x) = \max\{1 - |x|, 0\}$$

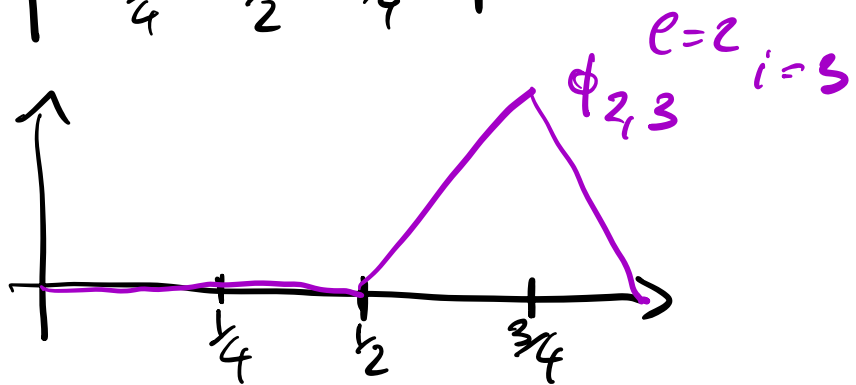
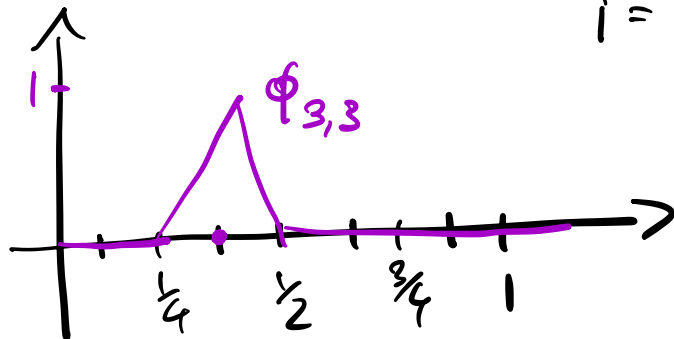
- ▶ Nodal point basis $\Phi_l = \{\phi_{l,i} : 1 \leq i < 2^l\}$

Visualize these basis functions on the board!



$$l = 3$$

$$i = 3$$



SG: Piecewise linear functions

Consider only functions $u : [0, 1] \rightarrow \mathbb{R}$ with $u(0) = u(1) = 0$ in the following.

Need the following quantities

- ▶ mesh width $h_l = 2^{-l}$
- ▶ grid points $x_{l,i} = ih_l = i2^{-l}$
- ▶ Basis function

$$\phi_{l,i}(x) = \phi\left(\frac{x - x_{l,i}}{h_l}\right), \quad \phi(x) = \max\{1 - |x|, 0\}$$

- ▶ Nodal point basis $\Phi_l = \{\phi_{l,i} : 1 \leq i < 2^l\}$

Visualize these basis functions on the board!

The space $V_l = \text{span}(\Phi_l)$ is the space of piecewise linear, continuous functions with respect to the grid points $x_{l,i}$ for $i = 1, \dots, 2^l - 1$

SG: Hierarchical basis

For the hierarchical representation, we consider the hierarchical increment W_l , spanned by the basis functions $\phi_{l,i}$ such that

$$V_l = V_{l-1} \oplus W_l$$

is a direct sum (each $u_l \in V_l$ can be uniquely decomposed as $u_l = u_{l-1} + w_l$ with $u_{l-1} \in V_{l-1}$ and $w_l \in W_l \rightarrow$ remember the triangles in approach by Archimedes)

Because $\dim(V_l) = 2^l - 1$ and $\dim(V_l) = \dim(V_{l-1}) + \dim(W_l)$ we need $\dim(W_l) = 2^{l-1}$. These are given by $\phi_{l,i}$ with

$$i \in I_l = \{j : 1 \leq j < 2^l, \quad j \text{ odd}\}$$

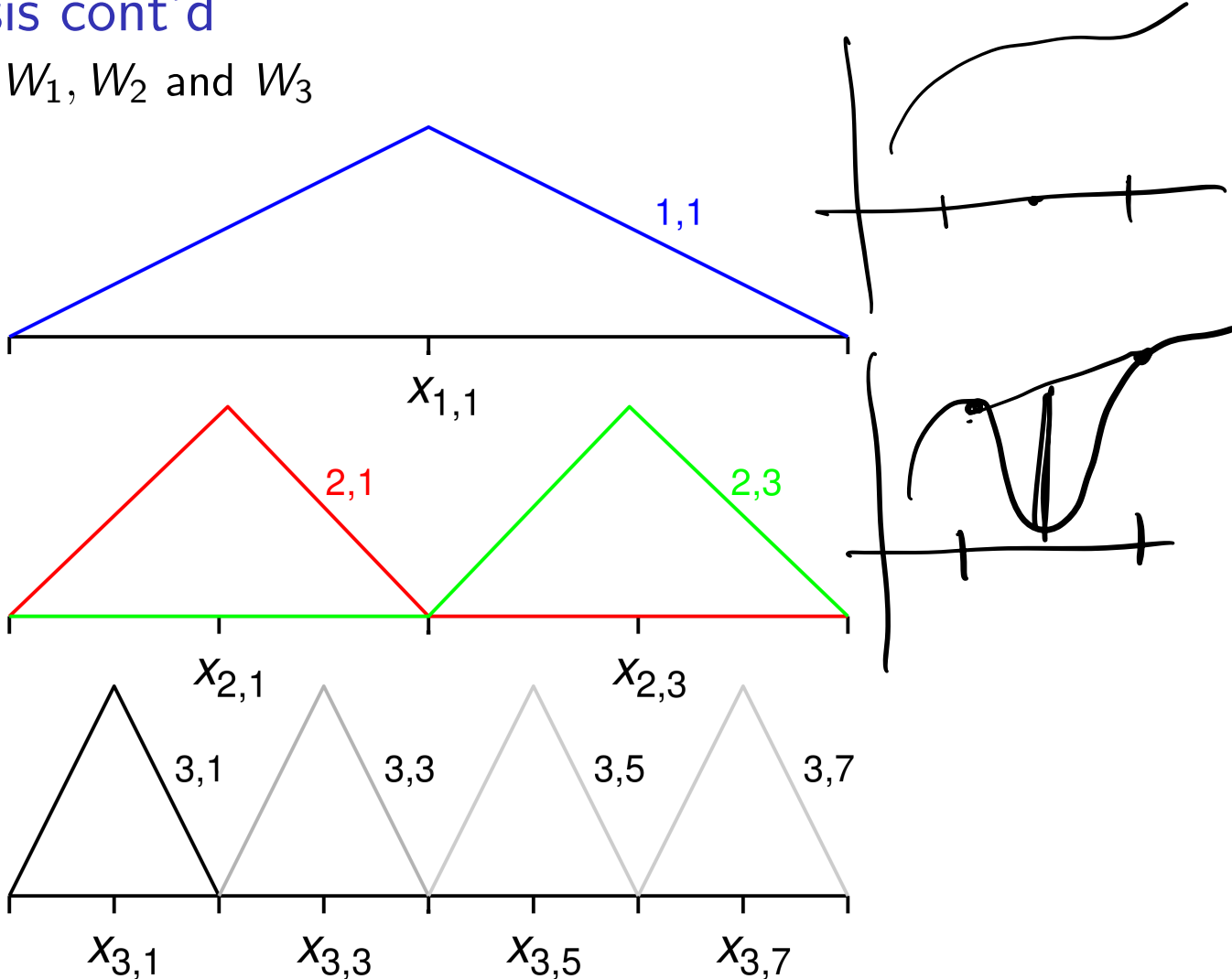
Then

$$W_l = \text{span}\{\phi_{l,i}, i \in I_l\}$$

Note that $W_1 = V_1$

SG: Hierarchical basis cont'd

The bases for spaces W_1, W_2 and W_3



SG: Hierarchical basis cont'd

Obtain

$$V_n = \bigoplus_{l=1}^n W_l,$$

so that there is a unique representation for each $u \in V_n$ as

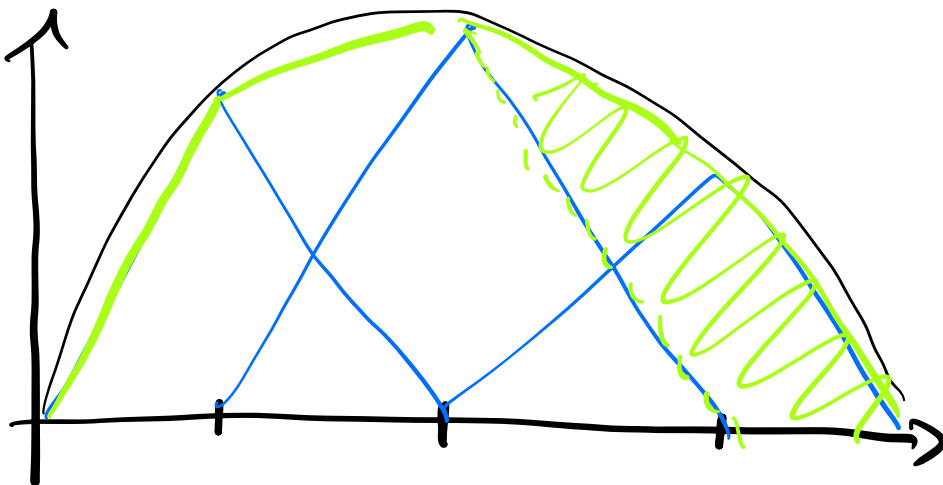
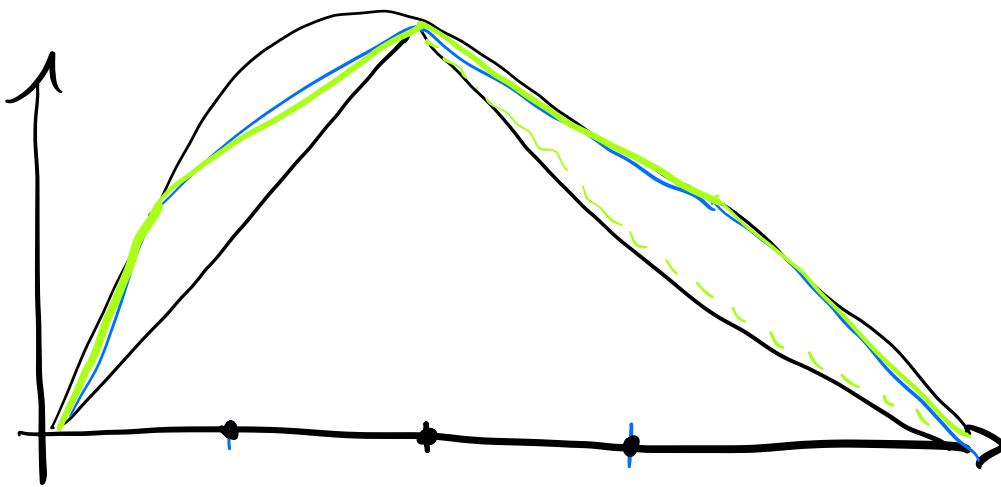
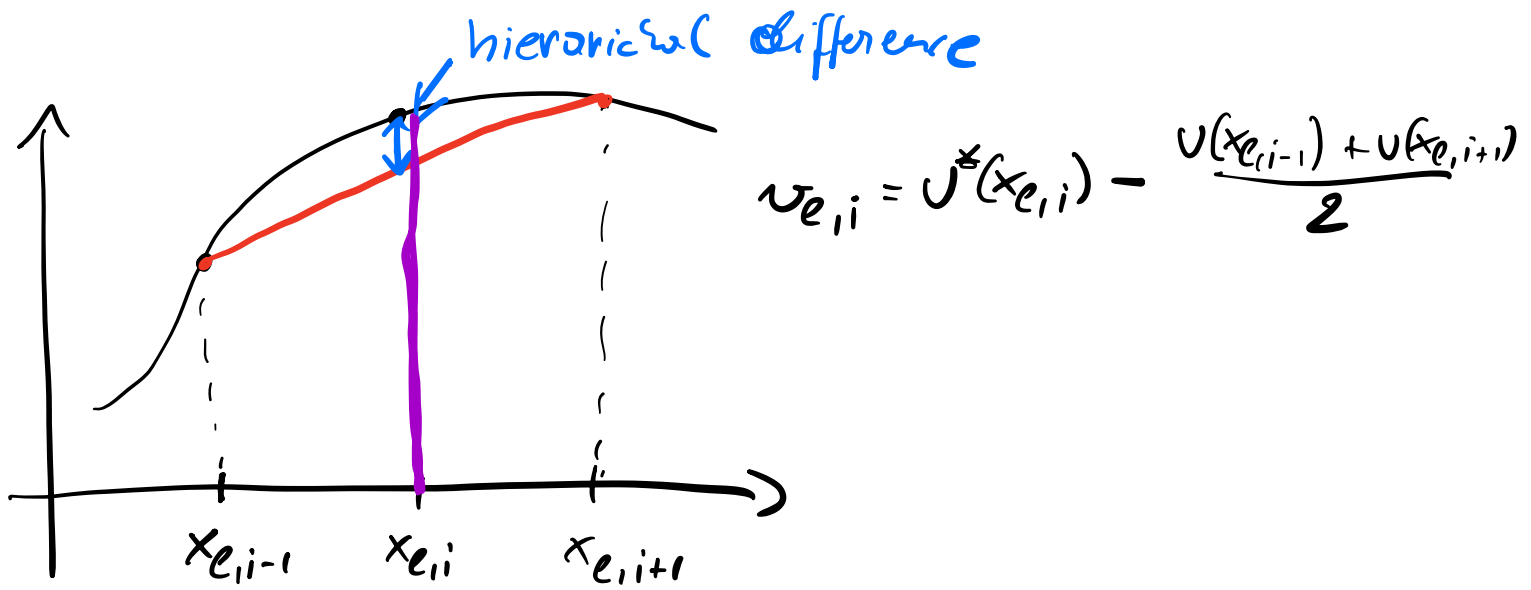
$$u = \sum_{l=1}^n w_l = \sum_{l=1}^n \sum_{i \in I_l} v_{l,i} \phi_{l,i}$$

Coefficients $v_{l,i}$ in case of interpolation \rightsquigarrow visualize on board

The coefficients $v_{l,i}$ are hierarchical differences

$$v_{l,i} = u^*(x_{l,i}) - \frac{u(x_{l,i-1}) + u(x_{l,i+1})}{2}$$

where u^* is the function to be interpolated



SG: Analysis of hierarchical decomposition

We now analyze the hierarchical decomposition

$$u_n = \sum_{l=1}^n w_l = \sum_{l=1}^n \sum_{i \in I_l} v_{l,i} \phi_{l,i}$$

when interpolating a u . **What would we like to obtain?**

SG: Analysis of hierarchical decomposition

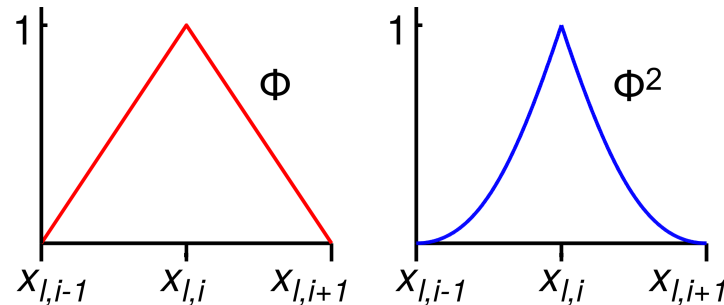
We now analyze the hierarchical decomposition

$$u_n = \sum_{l=1}^n w_l = \sum_{l=1}^n \sum_{i \in I_l} v_{l,i} \phi_{l,i}$$

when interpolating a u . **What would we like to obtain?**

For our decomposition, we first calculate the norms of the basis functions

$$\|\phi_{l,i}\|_{\infty} = 1, \quad \|\phi_{l,i}\|_2 = \sqrt{\frac{2h_l}{3}}$$



Show decay rate of

$$w_e = \sum_{i \in I_e} v_{e,i} \phi_{e,i}$$

Represent coefficients as [BG, Lemma 3.2]

$$v_{e,i} = \int_0^1 \psi_{e,i}(x) \frac{\partial^2}{\partial x^2} u(x) dx$$

$$\hookrightarrow \psi_{e,i} = -\frac{h_e}{2} \phi_{e,i}$$

\hookrightarrow requires twice differentiable u !!!
 \Rightarrow Strong smoothness assumption

Now bound coefficients

$$\begin{aligned} \underline{|v_{e,i}|} &\leq \frac{h_e}{2} \underbrace{\|\phi_{e,i}\|_2}_{\leq \sqrt{\frac{3}{h_e}}} \left\| \frac{\partial^2}{\partial x^2} u|_{T_i} \right\|_2 \\ &\leq \sqrt{\frac{h_e^3}{6}} \left\| \frac{\partial^2}{\partial x^2} u|_{T_i} \right\|_2 \end{aligned}$$

where $u|_{T_i}$ is u restricted to support

$$T_i = [x_{e,i-1}, x_{e,i+1}] \text{ of } \phi_{e,i}$$

Numerical Methods I

MATH-GA 2010.001 / CSCI-GA 2420.001

Benjamin Peherstorfer
Courant Institute, NYU

Today •

Last time

- ▶ Quadrature in higher dimensions
- ▶ Sparse grids

Today

- ▶ Sparse grids

Announcements

- ▶ Homework 6 is due Mon, Dec 2 before class
- ▶ Wed, Nov 27 will be a Q&A session. No new material will be discussed. All notes will be posted online. Send me questions via email by Tue, Nov 26, 4.55pm.

Recap: Curse of dimensionality

- ▶ The curse of dimensionality is a term coined by Bellmann (1961) that refers to an exponential increase of costs with the dimension of a problem
- ▶ For example, consider an approximation with a prescribed accuracy $\epsilon > 0$, let the costs of achieving this approximation scale as $\mathcal{O}(\epsilon^{-d})$ in d dimensions \rightsquigarrow exponential increase of the costs as we increase d
- ▶ Consider a simple uniform grid over the domain $\Omega = [0, 1]^d$. To have a mesh with mesh width $h = 1/9$ in $d = 1$, we need $N = 10$ grid points. In $d = 2$ dimensions, need $N^2 = 100$ grid points. In $d = 5$ dimensions, need $N^5 = 10^5 \rightsquigarrow$ exponential growth of cost and storage requirements as we increase dimension d while keeping mesh width h (“accuracy”) fixed

Recap: Curse of dimensionality

- ▶ The curse of dimensionality is a term coined by Bellmann (1961) that refers to an exponential increase of costs with the dimension of a problem
- ▶ For example, consider an approximation with a prescribed accuracy $\epsilon > 0$, let the costs of achieving this approximation scale as $\mathcal{O}(\epsilon^{-d})$ in d dimensions \rightsquigarrow exponential increase of the costs as we increase d
- ▶ Consider a simple uniform grid over the domain $\Omega = [0, 1]^d$. To have a mesh with mesh width $h = 1/9$ in $d = 1$, we need $N = 10$ grid points. In $d = 2$ dimensions, need $N^2 = 100$ grid points. In $d = 5$ dimensions, need $N^5 = 10^5 \rightsquigarrow$ exponential growth of cost and storage requirements as we increase dimension d while keeping mesh width h (“accuracy”) fixed
- ▶ The curse can be circumvented (to some extent) **if?**

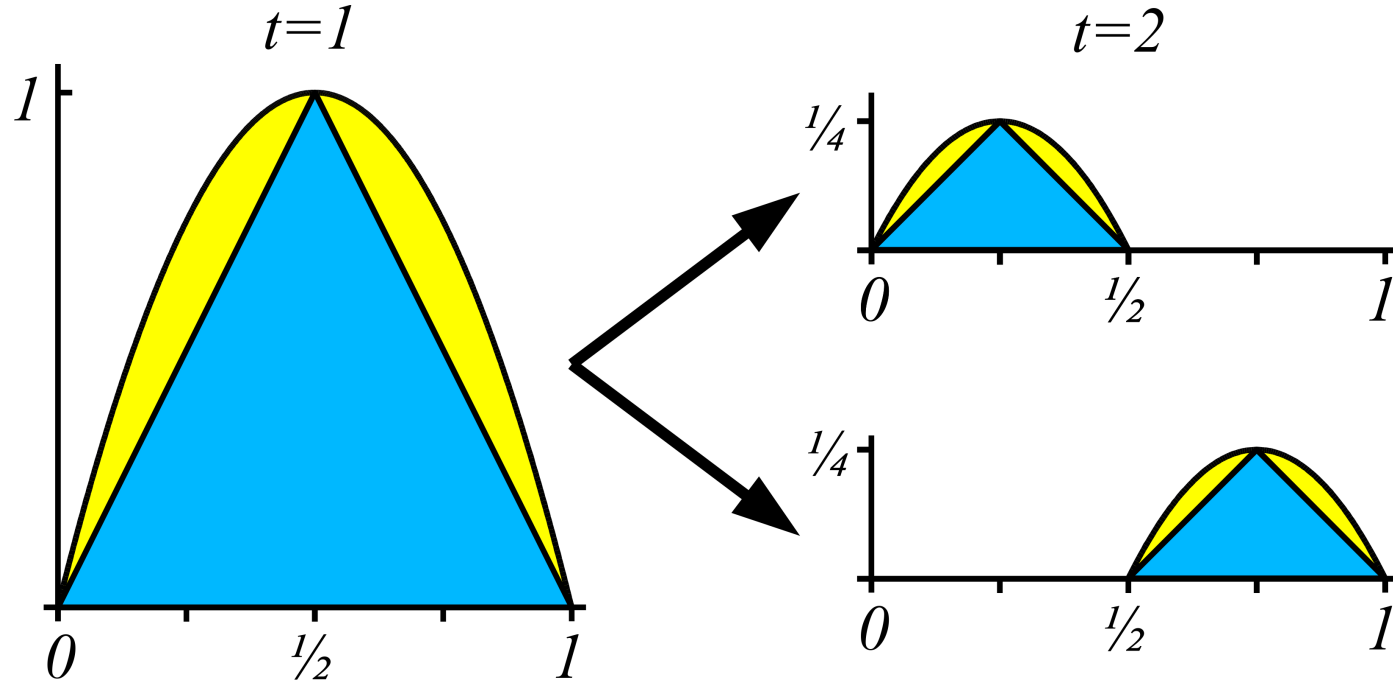
Recap: Curse of dimensionality

- ▶ The curse of dimensionality is a term coined by Bellmann (1961) that refers to an exponential increase of costs with the dimension of a problem
- ▶ For example, consider an approximation with a prescribed accuracy $\epsilon > 0$, let the costs of achieving this approximation scale as $\mathcal{O}(\epsilon^{-d})$ in d dimensions \rightsquigarrow exponential increase of the costs as we increase d
- ▶ Consider a simple uniform grid over the domain $\Omega = [0, 1]^d$. To have a mesh with mesh width $h = 1/9$ in $d = 1$, we need $N = 10$ grid points. In $d = 2$ dimensions, need $N^2 = 100$ grid points. In $d = 5$ dimensions, need $N^5 = 10^5 \rightsquigarrow$ exponential growth of cost and storage requirements as we increase dimension d while keeping mesh width h (“accuracy”) fixed
- ▶ The curse can be circumvented (to some extent) **if?** if we make stronger assumptions on the functions to approximate \rightsquigarrow topic of today

Recap: SG: Hierarchical decomposition 1D

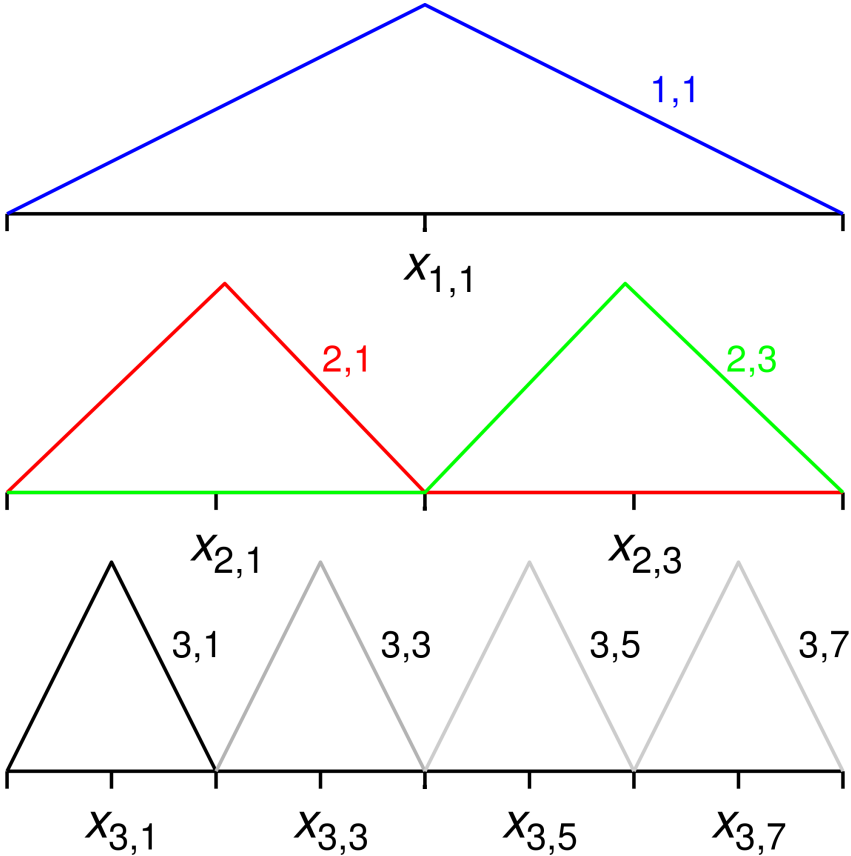
Archimedes Quadrature

$$\int_0^1 4x(1-x)dx = \frac{2}{3}$$



SG: Hierarchical basis cont'd

The bases for spaces W_1 , W_2 and W_3



SG: Hierarchical basis

For the hierarchical representation, we consider the hierarchical increment W_l , spanned by the basis functions $\phi_{l,i}$ such that

$$V_l = V_{l-1} \oplus W_l$$

is a direct sum (each $u_l \in V_l$ can be uniquely decomposed as $u_l = u_{l-1} + w_l$ with $u_{l-1} \in V_{l-1}$ and $w_l \in W_l \rightarrow$ remember the triangles in approach by Archimedes)

Because $\dim(V_l) = 2^l - 1$ and $\dim(V_l) = \dim(V_{l-1}) + \dim(W_l)$ we need $\dim(W_l) = 2^{l-1}$. These are given by $\phi_{l,i}$ with

$$i \in I_l = \{j : 1 \leq j < 2^l, \quad j \text{ odd}\}$$

Then

$$W_l = \text{span}\{\phi_{l,i}, i \in I_l\}$$

Note that $W_1 = V_1$

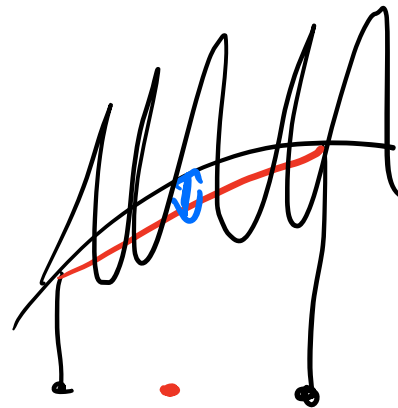
SG: Hierarchical basis cont'd

Obtain

$$V_n = \bigoplus_{l=1}^n W_l,$$

so that there is a unique representation for each $u \in V_n$ as

$$u = \sum_{l=1}^n w_l = \sum_{l=1}^n \sum_{i \in I_l} v_{l,i} \phi_{l,i}$$



Coefficients $v_{l,i}$ in case of interpolation \rightsquigarrow visualize on board

The coefficients $v_{l,i}$ are hierarchical differences

$$v_{l,i} = u^*(x_{l,i}) - \frac{u(x_{l,i-1}) + u(x_{l,i+1})}{2}$$

where u^* is the function to be interpolated

SG: Analysis of hierarchical decomposition

We now analyze the hierarchical decomposition

$$u_n = \sum_{l=1}^n w_l = \sum_{l=1}^n \sum_{i \in I_l} v_{l,i} \phi_{l,i}$$

when interpolating a u . **What would we like to obtain?**

SG: Analysis of hierarchical decomposition

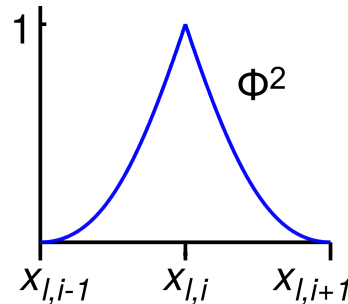
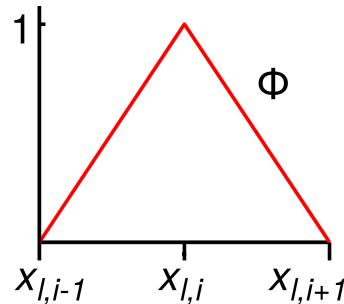
We now analyze the hierarchical decomposition

$$u_n = \sum_{l=1}^n w_l = \sum_{l=1}^n \sum_{i \in I_l} v_{l,i} \phi_{l,i}$$

when interpolating a u . **What would we like to obtain?**

For our decomposition, we first calculate the norms of the basis functions

$$\|\phi_{l,i}\|_{\infty} = 1, \quad \|\phi_{l,i}\|_2 = \sqrt{\frac{2h_l}{3}}$$



Show decay of

$$w_e = \sum_{i \in I_e} v_{e,i} \phi_{e,i}$$

Represent hierarchical coefficient as

$$v_{e,i} = \int_0^1 \psi_{e,i}(x) \frac{\partial^2}{\partial x^2} u(x) dx \quad \left. \vphantom{\int_0^1} \right\} \text{Lemma 3.2}$$

$$\psi_{e,i} = -\frac{h_e}{2} \phi_{e,i}$$

→ requires twice differentiable u

⇒ strong smoothness assumption

$$|v_{e,i}| \leq \frac{h_e}{2} \|\phi_{e,i}\|_2 \left\| \frac{\partial^2}{\partial x^2} u|_{T_i} \right\|_2 \quad T_i = [x_{e,i-1}, x_{e,i+1}]$$

$$\leq \sqrt{\frac{h_e^3}{6}} \underbrace{\left\| \frac{\partial^2}{\partial x^2} u|_{T_i} \right\|_2}$$

⇒ requires bounded 2nd derivatives

$$\text{define } \mu_2(u|_{T_i}) = \left\| \frac{\partial^2}{\partial x^2} u|_{T_i} \right\|_2$$

Now increment

$$w_e = \sum_{i \in I_e} v_{e,i} \phi_{e,i}$$

$$\begin{aligned}
\|w_e\|_2^2 &= \int_0^1 w_e^2(x) dx = \\
&= \int_0^1 \left(\sum_{i \in I_e} v_{ei} \phi_{ei}(x) \right)^2 dx \\
&\leq \sum_{i \in I_e} v_{ei}^2 \|\phi_{ei}\|_2^2 \\
&\leq \frac{\hbar_e^3}{6} \frac{2\hbar_e}{3} \sum_{i \in I_e} \mu_2(v(T_i))^2 \\
&= \frac{\hbar_e^4}{9} \mu_2(v)^2
\end{aligned}$$

$$\Rightarrow \|w_e\|_2 \in \mathcal{O}(\hbar_e^2)$$

\Rightarrow ^{norm of} increments decay quadratically with level e
if v bounded 2nd-order derivatives

\Rightarrow Represent v as series

$$v = \sum_{e=1}^{\infty} w_e$$

Recall

$$u_n = \sum_{l=1}^n w_l = \sum_{l=1}^n \sum_{i \in I_l} v_{l,i} \phi_{l,i}$$

If u is twice differentiable, the $\|\cdot\|_2$ norm of the increments

$$w_l = \sum_{i \in I_l} v_{l,i} \phi_{l,i}$$

decays as $\mathcal{O}(h_l^2) \rightsquigarrow$ **board**

Recall

$$u_n = \sum_{l=1}^n w_l = \sum_{l=1}^n \sum_{i \in I_l} v_{l,i} \phi_{l,i}$$

If u is twice differentiable, the $\|\cdot\|_2$ norm of the increments

$$w_l = \sum_{i \in I_l} v_{l,i} \phi_{l,i}$$

decays as $\mathcal{O}(h_l^2) \rightsquigarrow$ **board**

This means, we can write a twice differentiable u as a series

$$u = \sum_{l=1}^{\infty} w_l,$$

that converges because $\|w_l\|_2 \in \mathcal{O}(h_l^2)$. In particular, obtain

$$u - u_n = \sum_{l=n+1}^{\infty} w_l$$

\rightsquigarrow decay of $\|w_l\|_2$ helps us to understand error $u - u_n$

In the 1-dimensional case, the benefit of the hierarchical basis is limited:

In the 1-dimensional case, the benefit of the hierarchical basis is limited:

- ▶ We can leave out nodes and still get a reasonable approximation in contrast to a nodal-point basis

In the 1-dimensional case, the benefit of the hierarchical basis is limited:

- ▶ We can leave out nodes and still get a reasonable approximation in contrast to a nodal-point basis
- ▶ Adding a new level requires us to compute the coefficients of the new level only but keeps the coefficients at the previous levels unchanged

What will “can remove nodes and still get reasonable results” amount to in higher dimension?

In the 1-dimensional case, the benefit of the hierarchical basis is limited:

- ▶ We can leave out nodes and still get a reasonable approximation in contrast to a nodal-point basis
- ▶ Adding a new level requires us to compute the coefficients of the new level only but keeps the coefficients at the previous levels unchanged

What will “can remove nodes and still get reasonable results” amount to in higher dimension? In the multi-dimensional case, we will now see that many of the hierarchical increments have high costs and low benefit in terms of error → we will then remove those hierarchical increments and obtain *sparse grids*

SG: Hierarchical basis in multivariate case

Let now $\mathbf{x} = [x_1, \dots, x_d]$ with $d \in \mathbb{N}$ and $d > 1$.

Again consider the domain $\Omega = [0, 1]^d$ and functions $u|_{\partial\Omega} = 0$

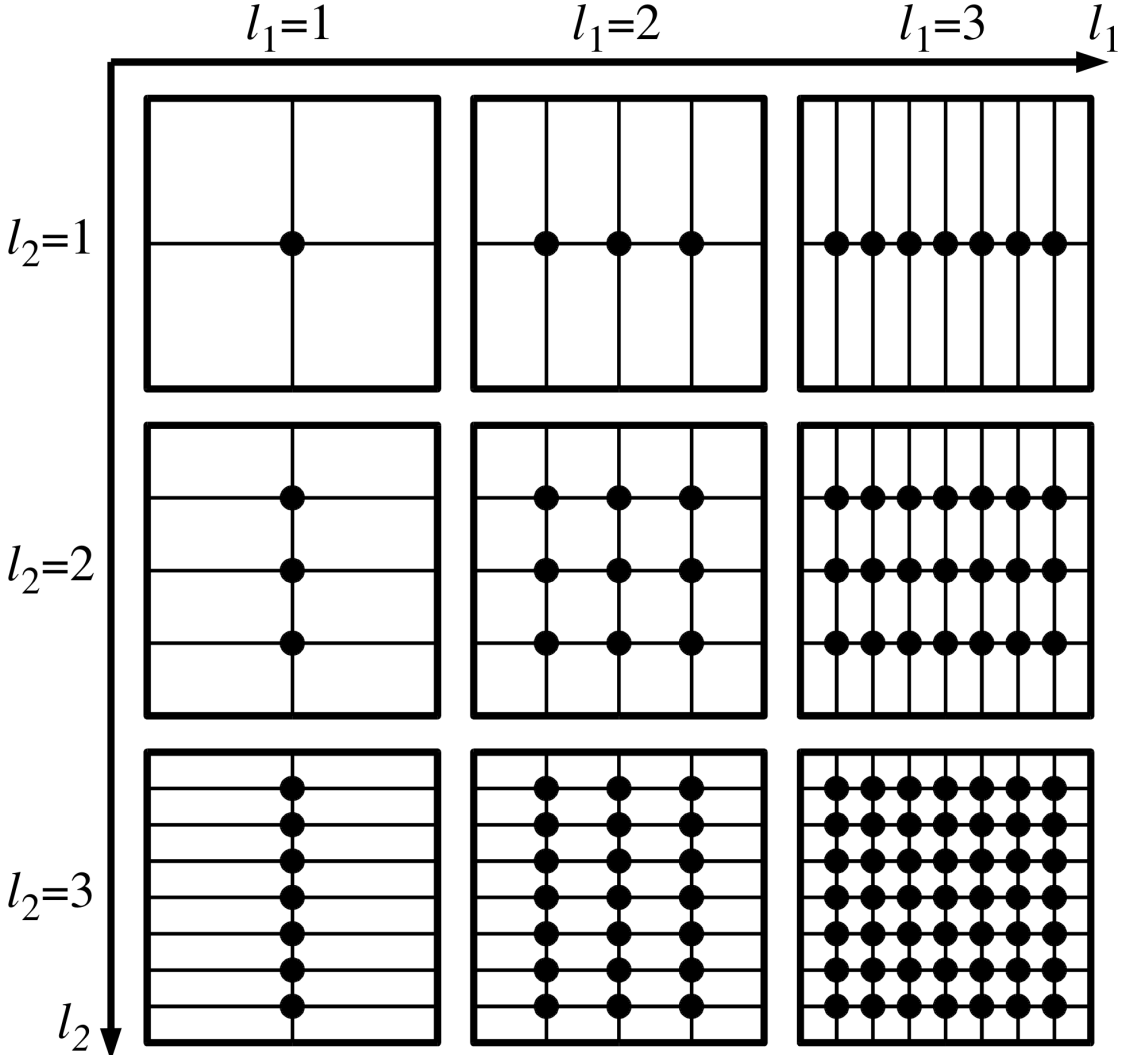
Multidimensional level $\mathbf{l} = [l_1, \dots, l_d] \in \mathbb{N}^d$

Multidimensional mesh width $h_{\mathbf{l}} = [h_1, \dots, h_d] = [2^{-l_1}, \dots, 2^{-l_d}]$. Note that different mesh width in different dimensions allowed

Define $|\mathbf{l}|_1 = l_1 + \dots + l_d$ and $|\mathbf{l}|_{\infty} = \max\{l_1, \dots, l_d\}$

Grid points are $\mathbf{x}_{\mathbf{l}, \mathbf{i}} = [i_1 h_1, \dots, i_d h_d]$

SG: Notation

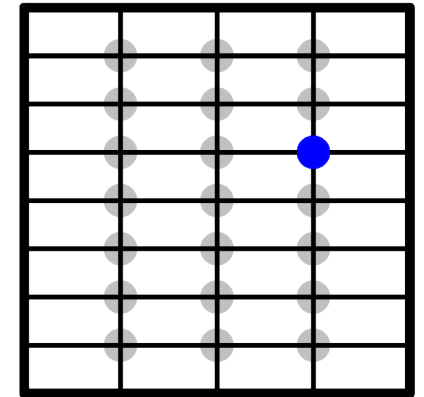
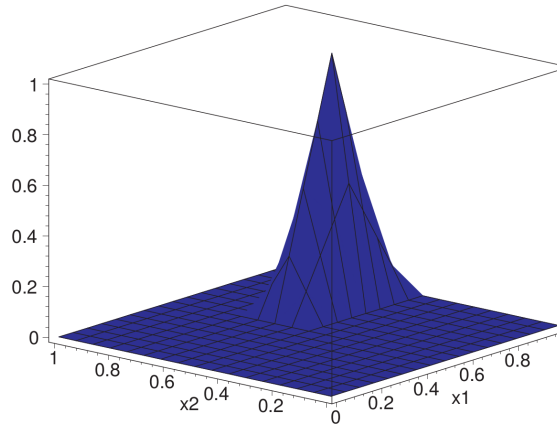
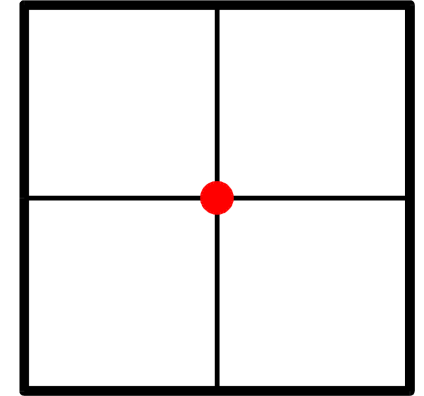
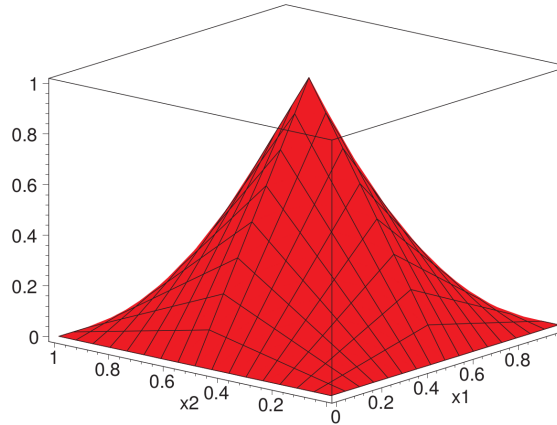


SG: Piecewise d -linear functions

Generalize continuous, piecewise linear functions to continuous, piecewise d -linear functions with respect to h_I :

$$\phi_{I,i}(\mathbf{x}) = \prod_{j=1}^d \phi_{I_j,i_j}(x_j)$$

For $d = 2$, the functions $\phi_{[1,1],[1,1]}$ and $\phi_{[2,3],[3,5]}$ are plotted on the right



SG: Spaces

Consider

$$\Phi_I = \{\phi_{I,i} : 1 \leq i < 2^I\},$$

where the \leq is to be read component-wise: each i_j must be at least 1 and at most $2^{l_j} - 1$

The space of piecewise d -linear functions is

$$V_I = \text{span}\{\Phi_I\}$$

with dimension

$$\dim(V_I) = (2^{l_1} - 1) \cdots (2^{l_d} - 1) \in \mathcal{O}(2^{\|I\|_1})$$

Special case $l_1 = \cdots = l_d$ set $V_n = V_{[n, \dots, n]}$

SG: Hierarchical increments

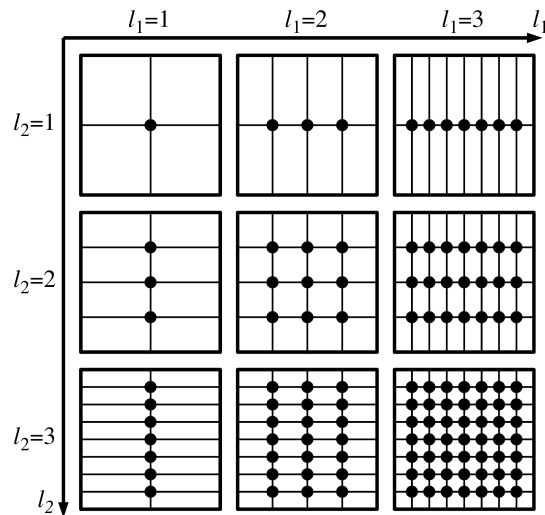
Define the hierarchical increment W_l as

$$W_l = \text{span}\{\phi_{l,i} : i \in I_l\},$$

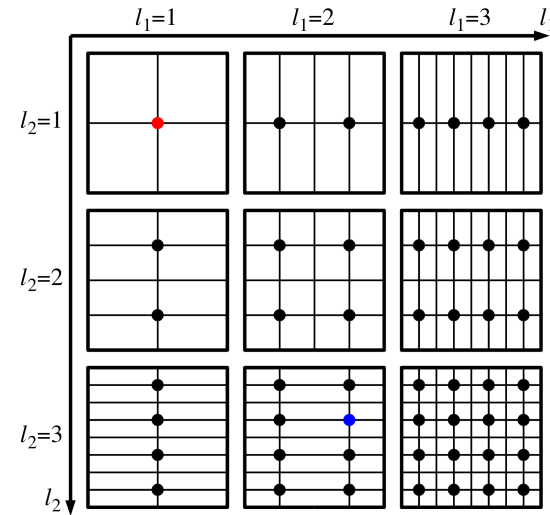
where $I_l = \{i : 1 \leq i < 2^l, \text{ all } i_j \text{ odd}\}$

Contains just those functions from V_l that vanish at all points of coarser grid

full grid



hierarchical increments



SG: Hierarchical subspace decomposition

Obtain unique representation of $u_I \in V_I$ for $I \in \mathbb{N}^d$ as

$$u_I = \sum_{I' \leq I} w_{I'} ,$$

with $w_{I'} \in W_{I'}$

Now it will be worthwhile to estimate the norm of w_I to understand which contribute most to the accuracy of the representation

↪ **board**

Multi-dimensional case:

Represent $U \in V_e \subset \mathbb{C} \in \mathbb{N}^d$ as

$$U = \sum_{e' \leq e} w_{e'} \quad , \quad w_{e'} \in W_{e'}$$

\Rightarrow estimate norm of $w_{e'}$

Hierarchical coefficient

$$w_{e'} = \int_{[0,1]^d} \psi_{e'} \partial^{2e'} U \, dx \quad , \quad \psi_{e'} = 2^{-|e'| - d} \phi_{e'}$$

$$\partial^{2e'} U = \frac{\partial^{2e'} U}{\partial x_1^2 \dots \partial x_d^2} \quad \text{mixed } 2d\text{-fold derivative}$$

\Rightarrow STRONG assumption

Obtain

$$\underline{\underline{\|w_e\|_2}} \leq 3^{-d} \underbrace{2^{-2|e|}}_{\text{decays fast with level } e} \|\partial^{2e} U\|_2$$

\Rightarrow understand cost/benefit of each of these increments w_e

costs: number of grid points of w_e

$$c(e) = 2^{|\mathcal{E}| - d}$$

benefit: Let $L \subseteq \mathcal{N}^d$ be a set of levels that are selected

$$U_L = \sum_{e \in L} w_e$$

then

$$U - U_L = \sum_{e \notin L} w_e$$

For each component have bound

$$\|w_e\|_2 \leq s(e) \mu(v)$$

$$\|U - U_L\|_2 \leq \sum_{e \notin L} \|w_e\|_2 \leq \left(\sum_{e \notin L} s(e) \right) \mu(v)$$

$$= \left[\sum_{e \in \mathcal{N}^d} s(e) - \sum_{e \in L} s(e) \right] \mu(v)$$

if select e with $s(e)$ big, then error

$\|U - U_L\|_2$ reduced

$\Rightarrow s(e)$ is benefit of w_e

With this new tool, let us analyze the cost/benefit of a "full grid",

$$L_n = \{e : |e| \leq h\}$$

Then calculations show

$$\sum_{e \in L_n} s(e) \geq \left(\frac{1}{3}\right)^d (1 - d 2^{-2n})$$

for $n \rightarrow \infty$

$$\sum_{e \in \mathbb{N}^d} s(e) = \left(\frac{1}{3}\right)^d$$

$$\begin{aligned} \sum_{e \in \mathbb{N}^d} s(e) - \sum_{e \in L_n} s(e) &\leq \left(\frac{1}{3}\right)^d - \left(\frac{1}{3}\right)^d (1 - d 2^{-2n}) \\ &\leq \frac{1}{3^d} d 2^{-2n} \end{aligned}$$

so that

$$\|u - u_L\|_2 \leq C \sum_{e \in L_n} s(e) \leq C \frac{d}{3^d} 2^{-2n} \in \mathcal{O}(h_n^2)$$

Details in [Bungartz et al., 2004]: Hierarchical coefficient is now

$$v_{l,i} = \int_{[0,1]^d} \psi_{l,i} \partial^{2d} u d\mathbf{x},$$

which now depends on **mixed** 2d-fold derivative

$$\partial^{2d} u = \frac{\partial^{2d}}{\partial x_1^2 \cdots \partial x_d^2}$$

and $\psi_{l,i} = 2^{-|l|_1 - d} \phi_{l,i}$

It is very important to note that the following holds *only* for functions with (L_2 -)bounded $\partial^{2d} u \Rightarrow$ strong assumption on function (smoothness)

Obtain

$$\|w_l\|_2 \leq 3^{-d} 2^{-2|l|_1} \|\partial^{2d} u\|_2$$

Now we select those subspaces from the subspace scheme that minimize cost and maximize benefit for approximation function $u : [0, 1]^d \rightarrow \mathbb{R}$ with sufficient smoothness

How should we measure costs?

Now we select those subspaces from the subspace scheme that minimize cost and maximize benefit for approximation function $u : [0, 1]^d \rightarrow \mathbb{R}$ with sufficient smoothness

How should we measure costs? We measure cost via the number of grid points

$$c(I) = 2^{|I|_1 - d}$$

How should we measure benefit?

Now we select those subspaces from the subspace scheme that minimize cost and maximize benefit for approximation function $u : [0, 1]^d \rightarrow \mathbb{R}$ with sufficient smoothness

How should we measure costs? We measure cost via the number of grid points

$$c(I) = 2^{|I|_1 - d}$$

How should we measure benefit? Measure benefit of subspace selection via introduced error if left out. Let $L \subset \mathbb{N}^d$ of levels that are selected, then obtain

$$u_L = \sum_{I \in L} w_I$$

and

$$u - u_L = \sum_{I \notin L} w_I$$

For each component w_I have derived bounds of the form

$$\|w_I\| \leq s(I)\mu(u)$$

where $\mu(u)$ was typically $\|\partial^{2d}u\|_2$. Then obtain

$$\begin{aligned} \|u - u_L\| &\leq \sum_{I \notin L} \|w_I\|_2 \leq \left(\sum_{I \notin L} s(I) \right) \mu(u) \\ &= \left[\left(\sum_{I \in \mathbb{N}^d} s(I) \right) - \left(\sum_{I \in L} s(I) \right) \right] \mu(u) \end{aligned}$$

Thus, if we select I with $s(I)$ big, then the error $\|u - u_L\|$ is reduced \rightarrow use $s(I)$ as the benefit of subspace W_I

SG: Quality of full-grid space

With this new tool in hand, let us analyze the cost/benefit of a “full grid”, i.e., a grid corresponding to the selection

$$L_n = \{I : \|I\|_\infty \leq n\}$$

In the L_2 norm, we have bounds of the order

$$s(I) = 2^{-2\|I\|_1}$$

Then, calculations show that

$$\sum_{I \in L_n} s(I) \geq \left(\frac{1}{3}\right)^d (1 - d2^{-2n})$$

and for $n \rightarrow \infty$

$$\sum_{I \in \mathbb{N}^d} s(I) = \left(\frac{1}{3}\right)^d$$

Thus

$$\sum_{I \in \mathbb{N}^d} s(I) - \sum_{I \in L_n} s(I) \leq \left(\frac{1}{3}\right)^d - \left(\frac{1}{3}\right)^d (1 - d2^{-2n}) \leq \frac{d}{3^d} 2^{-2n}$$

SG: Approximation quality of full-grid space

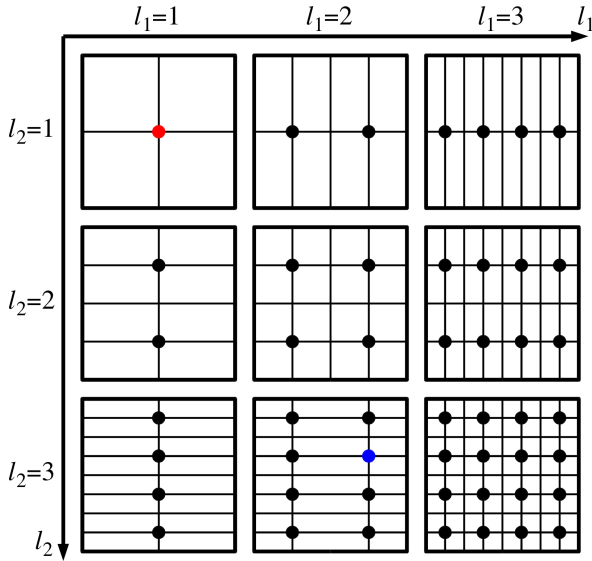
Benefit: $e^{-2|l|}$

Obtain

$$\|u - u_{L_n}\|_2 \leq C \sum_{I \notin L_n} s(I) \leq \frac{Cd}{3^d} 2^{-2n} \in \mathcal{O}(h_n^2)$$

This is what we expect from a piecewise linear approximation

What additional insights have we obtained?



SG: Approximation quality of full-grid space

Obtain

$$\|u - u_{L_n}\|_2 \leq C \sum_{I \notin L_n} s(I) \leq \frac{Cd}{3^d} 2^{-2n} \in \mathcal{O}(h_n^2)$$

This is what we expect from a piecewise linear approximation

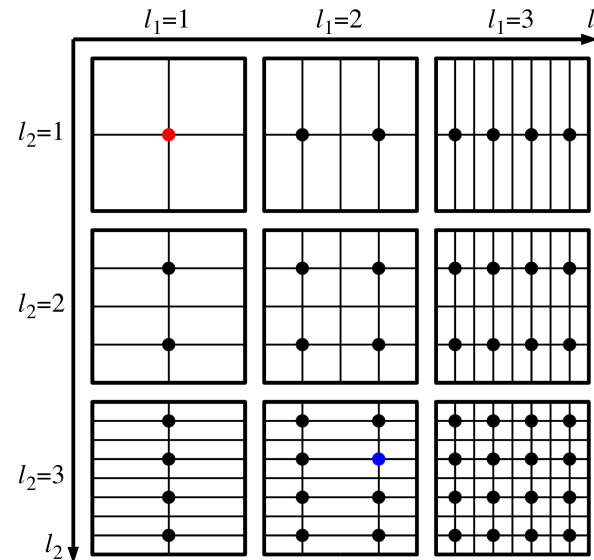
What additional insights have we obtained?

- ▶ The sum of local benefits

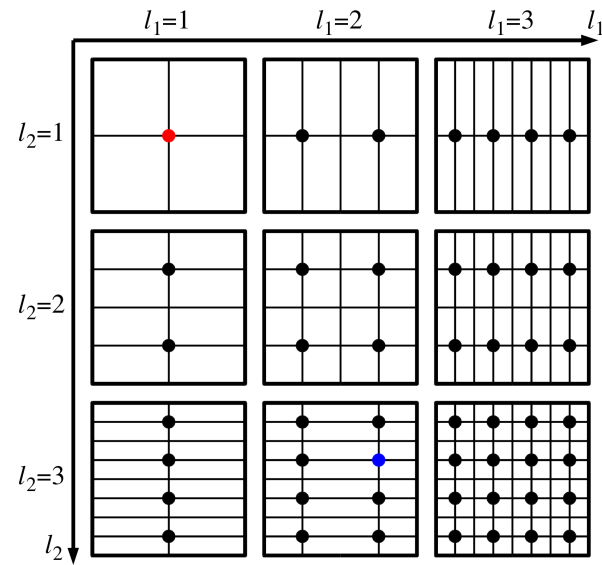
$$\sum_{I \in L_n} 2^{-2|I|_1},$$

means that subspace on diagonals (i.e., with constant $|I|_1$) have the same benefit.

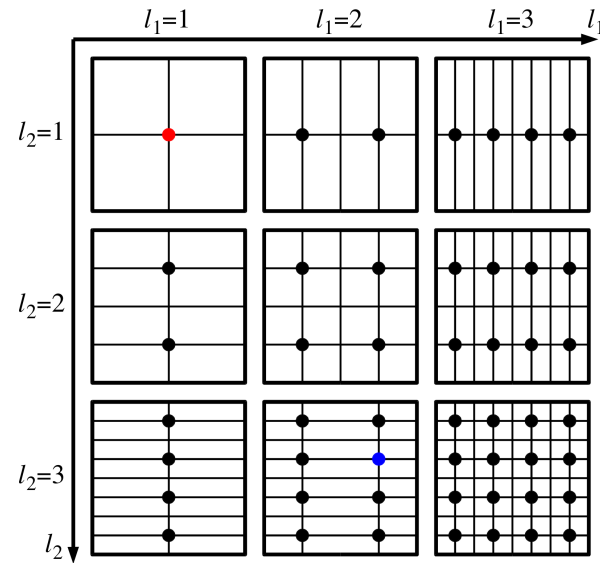
- ▶ Also, as we move further to the bottom right, the benefit gets less and less



If we now look at the costs $c(I) = 2^{|I|_1 - d}$ (which is the number of grid points), then we see that the costs are constant on diagonals as well

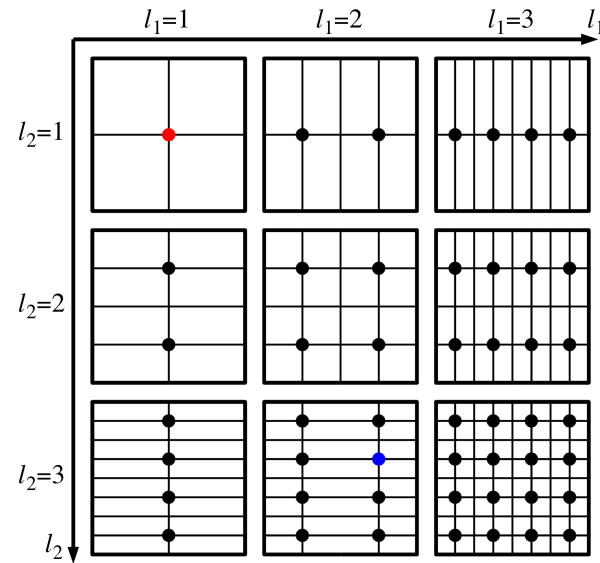


If we now look at the costs $c(I) = 2^{|I|_1 - d}$ (which is the number of grid points), then we see that the costs are constant on diagonals as well



Thus, the cost/benefit ratio $c(I)/s(I)$ is *constant* on diagonals. In particular, for lower-triangular diagonals, we add subspaces with worse and worse cost/benefit ratios
 \Rightarrow **what should we do with them?**

If we now look at the costs $c(I) = 2^{|I|_1 - d}$ (which is the number of grid points), then we see that the costs are constant on diagonals as well

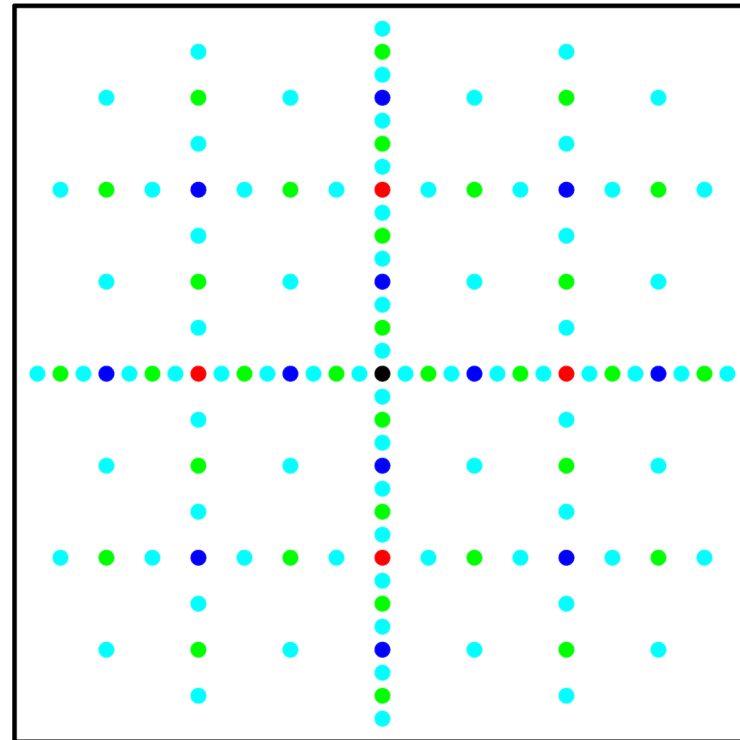
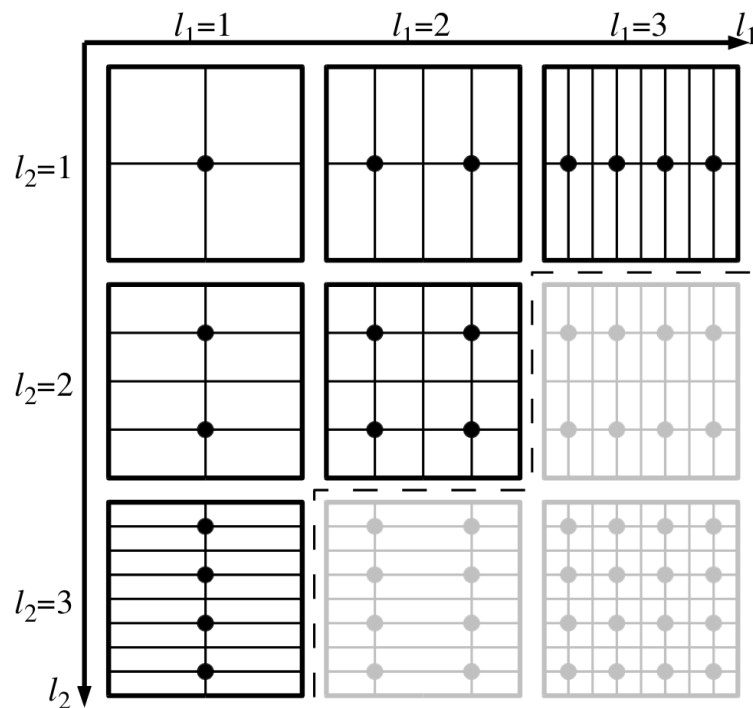


Thus, the cost/benefit ratio $c(I)/s(I)$ is *constant* on diagonals. In particular, for lower-triangular diagonals, we add subspaces with worse and worse cost/benefit ratios \Rightarrow **what should we do with them?** let's truncate them

Consider the diagonal cut $L_n^{(1)} = \{I : |I|_1 \leq n + d - 1\}$ and the sparse grid space

$$\mathbf{V}_n^{(1)} = \bigoplus_{|I|_1 \leq n+d-1} W_I$$

Here is an example of a sparse grid



SG: Properties of sparse grids

The number of grid points of a sparse grid grows as $\mathcal{O}(2^n n^{d-1})$ in contrast to $\mathcal{O}(2^{nd})$ of a full grid

If u has (L_2 -)bounded mixed derivatives up to order $2d$, then

$$\|u - u_n^{(1)}\|_2 \in \mathcal{O}(2^{-2n} n^{d-1}),$$

whereas a full-grid space achieves

$$\|u - u_n\|_2 \in \mathcal{O}(2^{-2n})$$

Sparse-grid spaces achieve slightly worse error than full-grid space but drastically reduced points in higher dimensions d

Comparing the number of grid points corresponding to full-grid and sparse-grid spaces:

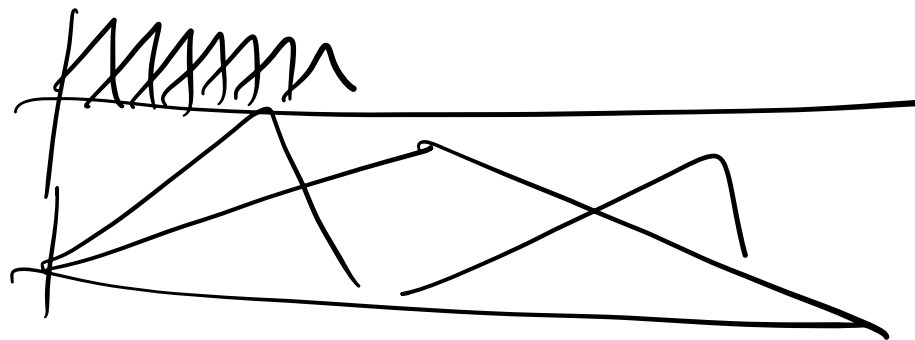
Dimension $d = 2$:

n	1	2	3	4	5	...	10
$\dim V_n = (2^n - 1)^2$	1	9	49	225	961	...	1 046 529
$\dim V_n^1 = 2^n(n - 1) + 1$	1	5	17	49	129	...	9 217

Dimension $d = 3$:

n	1	2	3	4	...	10
$\dim V_n = (2^n - 1)^3$	1	27	343	3 375	...	1 070 590 167
$\dim V_n^1 = 2^n \left(\frac{n^2}{2} - \frac{n}{2} + 1 \right) - 1$	1	7	31	111	...	47 103

Why does this work?



Why does this work?

Exploit the additional smoothness given by the assumption on the mixed derivatives of function u

The hierarchical basis is a key ingredient:

Why does this work?

Exploit the additional smoothness given by the assumption on the mixed derivatives of function u

The hierarchical basis is a key ingredient:

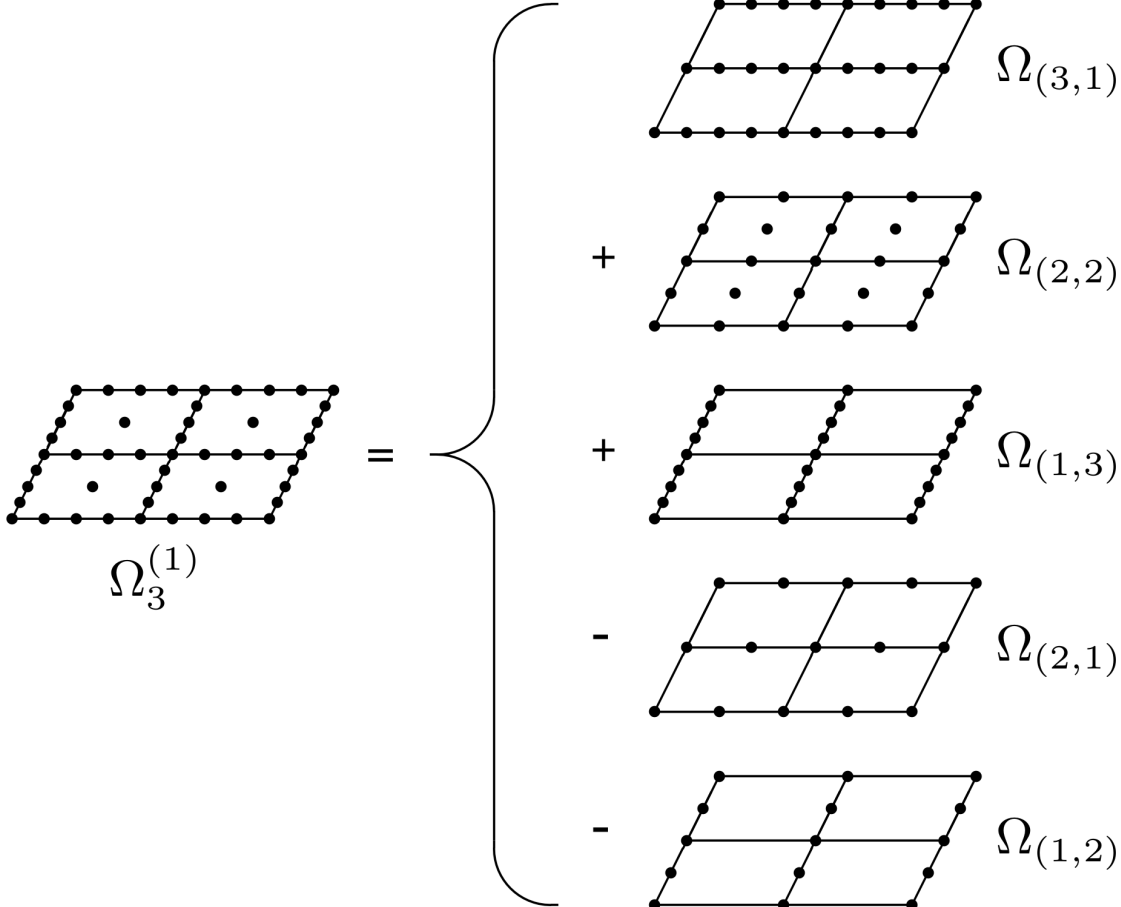
- ▶ Exploits smoothness by having “semi-global” support (i.e., function is smoother, so we can reach far over the domain and know it won’t change too much),
- ▶ Introduced a hierarchy/multilevel and the coefficients in this hierarchy/multilevel basis decay fast (\rightarrow multigrid, multilevel Monte Carlo)

(Logarithmic dependence can be avoided if measure error in energy norm)

Details: Bungartz, Griebel, Acta Numerica, 2004

SG: Combination technique

Formally, sparse grids are superpositions of coarser full grids



Numerical Methods I

MATH-GA 2010.001/CSCI-GA 2420.001

Benjamin Peherstorfer
Courant Institute, NYU

Based on slides by G. Stadler and A. Donev

- Equidistantly spaced nodes

$$h_i = h = \frac{b - a}{n}, \quad t_i = a + ih, \quad i = 0, \dots, n$$

then quadrature formulas are called the *Newton-Cotes formulas* with weights

$$\lambda_{in} = \frac{1}{b - a} \int_a^b \prod_{i \neq j} \frac{t - t_i}{t_i - t_j} dt = \frac{1}{n} \int_0^n \prod_{i \neq j} \frac{s - j}{i - j} ds$$

These weights are independent of the interval boundaries a and b and can be pre-computed once and for all:

Table 9.1. Newton-Cotes weights λ_{in} for $n = 1, \dots, 4$.

n	$\lambda_{0n}, \dots, \lambda_{nn}$	Error	Name
1	$\frac{1}{2} \quad \frac{1}{2}$	$\frac{h^3}{12} f''(\tau)$	Trapezoidal rule
2	$\frac{1}{6} \quad \frac{4}{6} \quad \frac{1}{6}$	$\frac{h^5}{90} f^{(4)}(\tau)$	Simpson's rule, Kepler's barrel rule
3	$\frac{1}{8} \quad \frac{3}{8} \quad \frac{3}{8} \quad \frac{1}{8}$	$\frac{3h^5}{80} f^{(4)}(\tau)$	Newton's 3/8-rule
4	$\frac{7}{90} \quad \frac{32}{90} \quad \frac{12}{90} \quad \frac{32}{90} \quad \frac{7}{90}$	$\frac{8h^7}{945} f^{(6)}(\tau)$	Milne's rule

Generic rule

$$\hat{I}(f) = (b-a) \sum_{i=0}^n \lambda_i f(t_i)$$

Trapezoidal

$$\hat{I}_T(f) = \underline{(b-a)} \left(\underline{\frac{1}{2}} f(a) + \underline{\frac{1}{2}} f(b) \right)$$

Simpson

$$\hat{I}_S(f) = \underline{(b-a)} \left(\underline{\frac{1}{6}} f(a) + \underline{\frac{4}{6}} f\left(\frac{a+b}{2}\right) + \underline{\frac{1}{6}} f(b) \right)$$

Lagrange polynomials: $x_0 = a$, $x_1 = b$

$$l_i(x) = \prod_{\substack{j=0 \\ j \neq i}}^n \frac{(x-x_j)}{(x_i-x_j)}$$

$$l_0(x) = \frac{(x-b)}{(a-b)} \quad \int_a^b l_0(x) dx = \frac{b-a}{2}$$

$$l_1(x) = \frac{(x-a)}{(b-a)} \quad \int_a^b l_1(x) dx = \frac{b-a}{2}$$

For $n+1$ pairwise distinct nodes t_0, \dots, t_n
there exists one and only one quadrature rule

$$\hat{I}(f) = (b-a) \sum_{i=0}^n \lambda_i f(t_i)$$

which is exact for $P \in \mathbb{P}_n$.

Have \hat{I} that is exact for $P \in \mathbb{P}_n$

\Rightarrow insert Lagrange polynomials for t_0, \dots, t_n

$$L_{i,n} \in \mathbb{P}_n$$

$$\begin{aligned} I(L_{i,n}) &= \hat{I}(L_{i,n}) = (b-a) \sum_{j=0}^n \lambda_j L_{i,n}(t_j) \\ &= (b-a) \lambda_i \end{aligned}$$

$$\lambda_i = \frac{1}{(b-a)} I(L_{i,n})$$

Newton-Cotes formulas continued

Given **fixed** nodes t_0, \dots, t_n , use polynomial approximation

$$\hat{f} = P_f(t|t_0, \dots, t_n) = \sum_{i=0}^n f(t_i)L_{in}(t)$$

with Lagrange polynomials L_{0n}, \dots, L_{nn}

Thus:

$$\hat{I}(f) = (b - a) \sum_{i=0}^n \lambda_{in} f(t_i),$$

where

$$\lambda_{in} = \frac{1}{b - a} \int_a^b L_{in}(t) dt$$

Gauss-Christoffel quadrature (cont'd)

Can we say something about the nodes $\tau_{0n}, \dots, \tau_{nn}$?

Gauss-Christoffel quadrature (cont'd)

Can we say something about the nodes $\tau_{0n}, \dots, \tau_{nn}$?

Theorem: For $n \in \mathbb{N}$, consider nodes $\tau_{0n}, \dots, \tau_{nn}$. For an n , define \hat{I} as

$$\hat{I}_n(f) = \sum_{i=0}^n \lambda_{in} f(\tau_{in})$$

and let it be exact for polynomials $p \in \mathbb{P}_{2n+1}$

$$\hat{I}_n(p) = \int_a^b \omega(t)p(t)dt.$$

Then, the polynomials $\{P_n\}$ given by

$$P_{n+1}(t) = (t - \tau_{0n}) \cdot \dots \cdot (t - \tau_{nn}) \in \mathbb{P}_{n+1}$$

are orthogonal with respect to the scalar product induced by $\omega(t)$

Gauss-Christoffel quadrature (cont'd)

Therefore, the nodes τ_{in} have to be roots of pairwise orthogonal polynomials $\{P_j\}$ of degree $\deg(P_j) = j$

For a given ω , we already know that the set of orthogonal polynomials $\{P_j\}$ is unique (if leading coefficient is 1)

We also know that the roots of these polynomials are real and have to lie in (a, b)
For each ω and $n \in \mathbb{N}$, this gives us a unique set of nodes, namely the roots of the corresponding orthogonal polynomial P_{n+1}

We thus have fixed $n + 1$ degrees of freedom so far...

What do we know about the weights for fixed t_0, \dots, t_n when we want to integrate exactly polynomials up to degree n ?

What do we know about the weights for fixed t_0, \dots, t_n when we want to integrate exactly polynomials up to degree n ?

Because we want to be able to exactly integrate polynomials up to degree $2n + 1$, we already know that to even achieve exactness up to degree n , the weights are fixed for given nodes t_0, \dots, t_n :

$$\lambda_{in} = \frac{1}{b-a} \int_a^b L_{in}(t) dt, \quad \text{Lagrange poly } L_{in}(\tau_{jn}) = \delta_{ij}$$

(“For $n + 1$ pairwise distinct nodes, there exists only one quadrature formula that exactly integrates polynomials up to degree n .”)

Theorem: Let $\tau_{0n}, \dots, \tau_{nn}$ be the roots of the $(n + 1)$ st orthogonal polynomial for the weight ω . Then any quadrature formula \hat{I} is exact for polynomials up to order n if and only if it is exact up to order $2n + 1$.

Proof

Let τ_0, \dots, τ_n be the roots of the $(n+1)$ st
ortho. polynomial with weight w . Then any
rule

$$\hat{I}_n(f) = \sum_{i=0}^n \lambda_i f(\tau_i)$$

satisfies

$$\hat{I}_n \text{ exact on } P_n \iff \hat{I}_n \text{ exact on } P_{2n+1}$$

$\leftarrow \Rightarrow$:

Suppose $P \in P_{2n+1}$, then there exists

$$Q, R \in P_n$$

$$P = Q P_{n+1} + R$$

P_{n+1} orthogonal to P_n , so that

$$\int_0^b w P = \underbrace{\int_0^b w Q P_{n+1}}_{=0} + \int_0^b w R = \int_0^b w R = \hat{I}(R) =$$

$$-\hat{I}(R) = \sum_{i=0}^n \lambda_i R(\tau_i)$$

$$= \sum_{i=0}^n \lambda_i \underbrace{(Q(\tau_i) P_{n+1}(\tau_i))}_{=0} + R(\tau_i)$$

$$= \hat{I}(P)$$

Lagrange vs Newton

$$l_i(x) = \prod_{\substack{j=0 \\ i \neq j}}^n \frac{(x - x_j)}{(x_i - x_j)}$$

$$w_i(x) = \prod_{j=0}^{i-1} (x - x_j)$$

Recap: Polynomial interpolation in Newton basis

The **Newton basis** $\omega_0, \dots, \omega_n$ is given by

$$\omega_i(x) := \prod_{j=0}^{i-1} (x - x_j) \in \mathbb{P}_i.$$

Would like to find coefficients c_0, c_1, \dots, c_n of interpolating polynomial in Newton basis

$$P_f(x|x_0, \dots, x_n) = c_0\omega_0(x) + c_1\omega_1(x) + \dots + c_n\omega_n(x)$$

Today •

Last time

- ▶ Function approximation in higher dimensions
- ▶ Quadrature in higher dimensions

Today

- ▶ Monte Carlo

Announcements

- ▶ Homework 7 is posted and due Mon, Dec 9 before class (1 week)
- ▶ Next week, Mon Dec 9, recap of important topics—highly recommended!
- ▶ Email me by Wed, Dec 4 if you need special accommodations for the final exam

There are many opportunities for asking questions

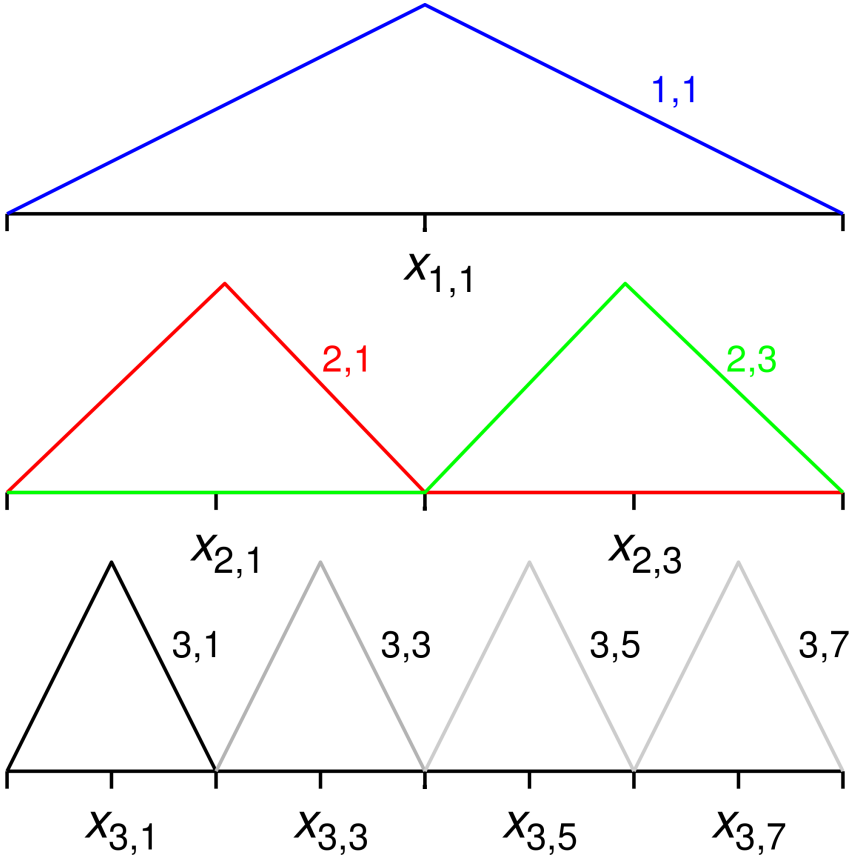
Mon, Dec 2	office hour
Wed, Dec 4	lecture (virtual) [not part of final exam]
Fri, Dec 6	office hour (grader)
Mon, Dec 9	recap and Q&A
Mon, Dec 9	office hour
Wed, Dec 11	extra office hour, 6.10pm ET (WWH 421) (all HWs graded; no HW re-grading after Thu, Dec 12)
Mon, Dec 16	final exam

- ▶ Double check that all homeworks are entered correctly in Brightspace
- ▶ Will go through example final exam problems next week
- ▶ Expect to write *some* code in the final exam
- ▶ Be in the room 10min early on the day of the final exam (**most likely room changed to what is in Albert now**)

High-dimensional interpolation and quadrature

SG: Hierarchical basis cont'd

The bases for spaces W_1 , W_2 and W_3



SG: Hierarchical basis cont'd

Obtain

$$V_n = \bigoplus_{l=1}^n W_l,$$

so that there is a unique representation for each $u \in V_n$ as

$$u = \sum_{l=1}^n w_l = \sum_{l=1}^n \sum_{i \in I_l} v_{l,i} \phi_{l,i}$$

Coefficients $v_{l,i}$ in case of interpolation \rightsquigarrow visualize on board

The coefficients $v_{l,i}$ are hierarchical differences

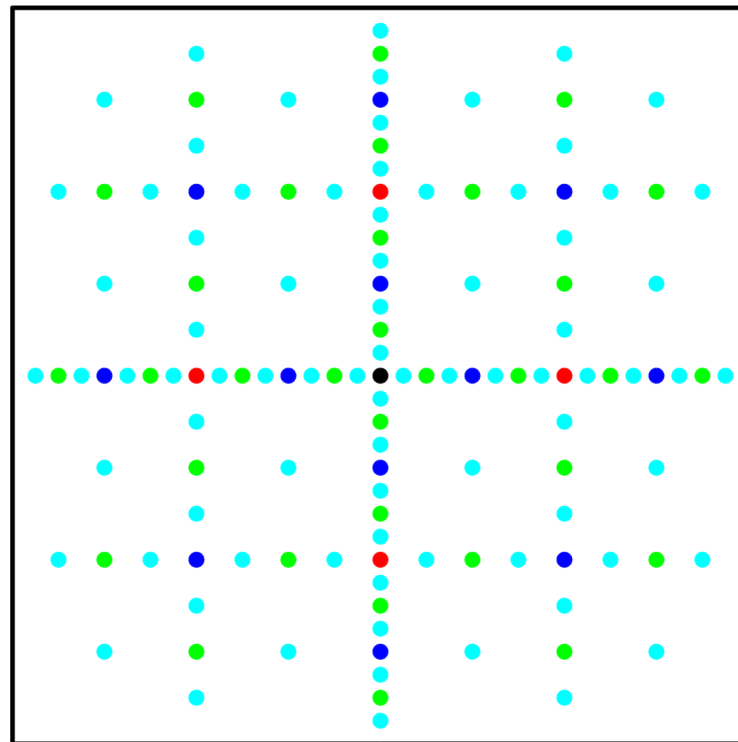
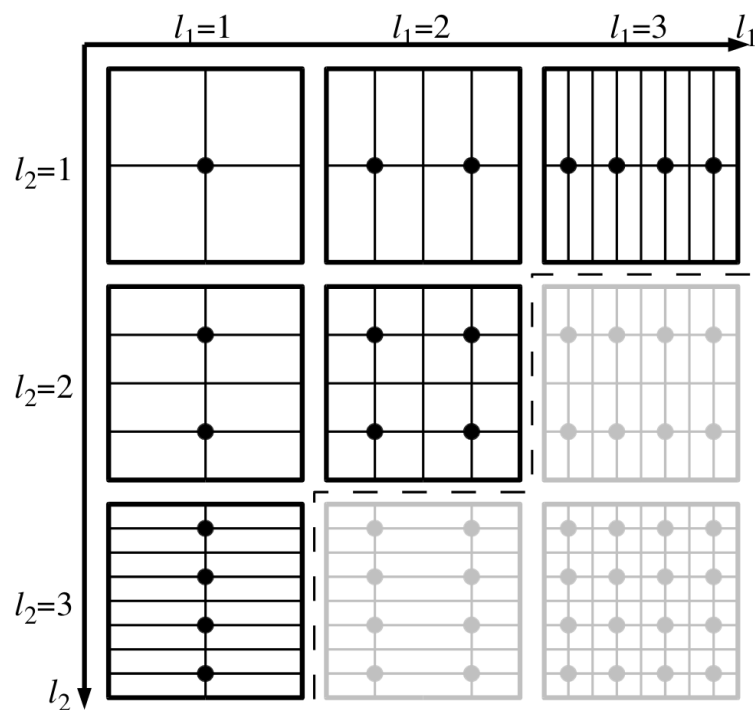
$$v_{l,i} = u^*(x_{l,i}) - \frac{u(x_{l,i-1}) + u(x_{l,i+1})}{2}$$

where u^* is the function to be interpolated

Consider the diagonal cut $L_n^{(1)} = \{I : |I|_1 \leq n + d - 1\}$ and the sparse grid space

$$\mathbf{V}_n^{(1)} = \bigoplus_{|I|_1 \leq n+d-1} W_I$$

Here is an example of a sparse grid



SG: Properties of sparse grids

The number of grid points of a sparse grid grows as $\mathcal{O}(2^n n^{d-1})$ in contrast to $\mathcal{O}(2^{nd})$ of a full grid

If u has (L_2 -)bounded mixed derivatives up to order $2d$, then

$$\|u - u_n^{(1)}\|_2 \in \mathcal{O}(2^{-2n} n^{d-1}),$$

whereas a full-grid space achieves

$$\|u - u_n\|_2 \in \mathcal{O}(2^{-2n})$$

Sparse-grid spaces achieve slightly worse error than full-grid space but drastically reduced points in higher dimensions d

Comparing the number of grid points corresponding to full-grid and sparse-grid spaces:

Dimension $d = 2$:

n	1	2	3	4	5	...	10
$\dim V_n = (2^n - 1)^2$	1	9	49	225	961	...	1 046 529
$\dim V_n^1 = 2^n(n - 1) + 1$	1	5	17	49	129	...	9 217

Dimension $d = 3$:

n	1	2	3	4	...	10
$\dim V_n = (2^n - 1)^3$	1	27	343	3 375	...	1 070 590 167
$\dim V_n^1 = 2^n \left(\frac{n^2}{2} - \frac{n}{2} + 1 \right) - 1$	1	7	31	111	...	47 103

Why does this work?

Why does this work?

Exploit the additional smoothness given by the assumption on the mixed derivatives of function u

The hierarchical basis is a key ingredient:

Why does this work?

Exploit the additional smoothness given by the assumption on the mixed derivatives of function u

The hierarchical basis is a key ingredient:

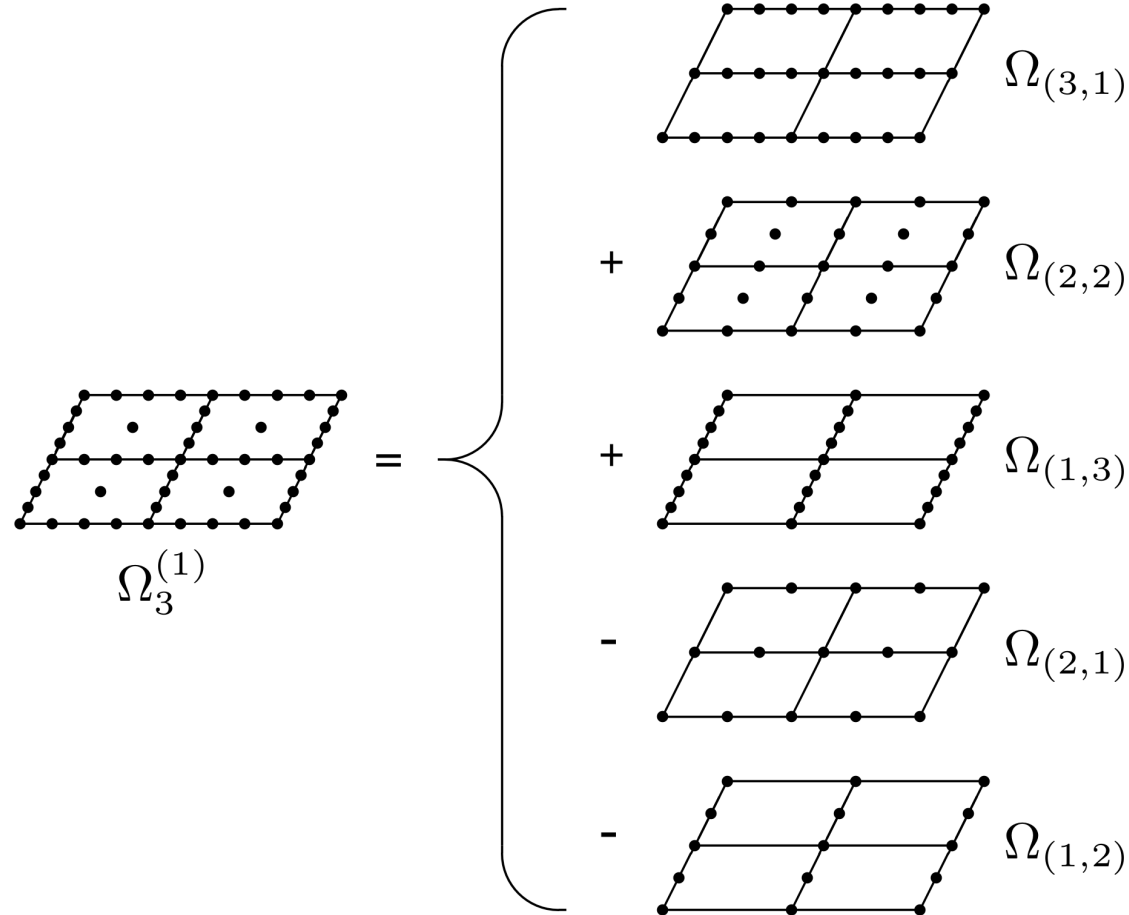
- ▶ Exploits smoothness by having “semi-global” support (i.e., function is smoother, so we can reach far over the domain and know it won’t change too much),
- ▶ Introduced a hierarchy/multilevel and the coefficients in this hierarchy/multilevel basis decay fast (\rightarrow multigrid, multilevel Monte Carlo)

(Logarithmic dependence can be avoided if measure error in energy norm)

Details: Bungartz, Griebel, Acta Numerica, 2004

SG: Combination technique

Formally, sparse grids are superpositions of coarser full grids



SG: Combination technique (cont'd)

This works for grids and also for functions (in certain situations)

- ▶ interpolation
- ▶ quadrature
- ▶ solutions of partial differential equations → limited

⇒ we are interested in quadrature

For quadrature, Smolyak has developed a related approach already in 1963

SG: Smolyak quadrature - 1D

Set $D = [-1, 1]$ and consider a one-dimensional function $f : D \rightarrow \mathbb{R}$ and we are interested in

$$If = \int_D f(x) dx$$

One-dimensional quadrature rule

$$Q_I^1 f = \sum_{i=1}^{n_I} w_i f(x_i)$$

with weights w_1, \dots, w_{n_I} and points x_1, \dots, x_{n_I} and

$$X_I = \{x_i : 1 \leq i \leq n_I\}$$

Quadrature rules are nested if $X_I \subset X_{I+1}$

SG: Smolyak quadrature

Define the difference formula

$$\Delta_k f = (Q_k^1 - Q_{k-1}^1) f,$$

with $Q_0^1 f = 0$. What does the difference formula remind you of?

SG: Smolyak quadrature

Define the difference formula

$$\Delta_k f = (Q_k^1 - Q_{k-1}^1) f,$$

with $Q_0^1 f = 0$. **What does the difference formula remind you of?** hierarchical coefficients

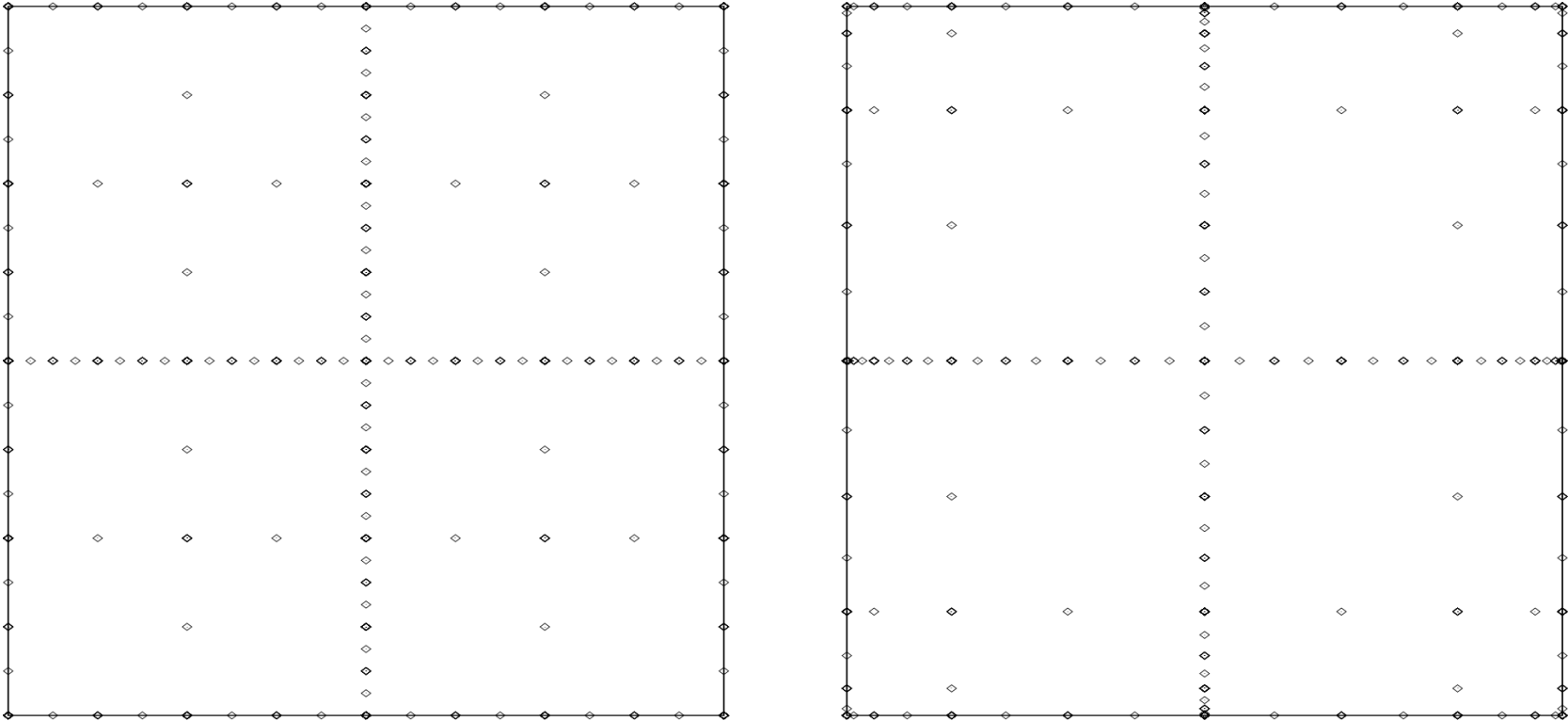
Smolyak's quadrature rule is

$$Q_l^d f = \sum_{|\mathbf{k}|_1 \leq n+d-1} (\Delta_{k_1}^1 \otimes \cdots \otimes \Delta_{k_d}^1) f,$$

where the tensor product of quadrature rules is

$$(\Delta_{k_1}^1 \otimes \cdots \otimes \Delta_{k_d}^1) f = \sum_{i_1=1}^{n_{k_1}} \cdots \sum_{i_d=1}^{n_{k_d}} w_{k_1, i_1} \cdots w_{k_d, i_d} f(x_{k_1, i_1}, \dots, x_{k_d, i_d})$$

SG: Smolyak grid w.r.t. Clenshaw-Curtis rule



Left: sparse grid w.r.t. trapezoidal rule (piecewise constant), right: sparse grid obtained with Clenshaw-Curtis rule

SG: Alternative representations of Smolyak

A non-hierarchical representation is

$$Q_l^d f = \sum_{n \leq |l|_1 \leq n+d-1} (-1)^{n+d-1-|l|_1} \binom{d-1}{|l|_1-n} (Q_{l_1}^1 \otimes \cdots \otimes Q_{l_d}^1) f$$

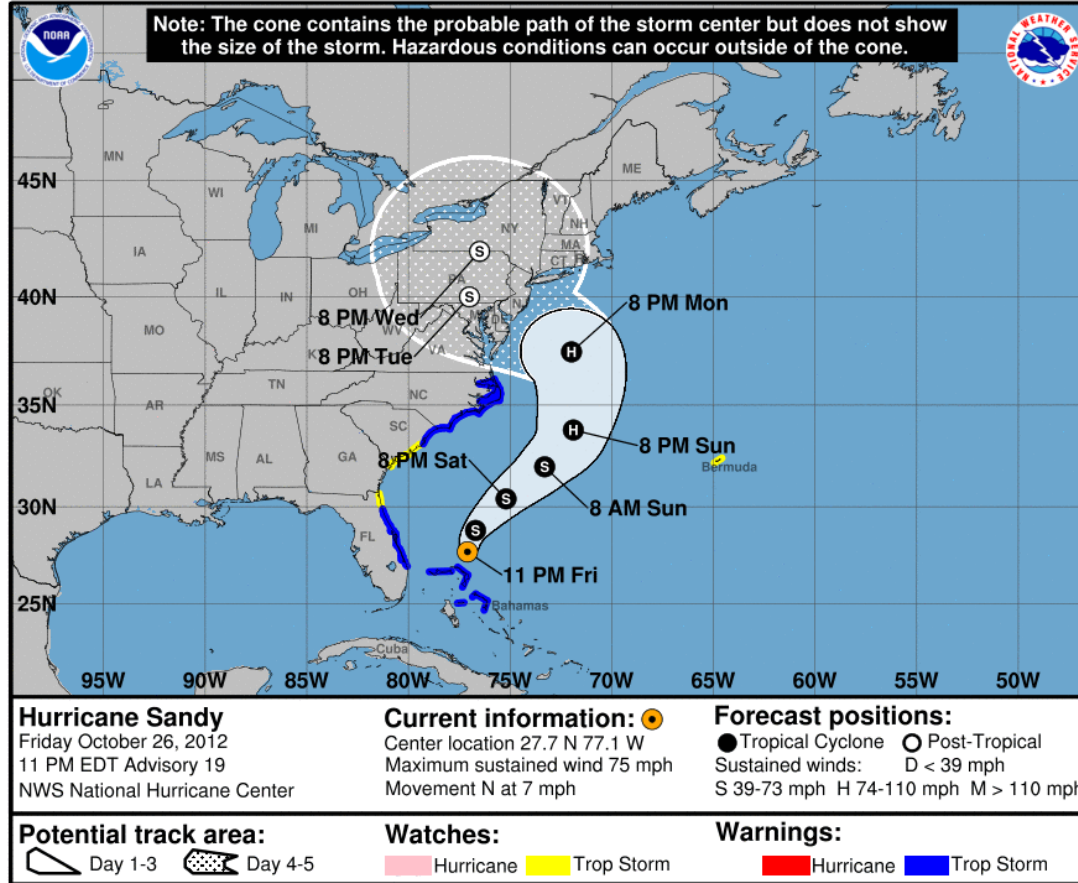
- ▶ Non-hierarchical: Can work on regular grids as with combination technique
- ▶ Simple to implement
- ▶ (Equivalence to hierarchical representation not obvious.)

Conclusions

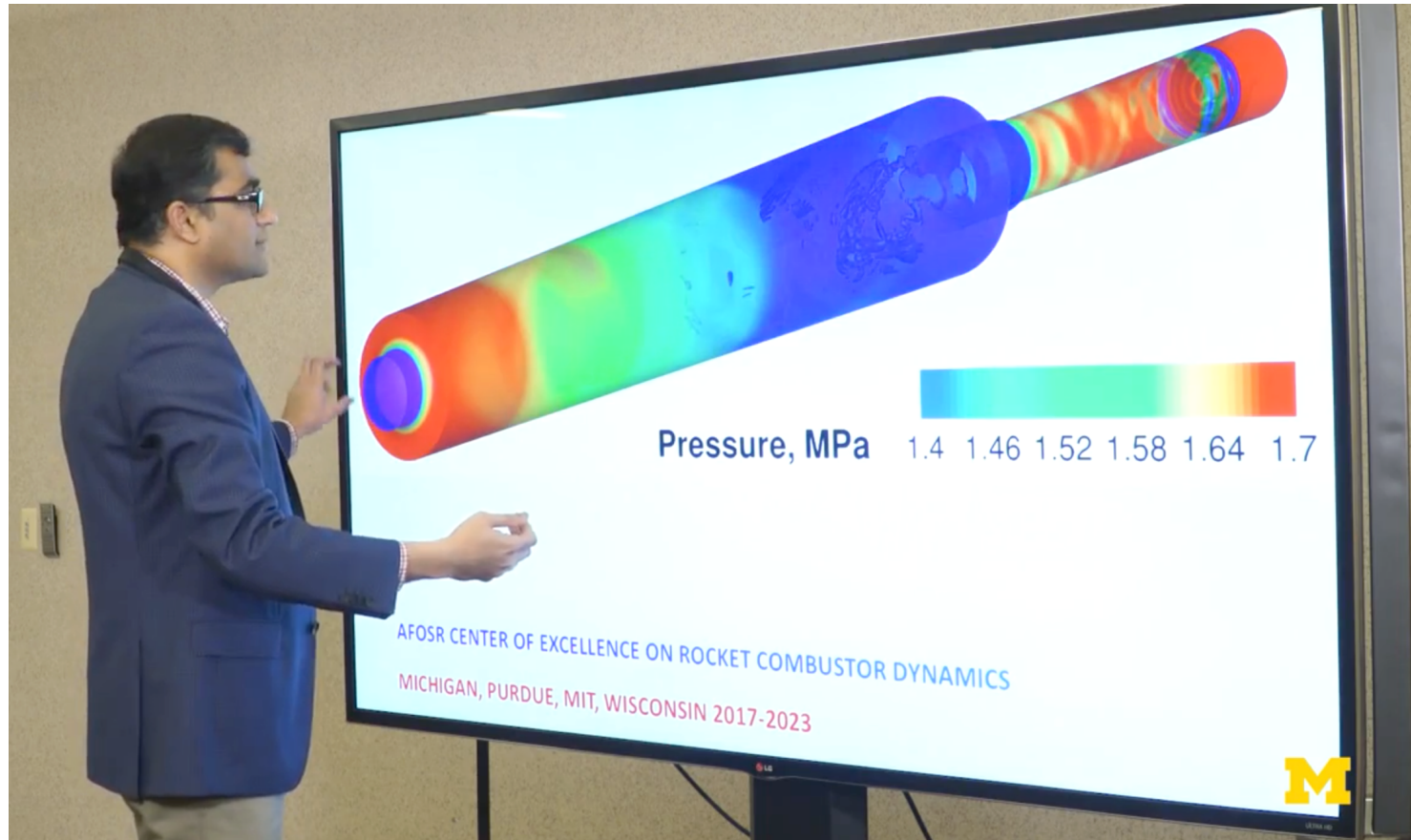
- ▶ Computations in higher dimensions are typically affected by the curse of dimensionality, which means that computational costs become exponentially more expensive as the dimension is increased.
- ▶ We have two options in high dimension. The first option is using randomized methods and Monte Carlo which can circumvent the curse
- ▶ The other option is exploiting additional structure in the problem that can help to circumvent the curse to some extent.
- ▶ In case of sparse grids, we can circumvent the curse by assuming additional smoothness of the function to be interpolated.
- ▶ Sparse grids are very useful for quadrature in moderately high dimensions (10 up to few 100s of dimensions). Quadrature based on sparse grids is sometimes called Smolyak quadrature.

Monte Carlo methods and numerical methods

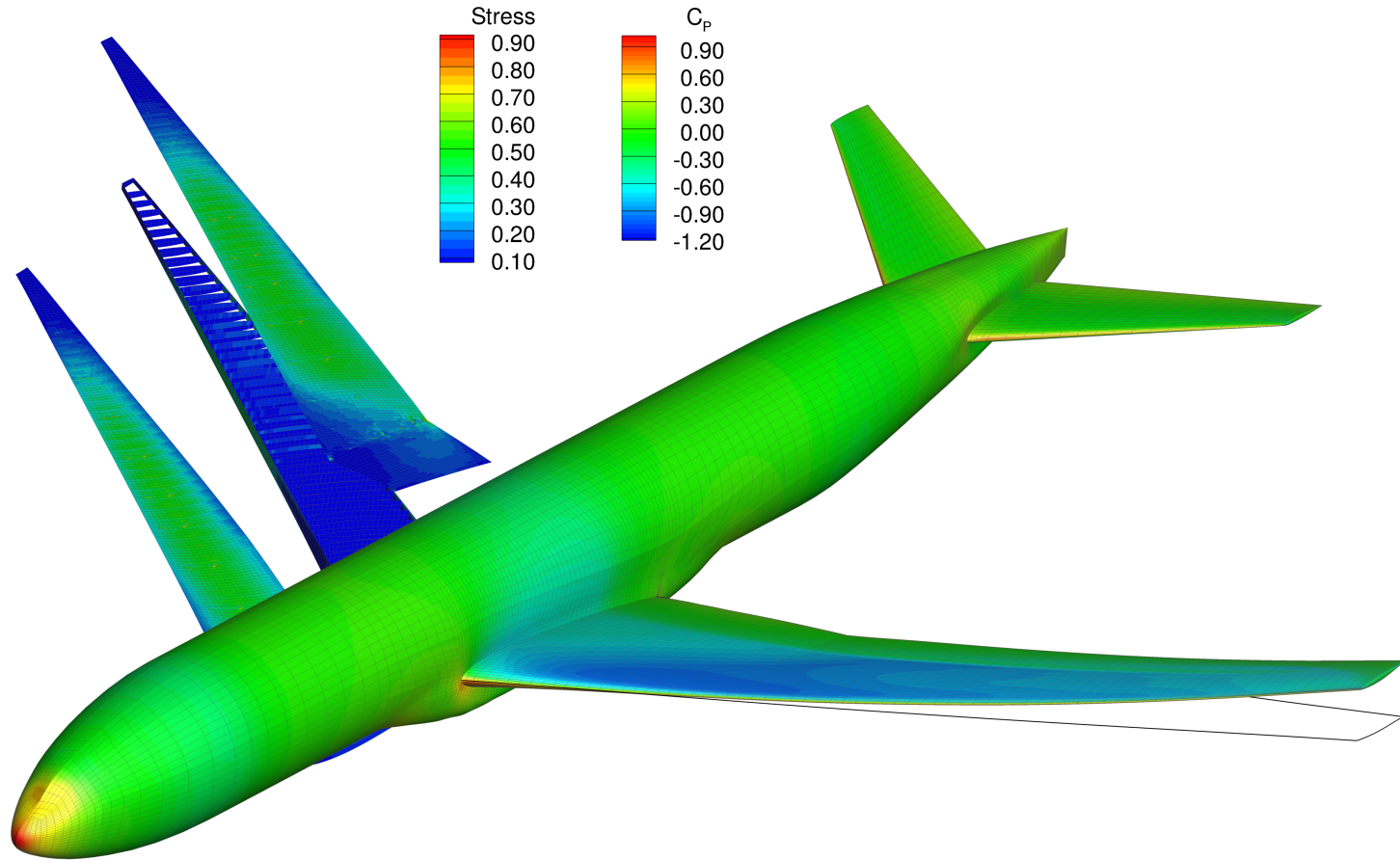
Known unknowns



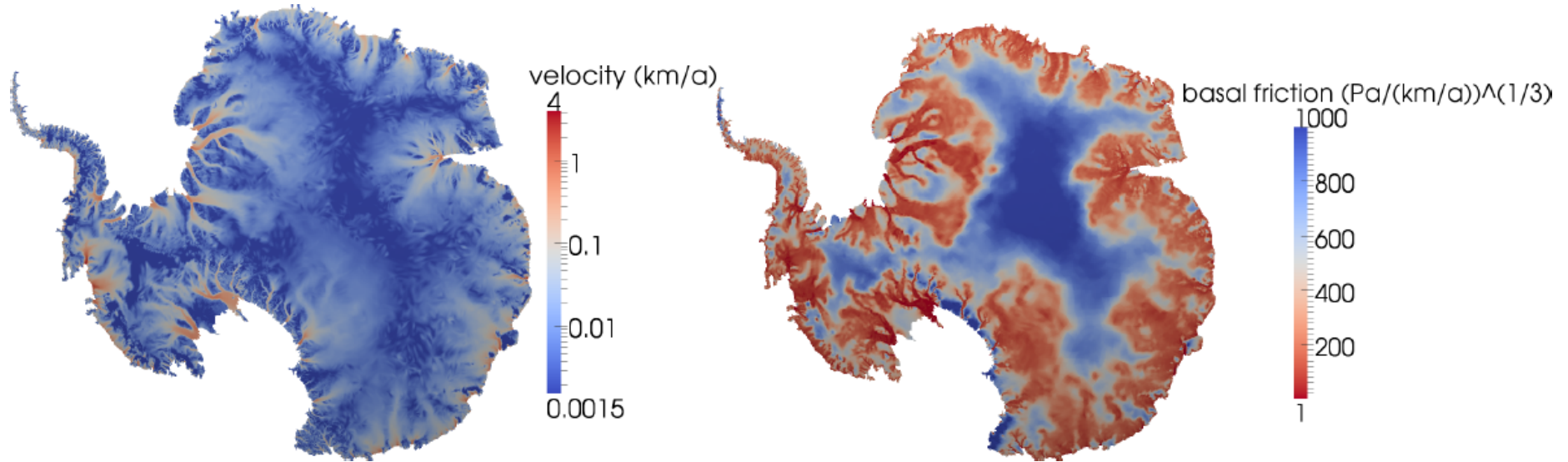
No hope to model physics exhaustively



Rapidly changing dynamics



Uncertainties due to data

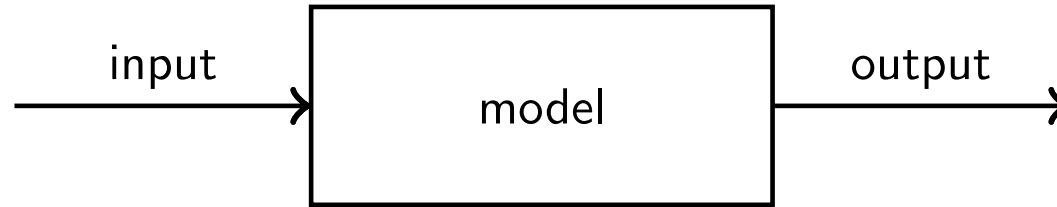


Figures: Petra, Ghattas, Isaac, Martin, Stadler, et al.

Intro: Model

Model of system of interest

- ▶ Model describes response of system to inputs, parameters, configurations
- ▶ Response typically is a quantity of interest
- ▶ Evaluating a model means numerically simulating the model
- ▶ Many models given in form of partial differential equations



Mathematical formulation

$$f : \mathcal{D} \rightarrow \mathcal{Y}$$

- ▶ Input domain \mathcal{D} and output domain \mathcal{Y}
- ▶ Maps $\mathbf{z} \in \mathcal{D}$ input onto $\mathbf{y} \in \mathcal{Y}$ output (quantity of interest)

Intro: Model - Navier-Stokes equations

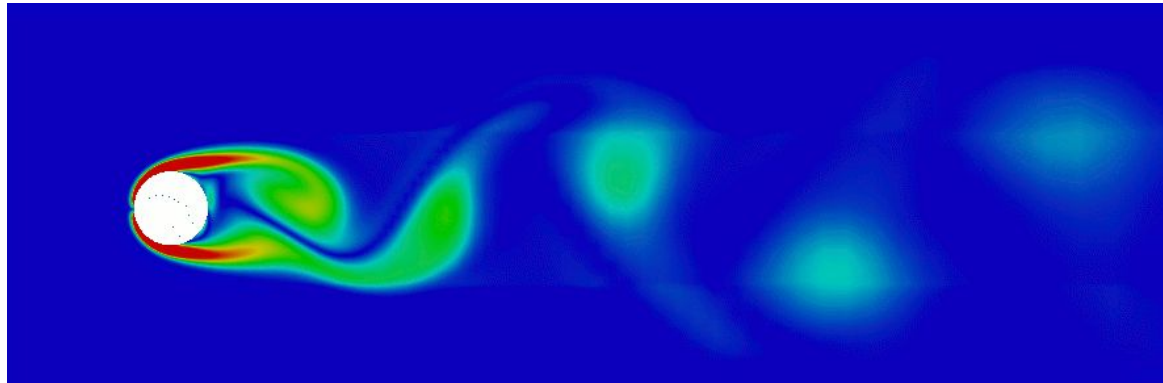
$$\rho \left(\frac{\partial u}{\partial t} + u \cdot \nabla u \right) = -\nabla p + \mu \Delta u + g$$

Examples of inputs

- ▶ Density ρ
- ▶ Dynamic viscosity μ

Examples of outputs (quantities of interest)

- ▶ Velocity at monitoring point
- ▶ Average pressure



Intro: Model - Diffusion-convection-reaction flow

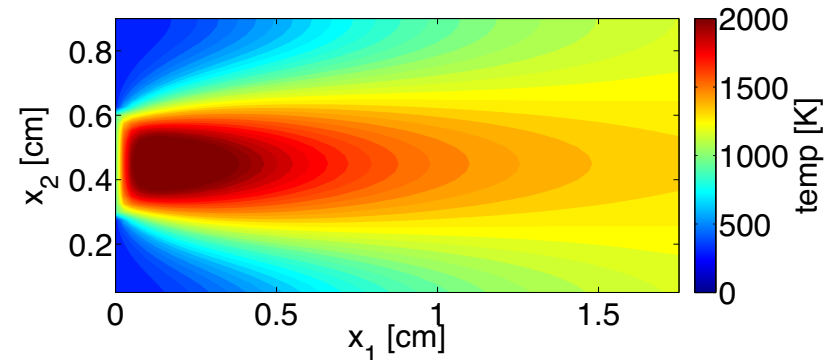
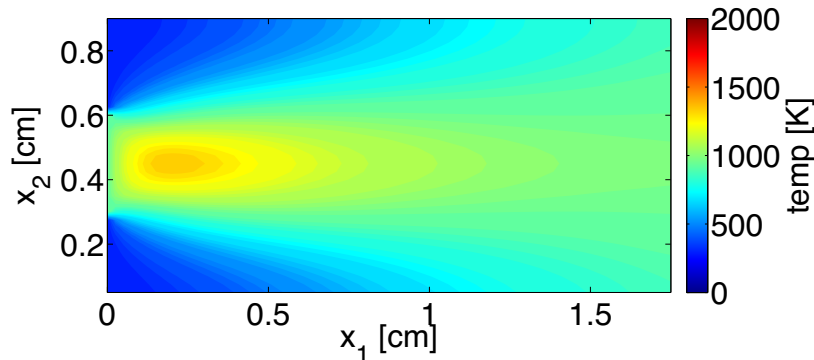
$$\frac{\partial u}{\partial t} = \Delta u - v \nabla u + g(u, \mu)$$

Examples of inputs

- ▶ Activation energy and pre-exponential factor (Arrhenius-type reaction)
- ▶ Temperature at boundary
- ▶ Ratio of fuel and oxidizer

Examples of outputs

- ▶ Average temperature in chamber



Intro: Uncertain inputs

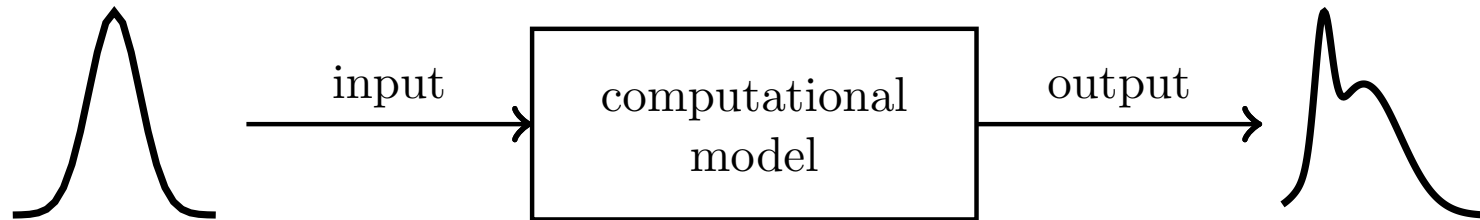
Inputs are uncertain

- ▶ Measurement errors in boundary conditions
- ▶ Manufacturing variations
- ▶ Model parameters determined by engineering judgment, etc.

Mathematically formulate uncertain inputs as random variables

$$Z : \Omega \rightarrow \mathcal{D}$$

Quantify effect of uncertainties in inputs on model outputs



Intro: General sampling-based approach

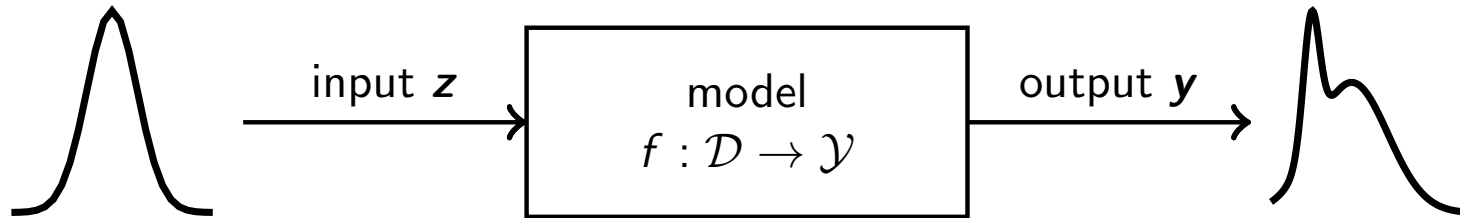
- ▶ Take many realizations of input random variable Z

$$\mathbf{z}_1, \dots, \mathbf{z}_n \in \mathcal{D}$$

- ▶ Evaluate model f at all $\mathbf{z}_1, \dots, \mathbf{z}_n$ realizations

$$\mathbf{y}_1 = f(\mathbf{z}_1), \dots, \mathbf{y}_n = f(\mathbf{z}_n)$$

- ▶ Estimate statistics from outputs $\mathbf{y}_1, \dots, \mathbf{y}_n$



Intro: General sampling-based approach

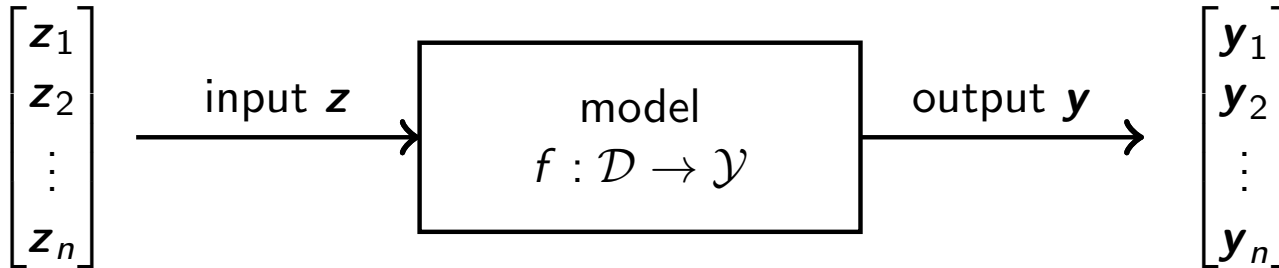
- ▶ Take many realizations of input random variable Z

$$\mathbf{z}_1, \dots, \mathbf{z}_n \in \mathcal{D}$$

- ▶ Evaluate model f at all $\mathbf{z}_1, \dots, \mathbf{z}_n$ realizations

$$\mathbf{y}_1 = f(\mathbf{z}_1), \dots, \mathbf{y}_n = f(\mathbf{z}_n)$$

- ▶ Estimate statistics from outputs $\mathbf{y}_1, \dots, \mathbf{y}_n$



Monte Carlo methods

- ▶ Given X_1, \dots, X_n iid random variables that are distributed as X , the basic Monte Carlo estimator of $\mathbb{E}[X]$ is

$$\bar{X}_n = \frac{1}{n} \sum_{i=1}^n X_i$$

- ▶ Note that each of the n X_1, \dots, X_n is a random variable, and thus \bar{X}_n is a random variable too
- ▶ Thus, it makes sense to consider $\mathbb{E}[\bar{X}_n]$ and $\text{Var}[\bar{X}_n]$
- ▶ Once the samples have been drawn/realized, the estimate \bar{X}_n is a real number (here; clash of terminology)

Monte Carlo estimators

- ▶ Unbiasedness of Monte Carlo estimator $\mathbb{E}[\bar{X}_n] = \mathbb{E}[X]$
- ▶ Variance is $\text{Var}[\bar{X}_n] = \sigma^2/n$
- ▶ With unbiasedness follows

$$\mathbb{E} \left[(\bar{X}_n - \mathbb{E}[X])^2 \right] = \text{Var} [\bar{X}_n] = \frac{\sigma^2}{n}$$

The mean-squared error rate is $1/n$.

VR: Control variates

The constant in the MC rate is the variance $\text{Var}[X]$ of the random variable from which MC samples are being drawn. By designing an equivalent MC approximation with lower variance, we can reduce the MSE

Auxiliary random variable

Let X be a random variable and Y be a another random variable that is correlated to X , i.e.,

$$|\rho(X, Y)| > 0.$$

The (Pearson) correlation coefficient $\rho(X, Y)$ is defined as

$$\rho(X, Y) = \frac{\text{Cov}(X, Y)}{\sqrt{\text{Var}[X] \text{Var}[Y]}}$$

with the covariance

$$\text{Cov}(X, Y) = \mathbb{E}[(X - \mathbb{E}[X])(Y - \mathbb{E}[Y])]$$

Control variates exploit the correlation of X and Y to reduce the MSE

Let X be a random variable and Y be an auxiliary random variable with known $\mathbb{E}[Y]$ and correlation coefficient $\rho = \rho(X, Y)$. Consider the control variate estimator

$$\theta = \bar{X}_m + \alpha(\mathbb{E}[Y] - \bar{Y}_m).$$

For

$$\alpha^* = \rho \sqrt{\frac{\text{Var}[X]}{\text{Var}[Y]}},$$

the MSE of the unbiased estimator θ of $\mathbb{E}[X]$ is

$$e(\theta) = (1 - \rho^2) \frac{\text{Var}[X]}{m}.$$

Proof \rightsquigarrow board

What is missing?

Which important point for us (numerical analysis!) does this analysis completely miss?

What is missing?

Which important point for us (numerical analysis!) does this analysis completely miss?

costs per sample/realization

Need to connect this analysis of Monte Carlo methods to our setting where each sample entails a numerical computation

VR: Cost complexity

Let $w_1 = c(X)$ be the costs of sampling X and let $w_2 = c(Y)$ be the costs of sampling the control variate. \rightsquigarrow board

VR: Cost complexity

Let $w_1 = c(X)$ be the costs of sampling X and let $w_2 = c(Y)$ be the costs of sampling the control variate. \rightsquigarrow board

Plain Monte Carlo achieves $\text{MSE} \leq \epsilon$ with costs

$$c(\theta^{\text{MC}}) \leq \frac{\text{Var}[X]}{\epsilon} w_1 .$$

Control variates estimator

$$c(\theta^{\text{CV}}) \leq \frac{\text{Var}[X]}{\epsilon} (1 - \rho^2)(w_1 + w_2)$$

The control variate estimator has lower costs than the plain Monte Carlo estimator if

$$1 - \rho^2 < \frac{w_1}{w_1 + w_2}$$

The inequality combines ρ and the costs \rightarrow cost vs. accuracy trade-off that we know from numerical analysis

X is RV, Y is RV with known $E[Y]$
correlation $\rho(X, Y)$

$$\Theta = \bar{X}_m + \alpha^* (E[Y] - \bar{Y}_m)$$

unbiasedness

$$\begin{aligned} E[\Theta] &= E[\bar{X}_m] + \alpha^* (E[Y] - E[\bar{Y}_m]) \\ &= E[\bar{X}_m] = E[X] \end{aligned}$$

MSE

$$e(\Theta) = \text{Var}[\Theta]$$

$$\begin{aligned} &= \text{Var}[\bar{X}_m] + \alpha^2 \text{Var}[E[Y] - \bar{Y}_m] + \\ &\quad 2 \text{Cov}[\bar{X}_m, (E[Y] - \bar{Y}_m)] \end{aligned}$$

$$= \dots = \frac{\text{Var}[X]}{m} + \alpha^2 \frac{\text{Var}[Y]}{m} - 2\alpha \frac{1}{m^2} \sum_{i,j=1}^m \text{Cov}[X_i, Y_j]$$

$$= \dots = \frac{\text{Var}[X]}{m} + \alpha^2 \frac{\text{Var}[Y]}{m} - 2\alpha \frac{1}{m} \rho \sqrt{\text{Var}[X]\text{Var}[Y]}$$

find α^* by minimizing $e(\Theta)$ w.r.t. α :

$$\partial_{\alpha} e(\Theta) = 2\alpha \frac{\text{Var}[Y]}{m} - 2 \frac{1}{m} \rho \sqrt{\text{Var}[X]\text{Var}[Y]}$$

$$\stackrel{!}{=} 0 : \alpha^* = \rho \sqrt{\frac{\text{Var}[X]}{\text{Var}[Y]}}$$

plug in: $e(\Theta) = \frac{\text{Var}[X]}{m} (1 - \rho^2)$

Example: Control variates (cont'd)

- ▶ $X \sim \mathcal{U}(0, 1)$ and $f(x) = 10 \sin(x)$
- ▶ Now consider $g(x) = 10(x - x^3/6)$ (Taylor approximation)
- ▶ Compute mean $\mathbb{E}[g(X)] = 10(1/2 - 1/24)$ analytically
- ▶ Set

$$F_i = f(X_i), \quad G_i = g(X_i), \quad i = 1, \dots, m$$

- ▶ Set synthetic costs

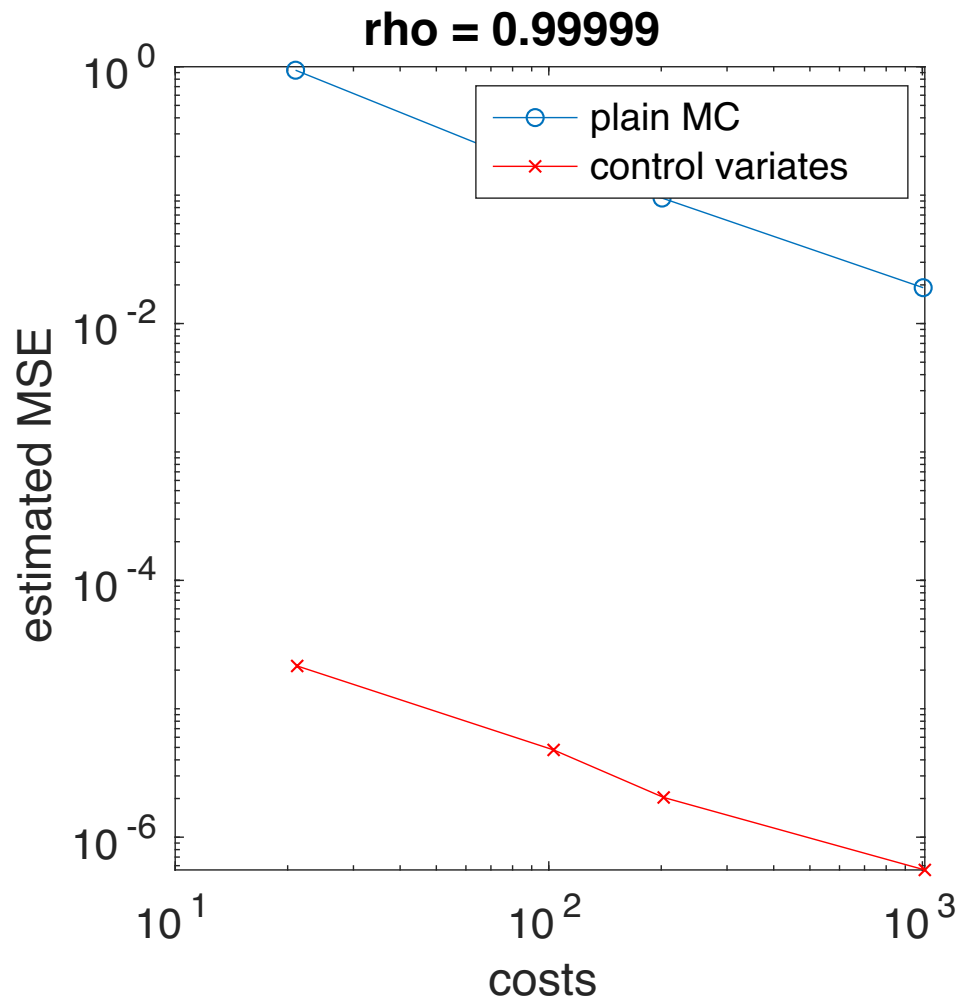
$$w_1 = 3, \quad w_2 = w_1/100$$

- ▶ Measure $\rho(f(X), g(X)) \approx 9.999861330445534e - 01$ and $\text{Var}[f(X)]$ and $\text{Var}[g(X)]$
- ▶ Control variate estimator

$$\theta = \frac{1}{m} \sum_{i=1}^m F_i + \alpha^* \left(\mathbb{E}[g(X)] - \frac{1}{m} \sum_{i=1}^m G_i \right)$$

- ▶ Compare MSE of θ to plain vanilla Monte Carlo estimator

Example: Numerical example



VR: Control variates with unknown expectation

Assumed we know $\mathbb{E}[g(X)]$, i.e., the mean of the control variate. Typically unavailable; for example if g is an approximation of f . Instead, now $\mathbb{E}[g(X)]$ is unknown and needs to be estimated. We obtain the control variate estimator

$$\theta^{CV} = \bar{F}_{m_1} + \alpha \left(\underbrace{\bar{G}_{m_2}}_{\text{estimates } \mathbb{E}[G]} - \bar{G}_{m_1} \right)$$

- ▶ There are two m 's now: m_1 plays the same role as m before; additionally have m_2 which gives the number of samples used to estimate the mean $\mathbb{E}[G]$ of the control variate
- ▶ Is this an unbiased estimator of $\mathbb{E}[F]$?
- ▶ Optimal choice of α ?
- ▶ Optimal choice of number of samples m_1 and m_2 to achieve MSE ϵ with minimal costs?

↪ board

$$\Theta = \frac{1}{m_1^*} \sum_{i=1}^{m_1^*} f(x_i) + \alpha \left(\frac{1}{m_2^*} \sum_{i=1}^{m_2^*} g(x_i) - \frac{1}{m_1^*} \sum_{i=1}^{m_1^*} g(x_i) \right)$$

Samples

$$x_1, \dots, x_{m_1^*}, x_{m_1^*+1}, \dots, x_{m_2^*} \quad m_1^* \leq m_2^*$$

Unbiasedness:

$$\Theta = \bar{F}_{m_1} + \alpha (\bar{G}_{m_2} - \bar{G}_{m_1})$$

$$E[\Theta] = E[\bar{F}_{m_1}] = E[f(x)]$$

MSE

$$\text{Var}[\Theta] = V[\bar{F}_{m_1}] + \alpha^2 V[\bar{G}_{m_2} - \bar{G}_{m_1}] +$$

$$2\alpha \text{Cov}[\bar{F}_{m_1}, \bar{G}_{m_2} - \bar{G}_{m_1}]$$

$$= \dots = V[\bar{F}_{m_1}] + \alpha^2 V[\bar{G}_{m_2}] + \alpha^2 V[\bar{G}_{m_1}]$$

$$- 2\alpha^2 C[\bar{G}_{m_2}, \bar{G}_{m_1}] - 2\alpha C[\bar{F}_{m_1}, \bar{G}_{m_1}]$$

$$+ 2\alpha C[\bar{F}_{m_1}, \bar{G}_{m_2}]$$

⋮

$$C[\bar{G}_{m_2}, \bar{G}_{m_1}] = \dots = \frac{1}{m_2} V[g(x)]$$

similar transformation:

$$\text{Var}[\Theta] = \frac{1}{m_1} (V[f(x)] + \alpha^2 V[g(x)] - 2\alpha \rho \sqrt{V[f(x)]V[g(x)]}) - \frac{1}{m_2} (\alpha^2 V[g(x)] - 2\alpha \rho \sqrt{V[f(x)]V[g(x)]})$$

$\min_{m_1, m_2, \alpha}$ $m_1 w_1 + m_2 w_2$ } this is the
 such that $V[\Theta_{m_1, m_2, \alpha}] \leq \epsilon$ } key step

Results: m_1^*, m_2^*, α^*

$$m_1^* = \frac{V[f(x)]}{\epsilon} (1 - (1 - \frac{1}{r^*}) \rho^2)$$

$$m_2^* = m_1^* r^*$$

$$C_\epsilon(\Theta) \leq m_1^* w_1 + m_2^* w_2 = \frac{V[f(x)]}{\epsilon} (1 - (1 - \frac{1}{r^*}) \rho^2) (w_1 + r^* w_2)$$

$$r^* = \sqrt{\frac{w_1 \rho^2}{w_2 (1 - \rho^2)}}$$

VR: Properties of control variate estimator

Consider two functions f and g with evaluation costs w_1 and w_2 , respectively. Assume $w_1 > w_2$. Consider the random variable X and set $\rho = \rho(f(X), g(X))$ and let σ_F and σ_G be the standard deviation of $f(X)$ and $g(X)$, respectively. Set

$$\alpha^* = \rho \frac{\sigma_F}{\sigma_G}, \quad r^* = \sqrt{\frac{w_1 \rho^2}{w_2 (1 - \rho^2)}}$$

The estimator

$$\theta = \frac{1}{m_1^*} \sum_{i=1}^{m_1^*} f(X_i) + \alpha^* \left(\frac{1}{m_2^*} \sum_{i=1}^{m_2^*} g(X_i) - \frac{1}{m_1^*} \sum_{i=1}^{m_1^*} g(X_i) \right)$$

is unbiased w.r.t. $\mathbb{E}[f(X)]$ if $X_1, \dots, X_{m_1}, X_{m_1+1}, \dots, X_{m_2} \sim X$ and achieves an MSE $e(\theta) \leq \epsilon$ with costs

$$c_\epsilon(\theta) \leq \frac{\sigma_F^2}{\epsilon} \left(1 - \left(1 - \frac{1}{r^*} \right) \rho^2 \right) (w_1 + r^* w_2)$$

if

$$m_1^* = \frac{\sigma_F^2}{\epsilon} \left(1 - \left(1 - \frac{1}{r^*} \right) \rho^2 \right), \quad m_2^* = m_1^* r^* .$$

The m_1^*, m_2^*, α^* are optimal in the sense that they minimize the costs $c_\epsilon(\theta)$.

Note that control variate reuses samples

VR: Comparison

Plain vanilla Monte Carlo

$$c_{\epsilon}(\theta) \leq \frac{\sigma_F^2}{\epsilon} w_1$$

Control variates if mean of control variate known

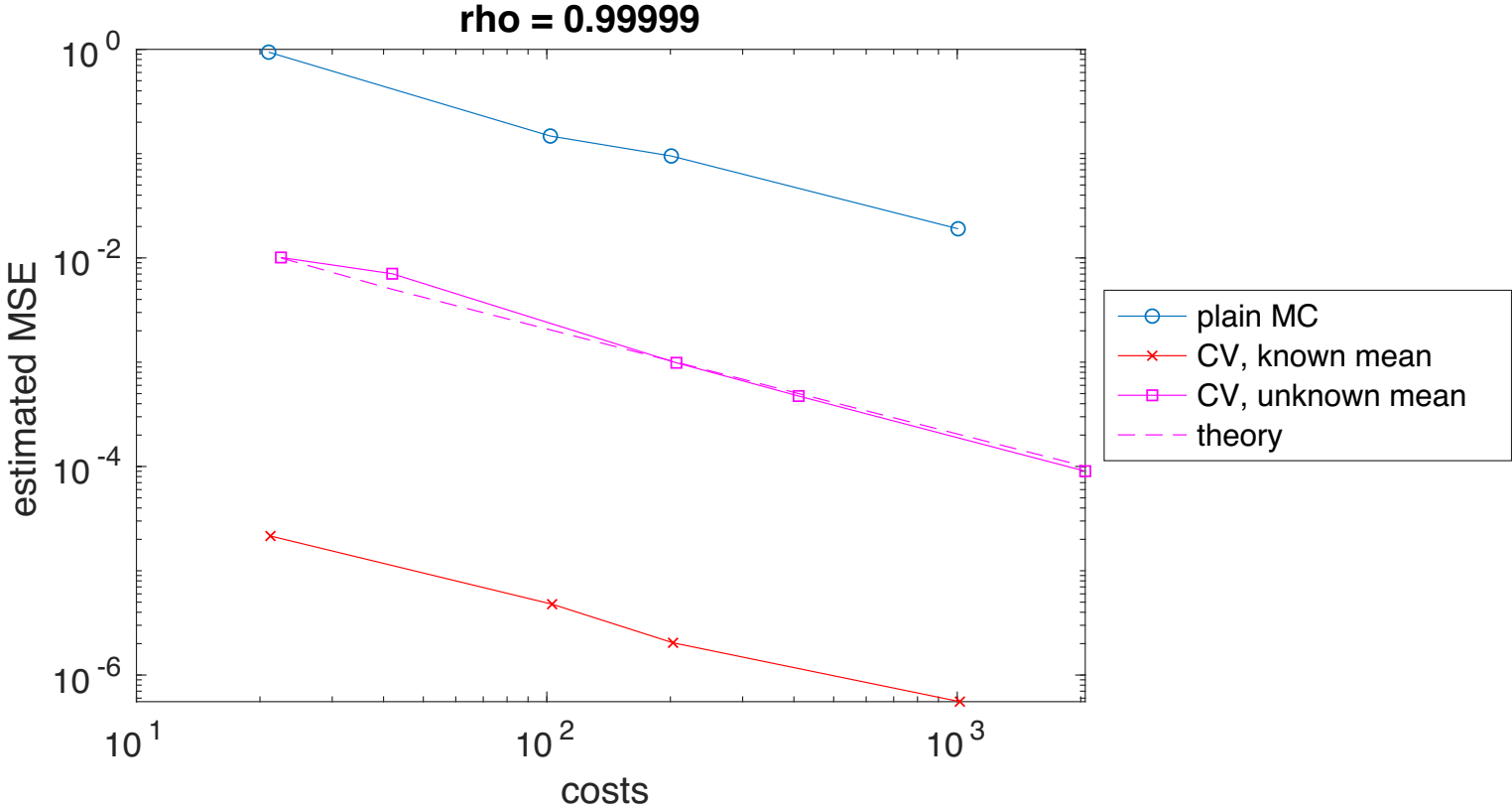
$$c_{\epsilon}(\theta) \leq \frac{\sigma_F^2}{\epsilon} (1 - \rho^2) (w_1 + w_2)$$

Control variates if mean of control variate unknown

$$c_{\epsilon}(\theta) \leq \frac{\sigma_F^2}{\epsilon} \left(1 - \left(1 - \frac{1}{r^*} \right) \rho^2 \right) (w_1 + r^* w_2)$$

Notice that the rate with respect to ϵ is the same but the constants change

VR: Example (cont'd)



Today

Last time

- ▶ Monte Carlo

Today

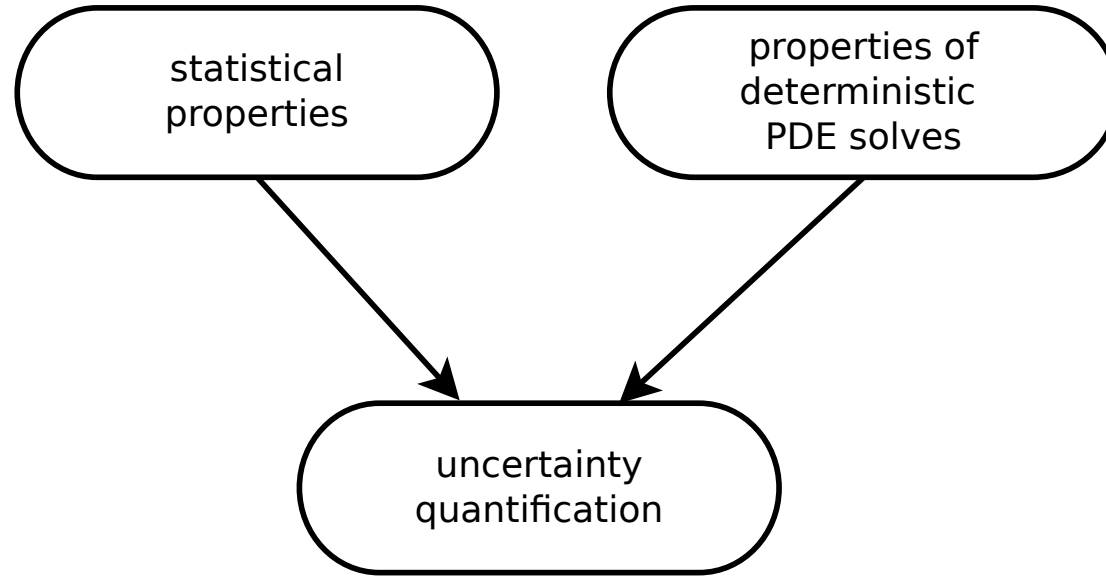
- ▶ Multi-level Monte Carlo

Announcements

- ▶ Homework 7 is posted and due Mon, Dec 9 before class (1 week)
- ▶ Next week, Mon Dec 9, recap of important topics—highly recommended!

What follows won't be on the final exam

Monte Carlo + Numerical Analysis

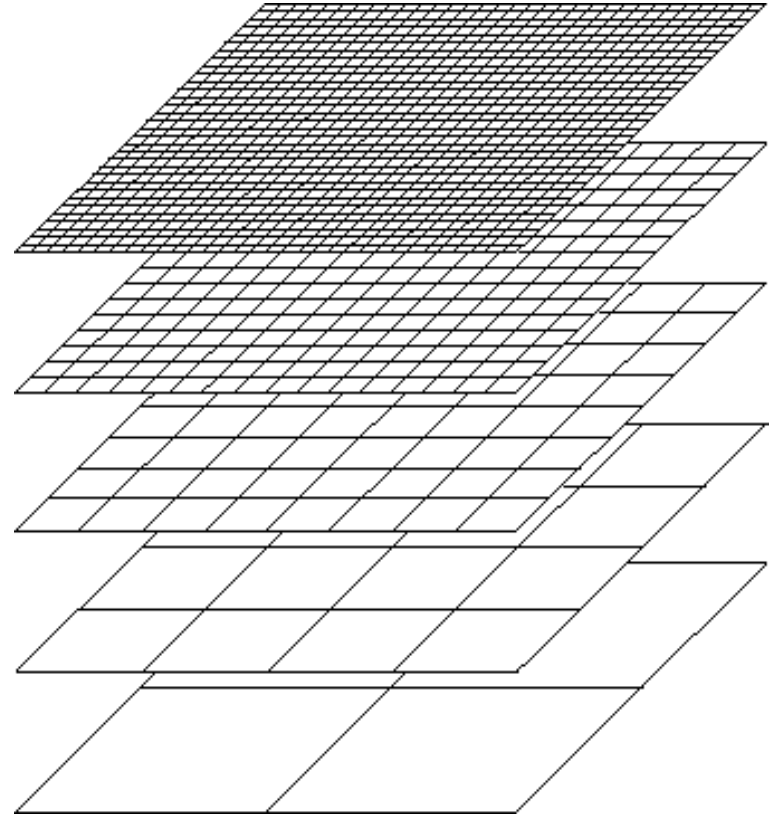


- ▶ Bounding **costs** instead of, e.g., number of samples
- ▶ Need to take into account numerical PDE solver costs

Level

Instead of describing the discretization with the number of degrees of freedom N , we describe it by the level $\ell \in \mathbb{N}$

- ▶ Level $\ell \in \mathbb{N}$, mesh width $h = 2^{-\ell}$
- ▶ Number of degrees of freedom in 1-dimensional spatial domain
 $N = 2^\ell + 1 \in \mathcal{O}(2^\ell)$
- ▶ In 2-dimensional spatial domain
 $N = (2^\ell + 1)^2 \in \mathcal{O}(4^\ell)$
- ▶ In 3-dimensional spatial domain
 $N = (2^\ell + 1)^3 \in \mathcal{O}(8^\ell)$
- ▶ General $N \in \mathcal{O}(s^\ell)$ with $s \in \mathbb{N}$



A first analysis of MC in the context of UQ

The following summary follows the analysis in [Cliffe et al., 2011].

To estimate the expectation $\mathbb{E}[Q]$ of a (random) quantity of interest Q , assume only approximations $Q_\ell \approx Q$ are computable, where $\ell \in \mathbb{N}$ is a discretization level such that

$$\lim_{\ell \rightarrow \infty} \mathbb{E}[Q_\ell] = \mathbb{E}[Q].$$

More precisely, we assume the error converges in mean (in distribution) with a rate $\alpha > 0$

$$|\mathbb{E}[Q - Q_\ell]| \lesssim s^{-\ell\alpha},$$

where \lesssim means up to constants independent of ℓ and α .

With increasing ℓ , the costs of solving the corresponding system may increase

$$c(Q_\ell) \lesssim s^{\ell\gamma}, \quad \gamma > 0$$

A first analysis of MC (cont'd)

- ▶ Let \bar{Q}_ℓ be an unbiased estimator of $\mathbb{E}[Q_\ell]$
- ▶ Summary:
 - ▶ Q is the QoI that is unavailable (“continuous solution of PDE solution”)
 - ▶ Q_ℓ is the approximation of Q (“numerical approximation”)
 - ▶ $\mathbb{E}[Q]$ is what we want to estimate (“no numerical error; no statistical error”)
 - ▶ $\mathbb{E}[Q_\ell]$ is mean of Q_ℓ (“no statistical error”)
 - ▶ \bar{Q}_ℓ is an unbiased estimator of $\mathbb{E}[Q_\ell]$ (“Monte Carlo estimator of $\mathbb{E}[Q_\ell]$ ”)
- ▶ The mean-squared error (MSE) of \bar{Q}_ℓ is

$$\mathbb{E}[(\bar{Q}_\ell - \mathbb{E}[Q])^2] = \underbrace{\text{Var}[\bar{Q}_\ell]}_{\text{variance}} + \underbrace{\mathbb{E}[Q_\ell - Q]^2}_{\text{bias}}$$

- ▶ Measure bias w.r.t. $\mathbb{E}[Q]$
 - ▶ Bias w.r.t. $\mathbb{E}[Q_\ell]$ is zero because \bar{Q}_ℓ unbiased by assumption
- ▶ How can we control the variance and bias?

A first analysis of MC (cont'd)

- ▶ Let \bar{Q}_ℓ be an unbiased estimator of $\mathbb{E}[Q_\ell]$
- ▶ Summary:
 - ▶ Q is the QoI that is unavailable (“continuous solution of PDE solution”)
 - ▶ Q_ℓ is the approximation of Q (“numerical approximation”)
 - ▶ $\mathbb{E}[Q]$ is what we want to estimate (“no numerical error; no statistical error”)
 - ▶ $\mathbb{E}[Q_\ell]$ is mean of Q_ℓ (“no statistical error”)
 - ▶ \bar{Q}_ℓ is an unbiased estimator of $\mathbb{E}[Q_\ell]$ (“Monte Carlo estimator of $\mathbb{E}[Q_\ell]$ ”)
- ▶ The mean-squared error (MSE) of \bar{Q}_ℓ is

$$\mathbb{E}[(\bar{Q}_\ell - \mathbb{E}[Q])^2] = \underbrace{\text{Var}[\bar{Q}_\ell]}_{\text{variance}} + \underbrace{\mathbb{E}[Q_\ell - Q]^2}_{\text{bias}}$$

- ▶ Measure bias w.r.t. $\mathbb{E}[Q]$
- ▶ Bias w.r.t. $\mathbb{E}[Q_\ell]$ is zero because \bar{Q}_ℓ unbiased by assumption
- ▶ How can we control the variance and bias?
 - ▶ Variance: Improve estimator \bar{Q}_ℓ of $\mathbb{E}[Q_\ell]$ → statistics
 - ▶ Bias: Improve accuracy of Q_ℓ w.r.t. Q → deterministic solver

Cost complexity of basic Monte Carlo

Assume that

- ▶ (A0) the variance $\text{Var}[Q_\ell]$ is constant w.r.t. ℓ ,
- ▶ (A1) $|\mathbb{E}[Q - Q_\ell]| \lesssim s^{-\ell\alpha}$, $\alpha > 0$,
- ▶ (A2) $c(Q_\ell) \lesssim s^{\ell\gamma}$, $\gamma > 0$.

Consider the Monte Carlo estimator $\bar{Q}_{\ell,m}^{\text{MC}}$ ($=\bar{Q}_{\ell,m}$) of $\mathbb{E}[Q_\ell]$ with $m \in \mathbb{N}$ iid copies $Q_\ell^{(1)}, \dots, Q_\ell^{(m)}$ of Q_ℓ . For $\epsilon > 0$, the Monte Carlo estimator $\bar{Q}_{\ell,m}^{\text{MC}}$ achieves the mean-squared error (MSE)

$$\mathbb{E}[(\bar{Q}_{\ell,m}^{\text{MC}} - \mathbb{E}[Q])^2] \leq \epsilon$$

with costs

$$c_\epsilon(\bar{Q}_{\ell,m}^{\text{MC}}) \lesssim \epsilon^{-1-\gamma/(2\alpha)}.$$

In terms of the root-mean-squared error (RMSE), the costs are bounded by $\epsilon^{-2-\gamma/\alpha}$.

know that

$$\text{Var}[\bar{Q}_{e,m}^{MC}] = \frac{\text{Var}[Q_e]}{m}$$

and thus

$$\mathbb{E}[(\bar{Q}_{e,m}^{MC} - \mathbb{E}[Q])^2] = \underbrace{\frac{\text{Var}[Q_e]}{m}}_{\substack{\text{depends on} \\ m \text{ only (A0)}}} + \underbrace{\mathbb{E}[Q_e - Q]^2}_{\substack{\text{only depends} \\ \text{on the level}}}$$

want to find l and m such that

$$\underline{\mathbb{E}[(\bar{Q}_{e,m}^{MC} - \mathbb{E}[Q])^2] \leq \varepsilon}$$

\Rightarrow balance bias and variance

(1) Variance

$$\frac{\text{Var}[Q_e]}{m} \leq \frac{\varepsilon}{2} \Leftrightarrow \underline{\underline{2 \text{Var}[Q_e] \varepsilon^{-1} \leq m}} \quad \text{(A0)}$$

(2) Bias

$$\mathbb{E}[Q - Q_e]^2 \stackrel{\text{(A1)}}{\leq} (c_1 s^{-l})^2 \stackrel{!}{\leq} \frac{\varepsilon}{2}$$

$$\Leftrightarrow \underline{\underline{\varepsilon^{-\frac{1}{2}} (\sqrt{2} c_1)^{\frac{1}{2}} \leq s^e}}$$

invoke (A2)

(A2)

$$C(\bar{Q}_{\epsilon, m}^{MC}) \leq \underbrace{m}_{\substack{\uparrow \\ \# \text{ samples}}} \underbrace{C_2 S^{\epsilon \eta}}_{\substack{\text{cost per sample}}}$$

to achieve $\leq \epsilon$, set $m = \frac{2 \text{Var}[Q_{\epsilon}] \epsilon^{-1}}{\epsilon^{\eta}}$ and
set level ℓ such that

$$S^{\ell} = \epsilon^{-\frac{1}{2\eta}} (\sqrt{2} c_1)^{\frac{1}{\eta}}$$

$$C_{\epsilon}(\bar{Q}_{\epsilon, m}^{MC}) \leq C_2 \underbrace{2 \text{Var}[Q_{\epsilon}] \epsilon^{-1}}_{= m} \underbrace{\epsilon^{-\frac{1}{2\eta}} (\sqrt{2} c_1)^{\frac{1}{\eta}}}_{= S^{\ell \eta}} \\ \leq \epsilon^{-1 - \frac{1}{2\eta}}$$

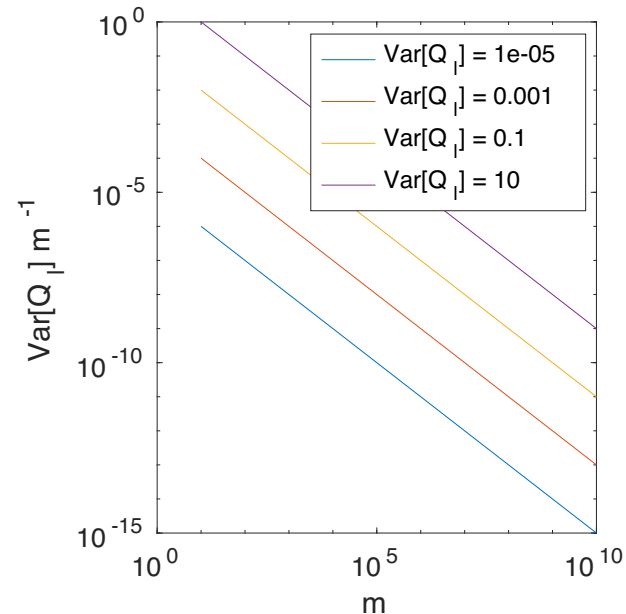
Improving plain Monte Carlo: Variance reduction

Convergence of plain vanilla Monte Carlo

$$\begin{aligned}\mathbb{E}[(\bar{Q}_{\ell,m}^{\text{MC}} - \mathbb{E}[Q])^2] &= \text{Var}[\bar{Q}_{\ell,m}^{\text{MC}}] + (\mathbb{E}[\bar{Q}_{\ell,m}^{\text{MC}}] - \mathbb{E}[Q])^2 \\ &= \underbrace{\frac{\text{Var}[Q_\ell]}{m}}_{\text{sampling error (variance)}} + \underbrace{(\mathbb{E}[Q_\ell] - \mathbb{E}[Q])^2}_{\text{model error (bias)}}\end{aligned}$$

- ▶ Convergence rate of plain vanilla Monte Carlo is $\mathcal{O}(m^{-1})$ in MSE
- ▶ The constant $\text{Var}[Q_\ell]$ plays a critical role

\Rightarrow it is worthwhile trying to reduce *constants* \Rightarrow variance reduction



MLMC: Multilevel Monte Carlo

- ▶ **Multilevel Monte Carlo** (MLMC) uses control variates in a judicious way
- ▶ First ideas for high-dimensional quadrature by Heinrich, 2000
Heinrich S. (2001) Multilevel Monte Carlo Methods. In: Margenov S., Waśniewski J., Yalamov P. (eds) Large-Scale Scientific Computing. LSSC 2001. Lecture Notes in Computer Science, vol 2179. Springer, Berlin, Heidelberg
- ▶ Independently discovered and popularized by Giles, 2007 in the context of stochastic differential equations in mathematical finance (one of the most influential papers in UQ!)
Michael B. Giles, Multilevel Monte Carlo Path Simulation, Operations Research 2008 56:3, 607-617
- ▶ First papers in the context of UQ
 - ▶ Cliffe, K.A., Giles, M.B., Scheichl, R. et al. Comput. Visual Sci. (2011) 14: 3.
<https://doi.org/10.1007/s00791-011-0160-x>
 - ▶ Barth, A., Schwab, C. & Zollinger, N. Numer. Math. (2011) 119: 123.
<https://doi.org/10.1007/s00211-011-0377-0>

MLMC: Multilevel Monte Carlo

- ▶ **Multilevel Monte Carlo** (MLMC) uses control variates in a judicious way
- ▶ First ideas for high-dimensional quadrature by Heinrich, 2000
Heinrich S. (2001) Multilevel Monte Carlo Methods. In: Margenov S., Waśniewski J., Yalamov P. (eds) Large-Scale Scientific Computing. LSSC 2001. Lecture Notes in Computer Science, vol 2179. Springer, Berlin, Heidelberg
- ▶ Independently discovered and popularized by Giles, 2007 in the context of stochastic differential equations in mathematical finance (one of the most influential papers in UQ!)
Michael B. Giles, Multilevel Monte Carlo Path Simulation, Operations Research 2008 56:3, 607-617
- ▶ First papers in the context of UQ
 - ▶ Cliffe, K.A., Giles, M.B., Scheichl, R. et al. *Comput. Visual Sci.* (2011) 14: 3. <https://doi.org/10.1007/s00791-011-0160-x>
 - ▶ Barth, A., Schwab, C. & Zollinger, N. *Numer. Math.* (2011) 119: 123. <https://doi.org/10.1007/s00211-011-0377-0>

MLMC: Multilevel estimator

- ▶ **Key idea:** use realizations of Q_ℓ on a **hierarchy of different levels**, i.e., for different values $\ell = 1, \dots, L$ of the discretization parameter (\rightarrow recall sparse grids)
- ▶ Make the following decomposition (telescoping sum)

$$\mathbb{E}[Q_{N_L}] = \mathbb{E}[Q_{N_0}] + \sum_{\ell=1}^L \mathbb{E}[Q_{N_\ell} - Q_{N_{\ell-1}}] = \sum_{\ell=0}^L \mathbb{E}[Y_\ell],$$

where $N_0 \in \mathbb{N}$ and $N_\ell = sN_{\ell-1}$ for $\ell = 1, \dots, L$ and $s \in \mathbb{N} \setminus \{1\}$ and

$$Y_0 = Q_{N_0}, \quad Y_\ell = Q_{N_\ell} - Q_{N_{\ell-1}}$$

- ▶ Given (unbiased) estimators \bar{Y}_{ℓ, m_ℓ} for $\mathbb{E}[Y_\ell]$, the estimator

$$\bar{Q}_{L, \mathbf{m}}^{\text{ML}} = \sum_{\ell=0}^L \bar{Y}_{\ell, m_\ell}$$

is a **multilevel estimator** of Q with $\mathbf{m} = [m_0, \dots, m_L]$

MLMC: Multilevel Monte Carlo estimator

- ▶ All estimators \bar{Y}_{ℓ, m_ℓ} sampled independently, then

$$\text{Var}[\bar{Q}_{L, \mathbf{m}}^{\text{ML}}] = \sum_{\ell=0}^L \text{Var}[\bar{Y}_{\ell, m_\ell}]$$

- ▶ If each \bar{Y}_ℓ is a plain Monte Carlo estimator

$$\bar{Y}_{0, m_0} = \frac{1}{m_0} \sum_{i=1}^{m_0} Q_{N_0}^{(i)}, \quad Q_{N_0}^{(i)} \sim Q_{N_0},$$

and

$$\bar{Y}_{\ell, m_\ell} = \frac{1}{m_\ell} \sum_{i=1}^{m_\ell} Q_{N_\ell}^{(i)} - Q_{N_{\ell-1}}^{(i)},$$

one obtains a **multilevel Monte Carlo estimator**

- ▶ The MSE of the multilevel Monte Carlo estimator is

$$\mathbb{E} \left[\left(\bar{Q}_{L, \mathbf{m}}^{\text{MLMC}} - \mathbb{E}[Q] \right)^2 \right] = \sum_{\ell=0}^L \frac{\text{Var}[Y_\ell]}{m_\ell} + \mathbb{E}[Q_L - Q]^2$$

MLMC: Variance reduction

$$\mathbb{E} \left[\left(\bar{Q}_{L,m}^{\text{MLMC}} - \mathbb{E}[Q] \right)^2 \right] = \sum_{\ell=0}^L \frac{\text{Var}[Y_\ell]}{m_\ell} + \mathbb{E}[Q_L - Q]^2$$

Why do we have hope for variance reduction (=lower cost for same variance)?

MLMC: Variance reduction

$$\mathbb{E} \left[\left(\bar{Q}_{L,m}^{\text{MLMC}} - \mathbb{E}[Q] \right)^2 \right] = \sum_{\ell=0}^L \frac{\text{Var}[Y_\ell]}{m_\ell} + \mathbb{E}[Q_L - Q]^2$$

Why do we have hope for variance reduction (=lower cost for same variance)?

- ▶ As we coarsen the problem, the cost per sample decays rapidly from level to level
→ observed this already in the control variates example → sampling gets cheaper and cheaper

MLMC: Variance reduction

$$\mathbb{E} \left[\left(\bar{Q}_{L,m}^{\text{MLMC}} - \mathbb{E}[Q] \right)^2 \right] = \sum_{\ell=0}^L \frac{\text{Var}[Y_\ell]}{m_\ell} + \mathbb{E}[Q_L - Q]^2$$

Why do we have hope for variance reduction (=lower cost for same variance)?

- ▶ As we coarsen the problem, the cost per sample decays rapidly from level to level → observed this already in the control variates example → sampling gets cheaper and cheaper
- ▶ Since $Q_{N_\ell} \rightarrow Q$, we have $\text{Var}[Y_\ell] = \text{Var}[Q_{N_\ell} - Q_{N_{\ell-1}}] \rightarrow 0$ going to 0 fast for $\ell \rightarrow \infty$, **allowing for smaller and smaller sample sizes m_ℓ to estimate the difference $Q_{N_\ell} - Q_{N_{\ell-1}}$ on finer and finer (more and more expensive) levels.**

MLMC: Cost complexity

Theorem (Giles 2008; Cliffe, Giles, Scheichl, Teckentrup 2011)

Let $\bar{Y}_\ell = \bar{Y}_{\ell, m_\ell}^{MC}$ and suppose that there are positive constants $\alpha, \beta, \gamma > 0$ such that $\alpha \geq \frac{1}{2} \min(\beta, \gamma)$ and

- ▶ (A1) $|\mathbb{E}[Q_{N_\ell} - Q]| \lesssim N_\ell^{-\alpha}$
- ▶ (A2) $\text{Var}[Y_\ell] \lesssim N_\ell^{-\beta}$
- ▶ (A3) $w_\ell \lesssim N_\ell^\gamma$.

Then, for any $\sqrt{\epsilon} < e^{-1}$, there exist a value $L \in \mathbb{N}$ and $\mathbf{m} = [m_0, \dots, m_L]$ such that

$$\mathbb{E} \left[\left(\bar{Q}_{L, \mathbf{m}}^{MLMC} - \mathbb{E}[Q] \right) \right] < \epsilon,$$

and

$$c_\epsilon(\bar{Q}_{L, \mathbf{m}}^{MLMC}) \lesssim \begin{cases} \epsilon^{-1} & \text{if } \beta > \gamma, \\ \epsilon^{-1} (\log \epsilon^{1/2})^2, & \text{if } \beta = \gamma, \\ \epsilon^{-1 - (\gamma - \beta)/(2\alpha)}, & \text{if } \beta < \gamma. \end{cases}$$

$$(A1) |E[Q_{N_e} - Q]| \lesssim c_1 N_e^{-\alpha}$$

$$(A2) \text{Var}[Y_e] \lesssim c_2 N_e^{-\beta}$$

$$(A3) w_e \lesssim c_3 N_e^{\gamma}$$

$$|E[(\bar{Q}_{L,m}^{\text{-MLMC}} - E[Q])^2]| < \varepsilon$$

and costs

$$c_\varepsilon(\bar{Q}_{L,m}^{\text{-MLMC}}) \lesssim$$

$$\left\{ \begin{array}{l} \varepsilon^{-1}, \quad \text{if } \underline{\beta} > \underline{\gamma} \\ \vdots \end{array} \right.$$

Recall $N_e = s N_{e-1}$, $e = 1, \dots, \underline{L}$, $N_0 = 1$

Approach: for each case, bound bias by $\frac{\varepsilon}{2}$ and variance by $\frac{\varepsilon}{2}$ and then bound costs

Bias: Set

$$\underline{L} = \left\lceil \alpha^{-1} \log_s (\varepsilon^{-1/2} \sqrt{2} c_1) \right\rceil$$

$$< \alpha^{-1} \log_s (\varepsilon^{-1/2} \sqrt{2} c_1) + 1$$

$$\underline{E}[Q_L - Q]^2 \leq \frac{\varepsilon}{2}$$

Variance

case $\beta > \gamma$: Set

$$\underline{m_e} = \left[\frac{1}{2} \varepsilon^{-1} c_2 (1 - s^{-(\beta-\gamma)/2})^{-1} s^{-(\beta+\gamma)/2} \right]$$

$$\sum_{e=0}^L \underline{\text{Var}[\bar{y}_{e,m}]} = \sum_{e=0}^L \underbrace{m_e^{-1}}_{\text{balance}} \underbrace{\text{Var}[y_e]}_{\text{balance}} \leq \sum_{e=0}^L m_e^{-1} c_2 s^{-\beta e}$$

$$\leq \frac{1}{2} \varepsilon (1 - s^{-(\beta-\gamma)/2}) \sum_{e=0}^L c_2 \underbrace{s^{-\beta e} s^{-(\beta+\gamma)/2}}_{s^{-(\beta-\gamma)/2}}$$

$$= \frac{1}{2} \varepsilon (1 - \cancel{s^{-(\beta-\gamma)/2}}) \sum_{e=0}^L (s^{-(\beta-\gamma)/2})^e$$

$$\frac{1 - \underline{s}^{-(\beta-\gamma)/2(L+1)}}{1 - \cancel{s^{-(\beta-\gamma)/2}}}$$

$$\leq \underline{\frac{1}{2} \varepsilon}$$

Costs:

$$C(\bar{Q}_{L,m}^{MLMC}) = \sum_{e=0}^L \underline{w_e} m_e$$

$$\leq \sum_{\ell=0}^L c_3 s^{\ell \gamma} (2 \varepsilon^{-1} c_2 (1 - s^{-(\beta-\gamma)/2})^{-1} s^{-(\beta+\gamma)\ell/2 + 1})$$

$$\leq c_3 2 \varepsilon^{-1} c_2 (1 - s^{-(\beta-\gamma)/2})^{-1}$$

$$\underbrace{\sum_{\ell=0}^L s^{-(\beta+\gamma)\ell/2} s^{\ell \gamma}}_{s^{-(\beta-\gamma)\ell/2}} + c_3 \sum_{\ell=0}^L s^{\ell \gamma}$$

geo. series

$$\left. \frac{1 - s^{-(\beta-\gamma)(L+1)/2}}{(1 - s^{-(\beta-\gamma)/2})} \right\}$$

$$= c_3 2 \varepsilon^{-1} c_2 (1 - s^{-(\beta-\gamma)/2})^{-2} \underbrace{(1 - s^{-(\beta-\gamma)(L+1)/2})}_{< 1} + c_3 \varepsilon \dots$$

> 0

$$\leq c \varepsilon^{-1} + c_3 \sum_{\ell=0}^L s^{\ell \gamma} \leq c \varepsilon^{-1} \approx \varepsilon^{-1}$$

because of rounding $\Gamma 7$

→ technical lemma

for the other cases: $\beta > \gamma, \beta = \gamma$

MLMC: Interpreting the cost complexity

For idealized flows through a porous medium (Darcy flow):

- ▶ Numerically observed deterministic error \Rightarrow rate $\alpha \approx 3/8$
- ▶ Numerically observed cost/sample \Rightarrow rate $\gamma \approx 1$

Consideration generalization of the problem in $d = 1, 2, 3$ dimensional spatial domain. It has been observed $\alpha \approx 3/4d^{-1}$, $\gamma \approx 1$, $\beta \approx 2\alpha$.

Using the complexity theorems of MC and MLMC and the costs per sample $\epsilon^{-\gamma/(2\alpha)}$ obtain

d	MC	MLMC	per sample on finest level
1	$\mathcal{O}(\epsilon^{-5/3})$	$\mathcal{O}(\epsilon^{-1})$	$\mathcal{O}(\epsilon^{-2/3})$
2	$\mathcal{O}(\epsilon^{-7/3})$	$\mathcal{O}(\epsilon^{-4/3})$	$\mathcal{O}(\epsilon^{-4/3})$
3	$\mathcal{O}(\epsilon^{-3})$	$\mathcal{O}(\epsilon^{-2})$	$\mathcal{O}(\epsilon^{-2})$

MLMC: Interpreting the cost complexity

For idealized flows through a porous medium (Darcy flow):

- ▶ Numerically observed deterministic error \Rightarrow rate $\alpha \approx 3/8$
- ▶ Numerically observed cost/sample \Rightarrow rate $\gamma \approx 1$

Consideration generalization of the problem in $d = 1, 2, 3$ dimensional spatial domain. It has been observed $\alpha \approx 3/4d^{-1}$, $\gamma \approx 1$, $\beta \approx 2\alpha$.

Using the complexity theorems of MC and MLMC and the costs per sample $\epsilon^{-\gamma/(2\alpha)}$ obtain

d	MC	MLMC	per sample on finest level
1	$\mathcal{O}(\epsilon^{-5/3})$	$\mathcal{O}(\epsilon^{-1})$	$\mathcal{O}(\epsilon^{-2/3})$
2	$\mathcal{O}(\epsilon^{-7/3})$	$\mathcal{O}(\epsilon^{-4/3})$	$\mathcal{O}(\epsilon^{-4/3})$
3	$\mathcal{O}(\epsilon^{-3})$	$\mathcal{O}(\epsilon^{-2})$	$\mathcal{O}(\epsilon^{-2})$

MLMC costs asymptotically the same as **one** deterministic solve to accuracy ϵ for $d > 1 \rightarrow$ “UQ is for free”

Today •

Today

- ▶ Review of *some* important topics

Announcements

- ▶ Be 10min earlier in the room for the final exam
- ▶ Mind that the room has been updated; check Albert
- ▶ Double check that all homeworks are entered correctly in brightspace

There are many opportunities for asking questions

Mon, Dec 2	office hour
Wed, Dec 4	lecture (virtual) [not part of final exam]
Fri, Dec 6	office hour (grader)
Mon, Dec 9	recap and Q&A
Mon, Dec 9	office hour
Wed, Dec 11	extra office hour, 6.10pm ET (WWH 421) (all HWs graded; no HW re-grading after Thu, Dec 12)
Mon, Dec 16	final exam

- ▶ Condition and stability
- ▶ Linear systems and linear least-squares problems
- ▶ Eigenproblems
- ▶ Iterative methods for linear systems and nonlinear systems
- ▶ Interpolation
- ▶ Quadrature

Condition of a problem

- ▶ Consider a generic problem: given F and data/input x , find output y such that

$$F(x, y) = 0$$

- ▶ Let's assume there is a unique solution so that we can write

$$y = f(x),$$

for a function f in the following

- ▶ Well-posed: Unique solution + If we perturb the input x a little bit, the solution y gets perturbed by a small amount.
- ▶ Otherwise, the problem is ill-posed; no numerical method can help with that.
(What should we do in such a situation?)

Condition of a problem

- ▶ Consider a generic problem: given F and data/input x , find output y such that

$$F(x, y) = 0$$

- ▶ Let's assume there is a unique solution so that we can write

$$y = f(x),$$

for a function f in the following

- ▶ Well-posed: Unique solution + If we perturb the input x a little bit, the solution y gets perturbed by a small amount.
- ▶ Otherwise, the problem is ill-posed; no numerical method can help with that. (What should we do in such a situation? \rightsquigarrow change the problem)

Condition of a problem (cont'd)

- ▶ Absolute condition number at x is

$$\kappa_{\text{abs}} = \lim_{\delta \rightarrow 0} \sup_{\|x - \hat{x}\| \leq \delta} \frac{\|f(x) - f(\hat{x})\|}{\|x - \hat{x}\|}$$

- ▶ Relative condition number at x is

$$\kappa_{\text{rel}} = \lim_{\delta \rightarrow 0} \sup_{\|x - \hat{x}\| \leq \delta} \frac{\|f(x) - f(\hat{x})\| / \|f(x)\|}{\|x - \hat{x}\| / \|x\|}$$

- ▶ If f is differentiable in x , then

$$\kappa_{\text{abs}} = \|f'(x)\| \quad \kappa_{\text{rel}} = \frac{\|x\|}{\|f(x)\|} \|f'(x)\|,$$

where $\|f'(x)\|$ is the norm of the Jacobian $f'(x)$ in the operator norm

$$\|A\| = \sup_{x \neq 0} \frac{\|Ax\|}{\|x\|} = \sup_{\|x\|=1} \|Ax\|$$

Revisiting stability

An algorithm \tilde{f} for a problem f is backward stable if for each $x \in X$ we have $\tilde{f}(x) = f(\tilde{x})$ for an \tilde{x} with

$$\frac{\|\tilde{x} - x\|}{\|x\|} \in \mathcal{O}(u),$$

where u is the roundoff unit

- ▶ Recall that, loosely speaking, the symbol $\mathcal{O}(u)$ means “on the order of the roundoff unit.”
- ▶ By allowing $u \rightarrow 0$ (which is implied here by the \mathcal{O}), we consider an idealization of a computer (in practice, u is fixed). So what we mean is that the error should decrease in proportion to u or faster.

Stability + condition

Suppose a backward stable algorithm is applied to solve a problem $f : X \rightarrow Y$ with relative condition number κ . Then, the relative errors satisfy

$$\frac{\|\tilde{f}(x) - f(x)\|}{\|f(x)\|} \in \mathcal{O}(\kappa(x)u).$$

Condition of solving system of linear equations

Recall that we derived the condition number $\kappa(\mathbf{A}) = \|\mathbf{A}\| \|\mathbf{A}^{-1}\|$ of a matrix \mathbf{A}

Problem 1: fix A , consider $b \mapsto A^{-1}b$

$$f: \mathbb{R}^n \rightarrow \mathbb{R}^n, \quad f(b) = A^{-1}b$$

$$k_{\text{obs}} = \|A^{-1}\|$$

$$k_{\text{rel}} = \frac{\|b\|}{\|A^{-1}b\|} \quad \|A^{-1}\| \leq \|A\| \|A^{-1}\| = \kappa(A)$$

Problem 2: fix b , $A \mapsto A^{-1}b$

$$k_{\text{rel}} \leq \|A\| \|A^{-1}\| = \kappa(A)$$

Problem 1: Show that $\rho(A) \leq \|A\|$ for any induced matrix norm, where $\rho(A)$ is the spectral radius of A (the largest absolute eigenvalue).

$$\|A\| = \sup_{x \neq 0} \frac{\|Ax\|}{\|x\|} \geq \frac{\|Ax_1\|}{\|x_1\|} = \frac{\|\rho(A)x_1\|}{\|x_1\|} = \rho(A)$$

Solving systems of linear equations

Gauss elimination and LU factorization

$$\begin{bmatrix} a_{11}^{(1)} & a_{12}^{(1)} & a_{13}^{(1)} \\ a_{21}^{(1)} & a_{22}^{(1)} & a_{23}^{(1)} \\ a_{31}^{(1)} & a_{32}^{(1)} & a_{33}^{(1)} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} = \begin{bmatrix} b_1^{(1)} \\ b_2^{(1)} \\ b_3^{(1)} \end{bmatrix}$$

multiply first row
by l_{i1} and subtract

$$\begin{bmatrix} b_1^{(1)} \\ b_2^{(1)} \\ b_3^{(1)} \end{bmatrix} \leftarrow \begin{bmatrix} l_{21} \\ l_{31} \end{bmatrix} = \begin{bmatrix} \frac{a_{21}^{(1)}}{a_{11}^{(1)}} \\ \frac{a_{31}^{(1)}}{a_{11}^{(1)}} \end{bmatrix}$$

$$\begin{bmatrix} a_{11}^{(1)} & a_{12}^{(1)} & a_{13}^{(1)} \\ a_{21}^{(1)} - l_{21} a_{11}^{(1)} & a_{22}^{(1)} - l_{21} a_{12}^{(1)} & a_{23}^{(1)} - l_{21} a_{13}^{(1)} \\ 0 & a_{32}^{(1)} - l_{31} a_{12}^{(1)} & a_{33}^{(1)} - l_{31} a_{13}^{(1)} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} = \begin{bmatrix} b_1^{(1)} \\ b_2^{(1)} - l_{21} b_1^{(1)} \\ b_3^{(1)} - l_{31} b_1^{(1)} \end{bmatrix}$$

$a_{21}^{(1)} - l_{21} a_{11}^{(1)} = a_{21}^{(1)} - \frac{a_{21}^{(1)}}{a_{11}^{(1)}} a_{11}^{(1)} = 0$

$$\begin{bmatrix} b_1^{(1)} \\ b_2^{(1)} - l_{21} b_1^{(1)} \\ b_3^{(1)} - l_{31} b_1^{(1)} \end{bmatrix}$$

- ▶ For any square (regular or singular) matrix \mathbf{A} , partial (row) pivoting ensures existence of

$$\mathbf{PA} = \mathbf{LU}$$

where \mathbf{P} is a permutation matrix

- ▶ Furthermore, pivoting (w.r.t. $\max |a_{ij}|$) leads to a backward stable algorithm.
- ▶ Once an LU factorization is available, solving a linear system is cheap:

$$\mathbf{Ax} = \mathbf{LUx} = \mathbf{L(Ux)} = \mathbf{Ly} = \mathbf{b}$$

or

$$\mathbf{PAx} = \mathbf{LUx} = \mathbf{L(Ux)} = \mathbf{Ly} = \mathbf{Pb}$$

- ▶ Solve for \mathbf{y} using *forward substitution*
- ▶ Solve for \mathbf{x} by using *backward substitution* $\mathbf{Ux} = \mathbf{y}$

Recap: Costs

For forward [backward] substitution at step k there are $\approx k [(n - k)]$ multiplications and subtractions plus a few divisions. The total over all n steps is

$$\sum_{k=1}^n k \in \mathcal{O}(n^2)$$

\rightsquigarrow the number of floating-point operations (FLOPs) scales as $\mathcal{O}(n^2)$

For Gaussian elimination, at step k , there are $\approx (n - k)^2$ operations. Thus, the total scales as

$$\sum_{k=1}^n (n - k)^2 \in \mathcal{O}(n^3)$$

Summary:

- ▶ Directly applying Gaussian elimination (=LU + fwd/bwd) scales as $\mathcal{O}(n^3)$
- ▶ Computing LU decomposition scales as $\mathcal{O}(n^3)$
- ▶ Forward/backward substitution scales as $\mathcal{O}(n^2)$
- ▶ LU + forward/backward scales as $\mathcal{O}(n^3)$ \rightsquigarrow can *reuse* LU for other b

Linear least-squares

Consider non-square matrices $A \in \mathbb{R}^{m \times n}$ with $m \geq n$ and $\text{rank}(A) = n$. Then the system

$$Ax = b$$

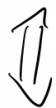
does, in general, not have a solution (more equations than unknowns). We thus instead solve a minimization problem

$$\min_x \|Ax - b\|^2 = \min_x \Phi(\mathbf{x})$$

$$\min_{x \in \mathbb{R}^n} \|Ax - b\|_2^2$$

geometric sketch

$$\bar{x} = \arg \min_x \|Ax - b\|_2^2$$



column
space of A

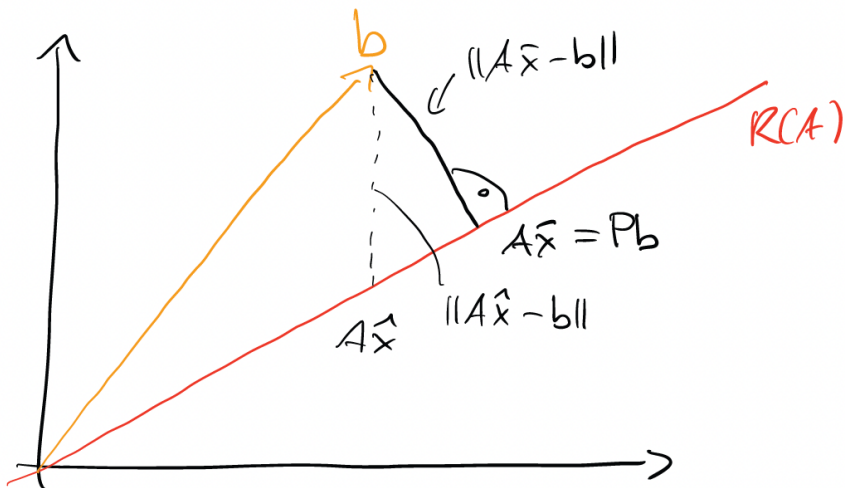
$$A\bar{x} - b \perp R(A)$$

$$\begin{cases} m \geq n \\ Ax = b \end{cases}$$

(idea: project b
onto the range
of A)

$$Pb$$

$$Ax = Pb$$



Linear least-squares problems

Now for the least-squares problem $\|\mathbf{Ax} - \mathbf{b}\|_2$. The relative condition number κ in the Euclidean norm is bounded by

- ▶ With respect to perturbations in \mathbf{b} :

$$\kappa \leq \frac{\kappa_2(A)}{\cos(\theta)}$$

- ▶ With respect to perturbations in \mathbf{A} :

$$\kappa \leq \kappa_2(A) + \kappa_2(A)^2 \tan(\theta)$$

Small residual problems, small angle θ $\cos(\theta) \approx 1$, $\tan(\theta) \approx 0$: behavior similar to linear system.

Large residual problems, large angle θ $\cos(\theta) \ll 1$, $\tan(\theta) \approx 1$: behavior very different from linear system because $\kappa_2(A)^2$ shows up

One would like to avoid the multiplication $A^T A$ (normal equations) and use a suitable factorization of A that avoids solving the normal equation directly:

One would like to avoid the multiplication $A^T A$ (normal equations) and use a suitable factorization of A that avoids solving the normal equation directly:

$$A = QR = [Q_1, Q_2] \begin{bmatrix} R_1 \\ 0 \end{bmatrix} = Q_1 R_1,$$

where $Q \in \mathbb{R}^{m \times m}$ is an orthonormal matrix ($QQ^T = I$), and $R \in \mathbb{R}^{m \times n}$ consists of an upper triangular matrix and a block of zeros.

How can the QR factorization be used to solve the least-squares problem?

One would like to avoid the multiplication $A^T A$ (normal equations) and use a suitable factorization of A that avoids solving the normal equation directly:

$$A = QR = [Q_1, Q_2] \begin{bmatrix} R_1 \\ 0 \end{bmatrix} = Q_1 R_1,$$

where $Q \in \mathbb{R}^{m \times m}$ is an orthonormal matrix ($QQ^T = I$), and $R \in \mathbb{R}^{m \times n}$ consists of an upper triangular matrix and a block of zeros.

How can the QR factorization be used to solve the least-squares problem?

$$\begin{aligned} \min_x \|Ax - b\|^2 &= \min_x \|Q^T(Ax - b)\|^2 &&= \min_x \left\| \begin{bmatrix} b_1 - R_1 x \\ b_2 \end{bmatrix} \right\|^2, \\ &= \min_x \|b_1 - R_1 x\|^2 + \|b_2\|^2 \end{aligned}$$

where $Q^T b = \begin{bmatrix} b_1 \\ b_2 \end{bmatrix}$.

Thus, the least squares solution is $x = R_1^{-1} b_1$ and the residual is $\|b_2\|$.

Problem 2: What is the following algorithm doing when applied to an $m \times n$ matrix A and is it a good idea to use it?

for $j = 1, \dots, n$

- ▶ $v_j = A_{:,j}$
- ▶ $R_{1:j-1,j} = Q_{:,1:j-1}^T A_{:,j}$
- ▶ $v_j = v_j - Q_{:,1:j-1} R_{1:j-1,j}$
- ▶ $R_{jj} = \|v_j\|_2$
- ▶ $Q_{:,j} = v_j / R_{jj}$

Problem 2: What is the following algorithm doing when applied to an $m \times n$ matrix A and is it a good idea to use it?

for $j = 1, \dots, n$

- ▶ $v_j = A_{:,j}$
- ▶ $R_{1:j-1,j} = Q_{:,1:j-1}^T A_{:,j}$
- ▶ $v_j = v_j - Q_{:,1:j-1} R_{1:j-1,j}$
- ▶ $R_{jj} = \|v_j\|_2$
- ▶ $Q_{:,j} = v_j / R_{jj}$

Computes the reduced QR decomposition with Gram-Schmidt. This is not modified Gram-Schmidt, so the columns in Q can be far from orthogonal.

Instead of directly computing

$$\mathbf{v}_j = \mathbf{a}_j - (\mathbf{q}_1^T \mathbf{a}_j) \mathbf{q}_1 - (\mathbf{q}_2^T \mathbf{a}_j) \mathbf{q}_2 - \cdots - (\mathbf{q}_{j-1}^T \mathbf{a}_j) \mathbf{q}_{j-1}$$

based on \mathbf{a}_j , the *modified* Gram-Schmidt procedure computes \mathbf{v}_j iteratively

$$\mathbf{v}_j^{(1)} = \mathbf{a}_j,$$

$$\mathbf{v}_j^{(2)} = \mathbf{v}_j^{(1)} - \mathbf{q}_1 \mathbf{q}_1^T \mathbf{v}_j^{(1)}, \quad \text{"subtract from } \mathbf{v}_j^{(1)} \text{ what is already in } \mathbf{q}_1 \text{"}$$

$$\mathbf{v}_j^{(3)} = \mathbf{v}_j^{(2)} - \mathbf{q}_2 \mathbf{q}_2^T \mathbf{v}_j^{(2)}, \quad \text{"subtract from } \mathbf{v}_j^{(2)} \text{ what is already in } \mathbf{q}_2 \text{"}$$

⋮

$$\mathbf{v}_j = \mathbf{v}_j^{(j)} = \mathbf{v}_j^{(j-1)} - \mathbf{q}_{j-1} \mathbf{q}_{j-1}^T \mathbf{v}_j^{(j-1)}$$

Computing a QR factorization with the modified Gram-Schmidt procedure is stabler than with the classical Gram-Schmidt procedure. However, even the modified Gram-Schmidt procedure can lead to vectors $\mathbf{q}_1, \dots, \mathbf{q}_n$ that are far from orthogonal if the condition number of \mathbf{A} is large (see, Golub et al., *Matrix Computations*, Section 5.2.9)

Eigenproblems

For a matrix $A \in \mathbb{C}^{n \times n}$ (potentially real), we want to find $\lambda \in \mathbb{C}$ and $\mathbf{x} \neq 0$ such that

$$A\mathbf{x} = \lambda\mathbf{x}.$$

Eigenproblems

For a matrix $A \in \mathbb{C}^{n \times n}$ (potentially real), we want to find $\lambda \in \mathbb{C}$ and $\mathbf{x} \neq 0$ such that

$$A\mathbf{x} = \lambda\mathbf{x}.$$

Let $\mathbf{A} \in \mathbb{C}^{n \times n}$ be diagonalizable matrix and λ_1 be a simple eigenvalue with

$$|\lambda_1| > |\lambda_2| \geq \dots \geq |\lambda_n|$$

Let \mathbf{x}_0 be an initial guess that is not orthogonal to the eigenspace of λ_1 , then \mathbf{x}_k obtained via the iterations

$$\mathbf{z}_{k+1} = \mathbf{A}\mathbf{x}_k \tag{1}$$

$$\mathbf{x}_{k+1} = \mathbf{z}_{k+1} / \|\mathbf{z}_{k+1}\|_2 \tag{2}$$

will converge to the normalized eigenvector of \mathbf{A} corresponding to λ_1 for $k \rightarrow \infty$.

This process is called the power method. **How did we proof convergence?**

Problem 3: Computing eigenvectors

Problem 3: Starting with $p_0 = q_0 = 1$, define the iteration

$$p_{n+1} = p_n + q_n$$

$$q_{n+1} = p_{n+1} + \underbrace{p_n}$$

for $n = 0, 1, 2, \dots$. The ratio q_n/p_n converges to $\sqrt{2}$ as $n \rightarrow \infty$.

Prove convergence of this algorithm using the power method. I.e., write the problem as an eigenvalue problem and show that the power method for this eigenvalue problem converges to an eigenvector and deduce from this eigenvector that $q_n/p_n \rightarrow \sqrt{2}$.

(*Hint:* $\det(A) = ad - bc$ for a 2×2 matrix with $A = \begin{bmatrix} a & b \\ c & d \end{bmatrix}$.)

$$q_0 = p_0 = 1$$

$$p_{n+1} = p_n + q_n$$

$$q_{n+1} = p_{n+1} + p_n$$

Write as matrix

$$\begin{bmatrix} p_{n+1} \\ q_{n+1} \end{bmatrix} = \begin{bmatrix} 1 & 1 \\ 2 & 1 \end{bmatrix} \begin{bmatrix} p_n \\ q_n \end{bmatrix}$$

Compute eigenvalues of A

$$\begin{aligned} \det(A - \lambda I) &= \det \begin{pmatrix} 1-\lambda & 1 \\ 2 & 1-\lambda \end{pmatrix} \\ &= (1-\lambda)(1-\lambda) - 1 \cdot 2 \\ &= \lambda^2 - 2\lambda - 1 \end{aligned}$$

$$\lambda_{1/2} = \dots = 1 \pm \sqrt{2}$$

\Rightarrow dominant EV: $|1 + \sqrt{2}| > |1 - \sqrt{2}|$

Compute EV for $1 + \sqrt{2} = \lambda_1$

$$0 = (A - \lambda_1 I) \vec{v}_1 = \begin{bmatrix} 1 - (1 + \sqrt{2}) & 1 \\ 2 & 1 - (1 + \sqrt{2}) \end{bmatrix} \begin{bmatrix} v_1 \\ v_2 \end{bmatrix}$$

$$= \begin{bmatrix} -\sqrt{2} & 1 \\ 2 & -\sqrt{2} \end{bmatrix} \begin{bmatrix} v_1 \\ v_2 \end{bmatrix}$$

$$v_1 = \frac{1}{\sqrt{2}} v_2 \quad (*)$$

obtain $\begin{bmatrix} v_1 \\ v_2 \end{bmatrix} = \begin{bmatrix} \frac{1}{\sqrt{2}} \\ 1 \end{bmatrix}$ as an EV

starting point $\begin{bmatrix} p_0 \\ q_0 \end{bmatrix} = \begin{bmatrix} 1 \\ 1 \end{bmatrix}$; not orthogonal to $\begin{bmatrix} v_1 \\ v_2 \end{bmatrix}$

power method converges to an eigen-vector

For all eigen-vectors have $(*)$

$$\frac{q_\infty}{p_\infty} = \frac{v_2}{v_1} = \frac{v_2}{\frac{1}{\sqrt{2}} v_2} = \sqrt{2}$$

Iterative methods for linear systems

Iterative solution of linear systems

Target problems: very **large** ($n = 10^5, 10^6, \dots$), A is usually **sparse** and has specific **properties**.

To solve

$$Ax = b$$

we construct a sequence

$$\mathbf{x}_1, \mathbf{x}_2, \dots$$

of iterates that converges fast to the solution \mathbf{x} , where \mathbf{x}_{k+1} can be cheaply computed from $\{\mathbf{x}_1, \dots, \mathbf{x}_k\}$ (e.g., one matrix-vector multiplication).

Thought experiment: If we can compute one iteration with cost $\mathcal{O}(n)$ (e.g., one matrix-vector multiplication with a sparse matrix) and need a constant $\mathcal{O}(1)$ number of iterations to reach desired precision, then we solve $Ax = b$ with costs $\mathcal{O}(n)$.

Intuitively, we cannot do better than that because we solve for n quantities and thus need to touch each at least once.

Let Q be invertible, then

$$\begin{aligned} \mathbf{Ax} = \mathbf{b} &\Leftrightarrow Q^{-1}(\mathbf{b} - \mathbf{Ax}) = 0 \\ &\Leftrightarrow (\mathbf{I} - Q^{-1}\mathbf{A})\mathbf{x} + Q^{-1}\mathbf{b} = \mathbf{x} \\ &\Leftrightarrow \mathbf{Gx} + \mathbf{c} = \mathbf{x} \end{aligned}$$

Leads to fixed-point iteration

$$\mathbf{x}_{k+1} = \mathbf{Gx}_k + \mathbf{c}$$

and with \mathbf{G} invertible obtain that $\mathbf{x} = \mathbf{A}^{-1}\mathbf{b}$ is a stationary point

Extreme cases for selecting Q

What are two extreme cases for selecting Q (need it to be invertible)?

Extreme cases for selecting Q

What are two extreme cases for selecting Q (need it to be invertible)?

Choose $Q = A^{-1}$, then our iteration becomes

$$\begin{aligned} Ax = \mathbf{b} &\Leftrightarrow A^{-1}(\mathbf{b} - Ax) = 0 \\ &\Leftrightarrow (I - A^{-1}A)\mathbf{x} + A^{-1}\mathbf{b} = \mathbf{x} \\ &\Leftrightarrow \mathbf{0} + \mathbf{x} = \mathbf{x} \end{aligned}$$

and we are done in just a single step

$$x_{k+1} = x$$

Thus, if we “know the solution” (in form of having the inverse A^{-1}) then no further work is needed here because we already did all the work when finding A^{-1}

The other extreme is setting $Q = I$, this leads to the Richardson method

$$\mathbf{x}_{k+1} = (\mathbf{I} - \mathbf{A})\mathbf{x}_k + \mathbf{b}$$

We have invested zero costs in finding Q and so we intuitively expect that $Q = I$ will require high costs in terms of number of iterations to converge in general, if it converges at all

Problem 3: Answer the following questions with 1-2 sentence explanations:

- ▶ **What is the problem with the Richardson method and what can we do about it?**

Problem 3: Answer the following questions with 1-2 sentence explanations:

- ▶ **What is the problem with the Richardson method and what can we do about it?**

Converges only for limited set of matrices. Use linear combination between new and previous iterate:

$$\mathbf{x}_{k+1} = \omega \underbrace{(G\mathbf{x}_k + c)}_{\mathbf{x}'_{k+1}} + (1 - \omega)\mathbf{x}_k = G_\omega\mathbf{x}_k + \omega c,$$

where $\omega > 0$ is a **damping/relaxation parameter**. Goal is to choose ω such that $\rho(G_\omega)$ is minimal.

- ▶ **What are other reasonable choices for Q ?**

Problem 3: Answer the following questions with 1-2 sentence explanations:

- ▶ **What is the problem with the Richardson method and what can we do about it?**

Converges only for limited set of matrices. Use linear combination between new and previous iterate:

$$\mathbf{x}_{k+1} = \omega \underbrace{(G\mathbf{x}_k + c)}_{\mathbf{x}'_{k+1}} + (1 - \omega)\mathbf{x}_k = G_\omega\mathbf{x}_k + \omega c,$$

where $\omega > 0$ is a **damping/relaxation parameter**. Goal is to choose ω such that $\rho(G_\omega)$ is minimal.

- ▶ **What are other reasonable choices for Q ?**

Jacobi uses $Q = D$ and Gauss-Seidel uses $Q = L + D$, which both can be computed quickly from A .

- ▶ **What property of A critically influences how quickly these relaxation methods converge?**

Problem 3: Answer the following questions with 1-2 sentence explanations:

- ▶ **What is the problem with the Richardson method and what can we do about it?**

Converges only for limited set of matrices. Use linear combination between new and previous iterate:

$$\mathbf{x}_{k+1} = \omega \underbrace{(G\mathbf{x}_k + c)}_{\mathbf{x}'_{k+1}} + (1 - \omega)\mathbf{x}_k = G_\omega\mathbf{x}_k + \omega c,$$

where $\omega > 0$ is a **damping/relaxation parameter**. Goal is to choose ω such that $\rho(G_\omega)$ is minimal.

- ▶ **What are other reasonable choices for Q ?**

Jacobi uses $Q = D$ and Gauss-Seidel uses $Q = L + D$, which both can be computed quickly from A .

- ▶ **What property of A critically influences how quickly these relaxation methods converge?** The spectral radius $\rho(A)$ of A

In the following A is symmetric positive definite.

Formulate solving $Ax = b$ as an optimization problem: Define

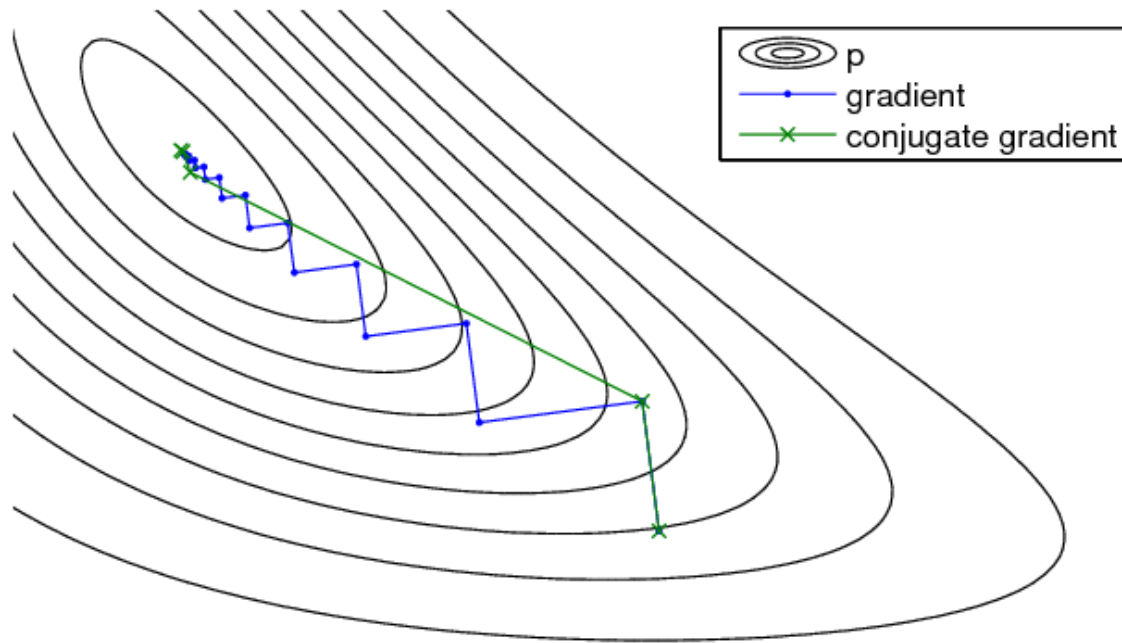
$$f(x) = \frac{1}{2}x^T Ax - b^T x,$$

and minimize

$$\min_{x \in \mathbb{R}^n} f(x)$$

Because A is positive definite, we have $f(x) = 0 \iff Ax = b$. It is sufficient to look at the gradient

$$\nabla f(x) = \frac{1}{2}A^T x + \frac{1}{2}Ax - b = Ax - b = -r(x) = 0 \iff Ax = b$$



[Figure: Kuusela et al., 2009]

The convergence behavior of steepest descent in this context can be poor: we eventually get arbitrarily close to the minimum but we can always destroy something of the already achieved when applying the update \rightsquigarrow can we find better search directions?

Conjugate gradient method

- ▶ All methods so far (relaxation, steepest descent) use information about x_{k-1} to get x_k . All information about earlier iterations is ignored.

Conjugate gradient method

- ▶ All methods so far (relaxation, steepest descent) use information about x_{k-1} to get x_k . All information about earlier iterations is ignored.
- ▶ The conjugate gradient (CG) method is a variation of steepest descent that *has a memory*.
- ▶ Let p_1, \dots, p_k be the directions up to step k , then CG uses the space

$$x_0 + \text{span}\{p_1, \dots, p_k\}, \quad x_0 \text{ starting point}$$

to find the next iterate x_k and thus

$$x_k = x_0 + \sum_{i=1}^k \alpha_i p_i$$

- ▶ (Recall that steepest descent uses only the search direction $p_k = r_{k-1} = -\nabla f(x_{k-1})$ to find the iterate x_k)

We want the following

- a The search directions p_1, \dots, p_k should be linearly independent ("we don't destroy what we have achieved")
- b We have ("we do the best we can at each step")

$$f(x_k) = \min_{x \in x_0 + \text{span}(p_1, \dots, p_k)} f(x)$$

- c The step x_k can be calculated easily from x_{k-1}

We then worked hard to derive the CG algorithm based on these three conditions

What was a critical step in CG?

We want the following

- a The search directions p_1, \dots, p_k should be linearly independent ("we don't destroy what we have achieved")
- b We have ("we do the best we can at each step")

$$f(x_k) = \min_{x \in x_0 + \text{span}(p_1, \dots, p_k)} f(x)$$

- c The step x_k can be calculated easily from x_{k-1}

We then worked hard to derive the CG algorithm based on these three conditions

What was a critical step in CG? \rightsquigarrow Gram-Schmidt orthogonalization to find the search direction

It can be shown that for $k \geq 1$ and $e_j \neq 0, j < k$ it holds

$$\|e_k\|_A \leq 2 \left(\frac{\sqrt{\kappa_2(A)} - 1}{\sqrt{\kappa_2(A)} + 1} \right)^k \|e_0\|_A$$

for spd matrices A . \rightsquigarrow Trefethen & Bau

For steepest descent, if A is spd, we obtained

$$\|x^* - x_k\|_A \leq \left(\frac{\kappa_2(A) - 1}{\kappa_2(A) + 1} \right)^k \|x^* - x_0\|_A,$$

where $\langle x, y \rangle_A = x^T A y$ and $\|\cdot\|_A = \sqrt{\langle \cdot, \cdot \rangle_A}$.

We keep finding we are limited by the condition number of A . What can be done about it (in Numerical Methods II)?

It can be shown that for $k \geq 1$ and $e_j \neq 0, j < k$ it holds

$$\|e_k\|_A \leq 2 \left(\frac{\sqrt{\kappa_2(A)} - 1}{\sqrt{\kappa_2(A)} + 1} \right)^k \|e_0\|_A$$

for spd matrices A . \rightsquigarrow Trefethen & Bau

For steepest descent, if A is spd, we obtained

$$\|x^* - x_k\|_A \leq \left(\frac{\kappa_2(A) - 1}{\kappa_2(A) + 1} \right)^k \|x^* - x_0\|_A,$$

where $\langle x, y \rangle_A = x^T A y$ and $\|\cdot\|_A = \sqrt{\langle \cdot, \cdot \rangle_A}$.

We keep finding we are limited by the condition number of A . What can be done about it (in Numerical Methods II)? Multigrid, multi-level

Solving systems of nonlinear equations

Solving nonlinear equations (“root finding”)

We want to solve the nonlinear equation

$$f(x) = 0, \quad x \in \mathbb{R}.$$

We could also have $n < \infty$ equations in n unknowns with $f : \mathbb{R}^n \rightarrow \mathbb{R}^n$

$$f(\mathbf{x}) = \mathbf{0}$$

In general, we will need an iterative approach that constructs x_1, x_2, x_3, \dots such that

$$\lim_{k \rightarrow \infty} x_k = x^*,$$

with $f(x^*) = 0$.

Reformulation as fixed point method so that x^* is fixed point

$$x^* = \Phi(x^*)$$

Let $F : \mathbb{R}^n \rightarrow \mathbb{R}^n$, $n \geq 1$ and solve

$$F(\mathbf{x}) = 0.$$

Truncated Taylor expansion of F about starting point \mathbf{x}^0 :

$$F(\mathbf{x}) \approx F(\mathbf{x}^0) + F'(\mathbf{x}^0)(\mathbf{x} - \mathbf{x}^0).$$

Hence:

$$\mathbf{x}^1 = \mathbf{x}^0 - F'(\mathbf{x}^0)^{-1}F(\mathbf{x}^0)$$

Newton iteration: Start with $\mathbf{x}^0 \in \mathbb{R}^n$, and for $k = 0, 1, \dots$ compute

$$F'(\mathbf{x}^k)\Delta\mathbf{x}^k = -F(\mathbf{x}^k), \quad \mathbf{x}^{k+1} = \mathbf{x}^k + \Delta\mathbf{x}^k$$

Requires that $F'(\mathbf{x}^k) \in \mathbb{R}^{n \times n}$ is invertible.

Problem 4: Answer the following questions with explanations of 1-2 sentences and/or proofs

- ▶ Describe two key requirements for Newton to converge quadratically.

Problem 4: Answer the following questions with explanations of 1-2 sentences and/or proofs

- ▶ Describe two key requirements for Newton to converge quadratically.

Need $F'(x)$ invertible and starting point x_0 close to solution x^*

- ▶ You have implemented Newton in Matlab and tried it on two starting points x_0 and y_0 for your problem. Newton converges for both but to a different solution. Does this mean there something wrong with your implementation?

Problem 4: Answer the following questions with explanations of 1-2 sentences and/or proofs

- ▶ Describe two key requirements for Newton to converge quadratically.

Need $F'(x)$ invertible and starting point x_0 close to solution x^*

- ▶ You have implemented Newton in Matlab and tried it on two starting points x_0 and y_0 for your problem. Newton converges for both but to a different solution. Does this mean there something wrong with your implementation?

No, can converge to local optima

- ▶ You want to solve $F(x) = 0$ but you have access only to a scrambled F through $(G_k \circ F) = G_k(F(x))$, where k is the Newton iteration, and thus you solve a different problem $G_k(F(x)) = 0$ in each iteration. Derive another condition on G_k than G_k being the identity so that Newton converges to x^* with $F(x^*) = 0$ in a neighborhood of x^* . Provide a proof.

$$G_{\eta} \circ F = A_{\eta} F \quad , \text{ regular } A_{\eta}$$

$$\gamma_{\eta+1} = \gamma_{\eta} - \left[(G_{\eta} \circ F)' \right]^{-1}(\gamma_{\eta}) (G_{\eta} \circ F)(\gamma_{\eta})$$

$$= \gamma_{\eta} - \underbrace{\left[(A_{\eta} F)' \right]^{-1}(\gamma_{\eta})}_{A_{\eta} F'} A_{\eta} F(\gamma_{\eta})$$

$$= \gamma_{\eta} - F'(\gamma_{\eta})^{-1} \underbrace{A_{\eta}^{-1} A_{\eta}} F(\gamma_{\eta})$$

$$= \gamma_{\eta} - F'(\gamma_{\eta})^{-1} F(\gamma_{\eta})$$

Problem 4: Answer the following questions with explanations of 1-2 sentences and/or proofs

- ▶ Describe two key requirements for Newton to converge quadratically.

Need $F'(x)$ invertible and starting point x_0 close to solution x^*

- ▶ You have implemented Newton in Matlab and tried it on two starting points x_0 and y_0 for your problem. Newton converges for both but to a different solution. Does this mean there something wrong with your implementation?

No, can converge to local optima

- ▶ You want to solve $F(x) = 0$ but you have access only to a scrambled F through $(G_k \circ F) = G_k(F(x))$, where k is the Newton iteration, and thus you solve a different problem $G_k(F(x)) = 0$ in each iteration. Derive another condition on G_k than G_k being the identity so that Newton converges to x^* with $F(x^*) = 0$ in a neighborhood of x^* . Provide a proof.

Newton is affine invariant. So as long as G_k is linear and the corresponding matrix is regular, it does not influence the Newton iterations. Proof is the same as for the affine invariance that we did in class.

Interpolation

Polynomial interpolation

Consider $n + 1$ pairs (x_i, y_i) , $i = 0, \dots, n$ of a function f with

$$y_i = f(x_i)$$

Let now \mathbb{P}_n be the set of all polynomials up to degree n over \mathbb{R} so that we have for all

$P \in \mathbb{P}_n$

$$P(x) = a_n x^n + a_{n-1} x^{n-1} + \dots + a_1 x + a_0, \quad a_n, \dots, a_0 \in \mathbb{R}$$

We would like to find a $P \in \mathbb{P}_n$ such that

$$P(x_i) = y_i, \quad i = 0, \dots, n$$

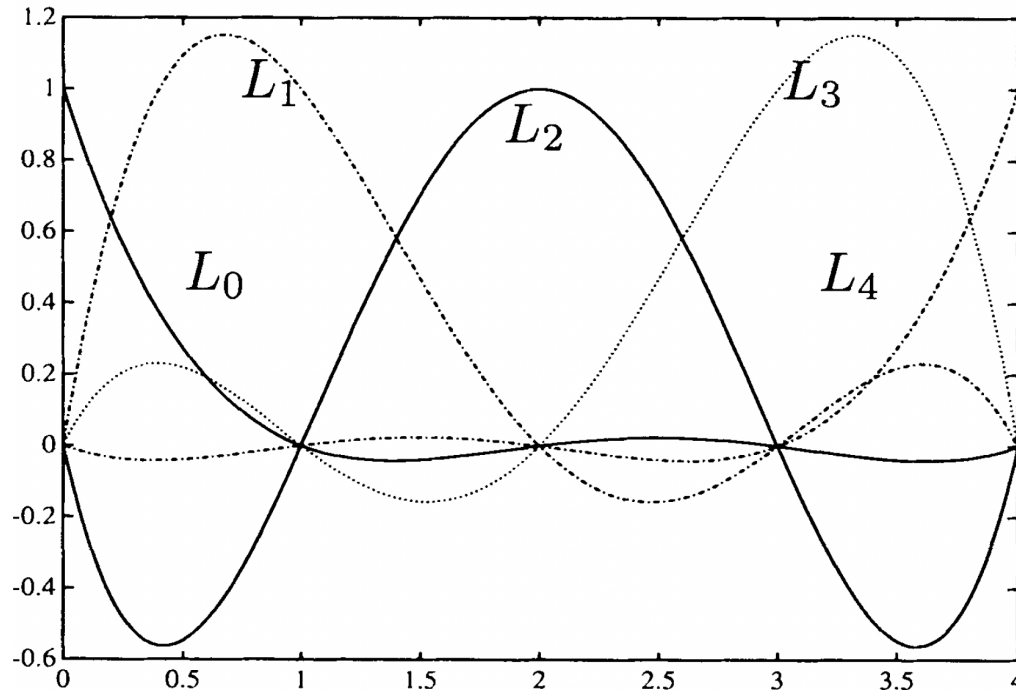
Theorem: Given $n + 1$ points (x_i, y_i) with pairwise distinct x_0, \dots, x_n , there exists a unique polynomial $P \in \mathbb{P}_n$ such that

$$P(x_i) = y_i, \quad i = 0, \dots, n$$

Lagrange basis

The Lagrange polynomials $L_0, \dots, L_n \in \mathbb{P}_n$ are uniquely defined for distinct x_0, \dots, x_n

$$L_i(x_j) = \delta_{ij}, \quad L_i \in \mathbb{P}_n.$$



Lagrange polynomials up to order $n = 4$ for equidistant x_0, \dots, x_4 . [Figure: Deuffhard]

The corresponding explicit formula is

$$L_i(x) = \prod_{\substack{j=0 \\ j \neq i}}^n \frac{x - x_j}{x_i - x_j}, \quad i = 0, \dots, n$$

What are the coefficients a_n, \dots, a_0 so that $P(x_i) = y_i$ for $i = 0, \dots, n$?

The corresponding explicit formula is

$$L_i(x) = \prod_{\substack{j=0 \\ j \neq i}}^n \frac{x - x_j}{x_i - x_j}, \quad i = 0, \dots, n$$

What are the coefficients a_n, \dots, a_0 so that $P(x_i) = y_i$ for $i = 0, \dots, n$?

$$P(x) = \sum_{i=0}^n y_i L_i(x)$$

because

$$P(x_j) = \sum_{i=0}^n y_i L_i(x_j) = \sum_{i=0}^n y_i \delta_{ij} = y_j$$

If we have the basis L_0, \dots, L_n , we obtain the polynomial P for free

The corresponding explicit formula is

$$L_i(x) = \prod_{\substack{j=0 \\ j \neq i}}^n \frac{x - x_j}{x_i - x_j}, \quad i = 0, \dots, n$$

What are the coefficients a_n, \dots, a_0 so that $P(x_i) = y_i$ for $i = 0, \dots, n$?

$$P(x) = \sum_{i=0}^n y_i L_i(x)$$

because

$$P(x_j) = \sum_{i=0}^n y_i L_i(x_j) = \sum_{i=0}^n y_i \delta_{ij} = y_j$$

If we have the basis L_0, \dots, L_n , we obtain the polynomial P for free but the cost of evaluating the polynomial is too high for practical computations

Polynomial interpolation in Newton basis

The **Newton basis** $\omega_0, \dots, \omega_n$ is given by

$$\omega_i(x) := \prod_{j=0}^{i-1} (x - x_j) \in \mathbb{P}_i.$$

Polynomial interpolation in Newton basis

The **Newton basis** $\omega_0, \dots, \omega_n$ is given by

$$\omega_i(x) := \prod_{j=0}^{i-1} (x - x_j) \in \mathbb{P}_i.$$

Would like to find coefficients c_0, c_1, \dots, c_n of interpolating polynomial in Newton basis

$$P_f(x|x_0, \dots, x_n) = c_0\omega_0(x) + c_1\omega_1(x) + \dots + c_n\omega_n(x)$$

Polynomial interpolation in Newton basis

The **Newton basis** $\omega_0, \dots, \omega_n$ is given by

$$\omega_i(x) := \prod_{j=0}^{i-1} (x - x_j) \in \mathbb{P}_i.$$

Would like to find coefficients c_0, c_1, \dots, c_n of interpolating polynomial in Newton basis

$$P_f(x|x_0, \dots, x_n) = c_0\omega_0(x) + c_1\omega_1(x) + \dots + c_n\omega_n(x)$$

The leading coefficient a_n of the interpolation polynomial

$$P_f(x|x_0, \dots, x_n) = a_n x^n + \dots + a_0$$

is called the *n-th divided difference*, $[x_0, \dots, x_n]f := a_n$.

The divided differences are the coefficients c_0, \dots, c_n : The interpolation polynomial $P_f(\cdot | x_0, \dots, x_n)$ for $x_0 \leq x_1 \leq \dots \leq x_n$ is given by

$$P(x) = \sum_{i=0}^n [x_0, \dots, x_i] f \omega_i(x).$$

The following recurrence relation holds for $x_i \neq x_j$:

$$[x_0, \dots, x_n] f = \frac{([x_0, \dots, \hat{x}_i, \dots, x_n] f - [x_0, \dots, \hat{x}_j, \dots, x_n] f)}{x_j - x_i}$$

which is helpful to compute the divided differences

Problem 5: Compute the polynomial p with $\underline{p(0) = 2}, \underline{p(1) = 3}, \underline{p(3) = 0}$ in the Newton basis

x_i	
0	$[x_0] p = 2 \quad P_1 \longrightarrow$
1	$[x_1] p = 3 \quad [x_0, x_1] p = \frac{[x_0] p - [x_1] p}{x_0 - x_1} = \frac{2 - 3}{0 - 1} = 1 \quad P_2$
3	$[x_0] p = 0$

\vdots
 \vdots
 \vdots

$$P_3(x) = 2 + x - \frac{5}{6} x(x-1)$$

Newton-Cotes formulas for quadrature

Given **fixed** nodes t_0, \dots, t_n , use polynomial approximation

$$\hat{f} = P_f(t|t_0, \dots, t_n) = \sum_{i=0}^n f(t_i) L_{in}(t)$$

with Lagrange polynomials L_{0n}, \dots, L_{nn}

Thus:

$$\hat{I}(f) = (b - a) \sum_{i=0}^n \lambda_{in} f(t_i),$$

where $\lambda_{in} = \frac{1}{b-a} \int_a^b L_{in}(t) dt \rightsquigarrow$ weights are unique

Quadrature formulas defined in this way are exact for polynomials $P \in \mathbb{P}_n$ of degree less than or equal to n

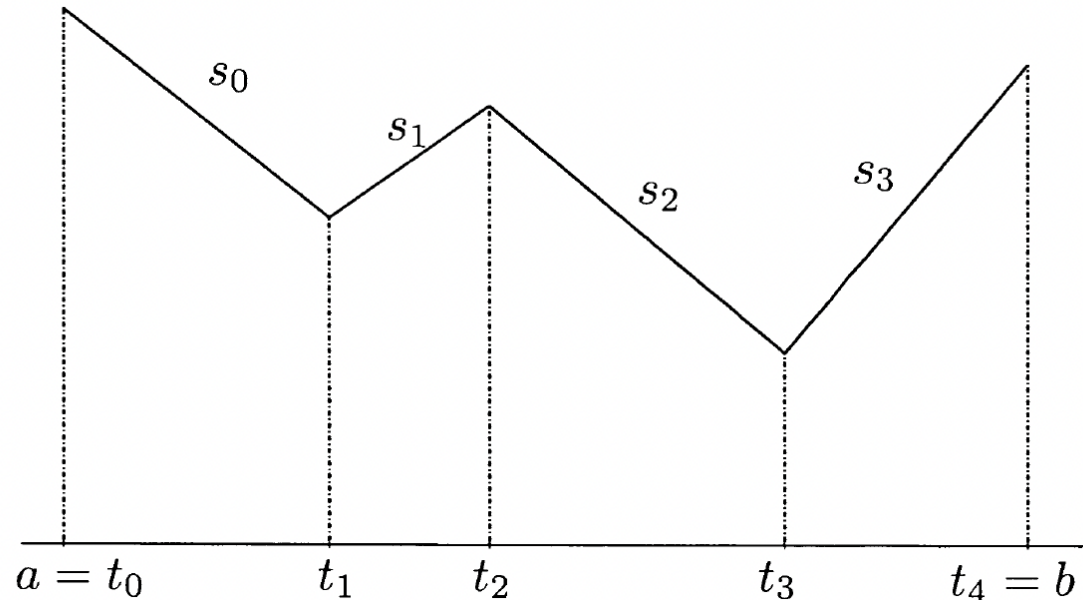
$$\hat{I}(P) = I(P_n(P)) = I(P), \quad \text{for all } P \in \mathbb{P}_n$$

Trapezoidal sums

To avoid poorly conditioned problems, let us split the integration interval $[a, b]$ into n sub-intervals $[t_{i-1}, t_i]$, $i = 1, \dots, n$. Then consider the rule

$$\hat{I}(f) = \sum_{i=1}^n \hat{I}_{t_{i-1}}^{t_i}(f),$$

where $\hat{I}_{t_{i-1}}^{t_i}$ is a quadrature formula on the interval $[t_{i-1}, t_i]$.



We have seen already the trapezoidal sum with $h = (b - a)/n$

$$T(h) = \sum_{i=1}^n T_i = h \left(\frac{1}{2}(f(a) + f(b)) + \sum_{i=1}^{n-1} f(a + ih) \right)$$

that has error

$$T(h) - \int_a^b f = \frac{(b-a)h^2}{12} f''(\tau), \quad \tau \in [a, b]$$

\rightsquigarrow we can increase n (and thus decrease h) to reduce the error without increasing the degree of the underlying polynomial

Conclusions

- ▶ This was a very selective review of the topics that we covered!
- ▶ Be prepared to answer True/False questions to show your basic understand of topics that we discussed.
- ▶ Be prepared to answer questions that require 1-2 sentence explanations or short proofs.
- ▶ You should be able to do quick calculations such as, e.g., approximating an integral with the trapezoidal rule and computing the LU decomposition of small matrices and computing eigenvalues of a 2×2 matrix etc.
- ▶ Recapping basic linear algebra topics will be beneficial (e.g., computing determinant of a 2×2 matrix)