

# Numerical Analysis / Scientific Computing

## CS450

Andreas Kloeckner

Fall 2023

# Outline

## Introduction to Scientific Computing

Notes

Notes (unfilled, with empty boxes)

Notes (source code on Github)

About the Class

Errors, Conditioning, Accuracy, Stability

Floating Point

## Systems of Linear Equations

## Linear Least Squares

## Eigenvalue Problems

## Nonlinear Equations

## Optimization

## Interpolation

## Numerical Integration and Differentiation

## Initial Value Problems for ODEs

## Boundary Value Problems for ODEs

## Partial Differential Equations and Sparse Linear Algebra

## Fast Fourier Transform

## Additional Topics

## What's the point of this class?

'*Scientific Computing*' describes a family of approaches to obtain approximate solutions to problems *once they've been stated mathematically*.

Name some applications:

- ▶ Engineering simulation
  - ▶ E.g. Drag from flow over airplane wings, behavior of photonic devices, radar scattering, ...
  - ▶ → Differential equations (ordinary and partial)
- ▶ Machine learning
  - ▶ Statistical models, with unknown parameters
  - ▶ → Optimization
- ▶ Image and Audio processing
  - ▶ Enlargement/Filtering
  - ▶ → Interpolation
- ▶ Lots more.

# What do we study, and how?

Problems with real numbers (i.e. *continuous* problems)

- ▶ As opposed to *discrete* problems.
- ▶ Including: How can we put a real number into a computer? (and with what restrictions?)

What's the general approach?

- ▶ Pick a *representation* (e.g.: a polynomial)
- ▶ Existence/uniqueness?

## What makes for *good* numerics?

How good of an answer can we expect to our problem?

- ▶ Can't even represent numbers exactly.
- ▶ Answers will always be *approximate*.
- ▶ So, it's natural to ask *how far off the mark* we really are.

*How fast* can we expect the computation to complete?

- ▶ A.k.a. what algorithms do we use?
- ▶ What is the cost of those algorithms?
- ▶ Are they efficient?  
(I.e. do they make good use of available machine time?)

## Implementation concerns

How do numerical methods *get implemented*?

- ▶ Like anything in computing: A layer cake of *abstractions* (“careful lies”)
- ▶ What tools/languages are available?
- ▶ Are the methods easy to implement?
- ▶ If not, how do we make use of existing tools?
- ▶ How robust is our implementation? (e.g. for error cases)

# Class web page

<https://bit.ly/cs450-f23>

- ▶ Assignments
  - ▶ HW1 (soon!)
  - ▶ Pre-lecture quizzes
  - ▶ In-lecture interactive content (bring computer or phone if possible)
- ▶ Textbook
- ▶ Exams
- ▶ Class outline (with links to notes/demos/activities/quizzes)
- ▶ Discussion forum
- ▶ Policies
- ▶ Video

## Programming Language: Python/numpy

- ▶ Reasonably readable
- ▶ Reasonably beginner-friendly
- ▶ Mainstream (top 5 in 'TIOBE Index')
- ▶ Free, open-source
- ▶ Great tools and libraries (not just) for scientific computing
- ▶ Python 2/3? 3!
- ▶ `numpy`: Provides an array datatype  
Will use this and `matplotlib` all the time.
- ▶ See class web page for learning materials

**Demo:** Sum the squares of the integers from 0 to 100. First without `numpy`, then with `numpy`.



## Supplementary Material

- ▶ [Numpy \(from the SciPy Lectures\)](#)
- ▶ [100 Numpy Exercises](#)
- ▶ [Dive into Python3](#)

## Sources for these Notes

- ▶ M.T. Heath, *Scientific Computing: An Introductory Survey*, Revised Second Edition. Society for Industrial and Applied Mathematics, Philadelphia, PA. 2018.
- ▶ [CS 450 Notes by Edgar Solomonik](#)
- ▶ Various bits of prior material by Luke Olson

## Open Source <3

These notes (and the accompanying demos) are open-source!

Bug reports and pull requests welcome:

<https://github.com/inducer/numerics-notes>

Copyright (C) 2020 Andreas Kloeckner

Permission is hereby granted, free of charge, to any person obtaining a copy of this software and associated documentation files (the "Software"), to deal in the Software without restriction, including without limitation the rights to use, copy, modify, merge, publish, distribute, sublicense, and/or sell copies of the Software, and to permit persons to whom the Software is furnished to do so, subject to the following conditions:

The above copyright notice and this permission notice shall be included in all copies or substantial portions of the Software.

THE SOFTWARE IS PROVIDED "AS IS", WITHOUT WARRANTY OF ANY KIND, EXPRESS OR IMPLIED, INCLUDING BUT NOT LIMITED TO THE WARRANTIES OF MERCHANTABILITY, FITNESS FOR A PARTICULAR PURPOSE AND NONINFRINGEMENT. IN NO EVENT SHALL THE AUTHORS OR COPYRIGHT HOLDERS BE LIABLE FOR ANY CLAIM, DAMAGES OR OTHER LIABILITY, WHETHER IN AN ACTION OF CONTRACT, TORT OR OTHERWISE, ARISING FROM, OUT OF OR IN CONNECTION WITH THE SOFTWARE OR THE USE OR OTHER DEALINGS IN THE SOFTWARE.

## What problems *can* we study in the first place?

To be able to compute a solution (through a process that introduces errors), the problem...

- ▶ Needs to *have* a solution
- ▶ That solution should be *unique*
- ▶ And *depend continuously* on the inputs

If it satisfies these criteria, the problem is called *well-posed*. Otherwise, *ill-posed*.

## Dependency on Inputs

We excluded discontinuous problems—because we don't stand much chance for those.

... what if the problem's input dependency is just *close to discontinuous*?

- ▶ We call those problems *sensitive* to their input data. Such problems are obviously trickier to deal with than non-sensitive ones.
- ▶ Ideally, the computational method will not *amplify* the sensitivity

# Approximation

*When* does approximation happen?

- ▶ Before computation
  - ▶ modeling
  - ▶ measurements of input data
  - ▶ computation of input data
- ▶ During computation
  - ▶ truncation / discretization
  - ▶ rounding

Demo: Truncation vs Rounding [cleared]

## Example: Surface Area of the Earth

Compute the surface area of the earth.

What parts of your computation are approximate?

*All of them.*

$$A = 4\pi r^2$$

- ▶ Earth isn't really a sphere
- ▶ What does radius mean if the earth isn't a sphere?
- ▶ How do you compute with  $\pi$ ? (By rounding/truncating.)

# Measuring Error

How do we measure error?

**Idea:** Consider all error as being *added onto* the result.

$$\text{Absolute error} = \text{approx value} - \text{true value}$$

$$\text{Relative error} = \frac{\text{Absolute error}}{\text{True value}}$$

**Problem:** True value not known

- ▶ Estimate
- ▶ 'How big at worst?' → Establish *Upper Bounds*



## Recap: Norms

What's a norm?

- ▶  $f(\mathbf{x}) : \mathbb{R}^n \rightarrow \mathbb{R}_0^+$ , returns a 'magnitude' of the input vector
- ▶ In symbols: Often written  $\|\mathbf{x}\|$ .

Define *norm*.

A function  $\|\mathbf{x}\| : \mathbb{R}^n \rightarrow \mathbb{R}_0^+$  is called a norm if and only if

1.  $\|\mathbf{x}\| > 0 \Leftrightarrow \mathbf{x} \neq 0$ .
2.  $\|\gamma\mathbf{x}\| = |\gamma| \|\mathbf{x}\|$  for all scalars  $\gamma$ .
3. Obeys triangle inequality  $\|\mathbf{x} + \mathbf{y}\| \leq \|\mathbf{x}\| + \|\mathbf{y}\|$

## Norms: Examples

Examples of norms?

The so-called *p-norms*:

$$\left\| \begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix} \right\|_p = \sqrt[p]{|x_1|^p + \cdots + |x_n|^p} \quad (p \geq 1)$$

$p = 1, 2, \infty$  particularly important

[Demo: Vector Norms \[cleared\]](#)

## Norms: Which one?

Does the choice of norm really matter much?

In finitely many dimensions, all *norms are equivalent*.

I.e. for fixed  $n$  and two norms  $\|\cdot\|$ ,  $\|\cdot\|^*$ , there exist  $\alpha, \beta > 0$  so that for all vectors  $\mathbf{x} \in \mathbb{R}^n$

$$\alpha \|\mathbf{x}\| \leq \|\mathbf{x}\|^* \leq \beta \|\mathbf{x}\|.$$

So: No, doesn't matter *that much*. Will start mattering more for so-called *matrix norms*—see later.

In these notes: If we write  $\|\cdot\|$  without any specifics, then the statement is true for any norm. If a specific norm is needed, the notation will indicate that.

## Norms and Errors

If we're computing a vector result, the error is a vector.  
That's not a very useful answer to 'how big is the error'.  
What can we do?

Apply a norm!

How? *Attempt 1:*

**Magnitude of error**  $\neq$   $\|\text{true value}\| - \|\text{approximate value}\|$

**WRONG!** (How does it fail?)

*Attempt 2:*

**Magnitude of error**  $= \|\text{true value} - \text{approximate value}\|$

## Forward/Backward Error

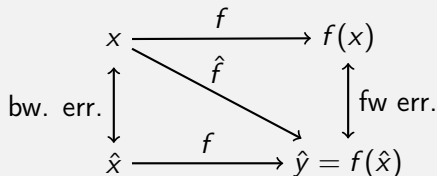
Suppose *want* to compute  $y = f(x)$ , but *approximate*  $\hat{y} = \hat{f}(x)$ .

What are the forward error and the backward error?

*Forward error:*  $\Delta y = \hat{y} - y$

*Backward error:* Imagine *all* error came from feeding the wrong input into a fully accurate calculation. Backward error is the difference between true and 'wrong' input. I.e.

- ▶ Find the  $\hat{x}$  closest to  $x$  so that  $f(\hat{x}) = \hat{y}$ .
- ▶  $\Delta x = \hat{x} - x$ .



## Forward/Backward Error: Example

Suppose you wanted  $y = \sqrt{2}$  and got  $\hat{y} = 1.4$ .  
What's the (magnitude of) the forward error?

$$|\Delta y| = |1.4 - 1.41421\dots| \approx 0.0142\dots$$

Relative forward error:

$$\frac{|\Delta y|}{|y|} = \frac{0.0142\dots}{1.41421\dots} \approx 0.01.$$

About 1 percent, or *two accurate significant digits*.

## Forward/Backward Error: Example

Suppose you wanted  $y = \sqrt{2}$  and got  $\hat{y} = 1.4$ .  
What's the (magnitude of) the backward error?

Need  $\hat{x}$  so that  $f(\hat{x}) = 1.4$ .

$$\sqrt{1.96} = 1.4, \quad \Rightarrow \quad \hat{x} = 1.96.$$

Backward error:

$$|\Delta x| = |1.96 - 2| = 0.04.$$

Relative backward error:

$$\frac{|\Delta x|}{|x|} \approx 0.02.$$

About 2 percent.

## Forward/Backward Error: Observations

What do you observe about the relative magnitude of the relative errors?

- ▶ In this case: Got smaller, i.e. variation *damped out*.
- ▶ Typically: Not that lucky: Input error *amplified*.
- ▶ If backward error is smaller than the input error: result “as good as possible”.

This amplification factor seems worth studying in more detail.



## Sensitivity and Conditioning

Consider a more general setting: An input  $x$  and its perturbation  $\hat{x}$ .

Want: the smallest number  $\kappa_{\text{rel}}$  such that

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

(rel. perturbation in output)  $\leq \kappa_{\text{rel}} \cdot$  (rel. perturbation in input)

Call this the (*relative*) *condition number*. Find it via:

$$\kappa_{\text{rel}} = \max_{x, \hat{x}} \frac{|f(x) - f(\hat{x})| / |f(x)|}{|x - \hat{x}| / |x|}.$$

- ▶ Technically: should use ‘supremum’.
- ▶ Must specify set of  $x, \hat{x}$  that are “of interest”.

## Absolute Condition Number

Can you also define an *absolute* condition number?

Certainly:

$$\kappa_{\text{abs}} = \max_{x, \hat{x}} \frac{|f(x) - f(\hat{x})|}{|x - \hat{x}|}$$

But: less commonly used than relative, because we *typically* care about relative error.

When not specified: Assume condition number means *relative*.

## Interpreting a Condition Number

What does it mean for condition numbers to be small/large?

If the condition number is...

- ▶ ...small: the problem *well-conditioned* or insensitive
- ▶ ...large: the problem *ill-conditioned* or sensitive

*Can* also talk about condition number for a single input  $x$ .

Relate the (relative) condition number back to the setting of (relative) backward error.

Realize  $\hat{x}$  in backward error is just a perturbation of  $x$ : Same setting as conditioning.

Therefore:

$$|\text{rel. fwd. err.}| \leq \kappa_{\text{rel}} \cdot |\text{rel. bwd. err.}|$$

## Example: Condition Number of Evaluating a Function

$y = f(x)$ . Assume  $f$  differentiable.

$$\kappa = \max_x \frac{|\Delta y| / |y|}{|\Delta x| / |x|}$$

Forward error:

$$\Delta y = f(x + \Delta x) - f(x) \approx f'(x)\Delta x$$

Condition number:

$$\kappa \geq \frac{|\Delta y| / |y|}{|\Delta x| / |x|} \approx \frac{|f'(x)| |\Delta x| / |f(x)|}{|\Delta x| / |x|} = \frac{|xf'(x)|}{|f(x)|}.$$

Demo: Conditioning of Evaluating tan [cleared]

## Stability and Accuracy

**Previously:** Considered *problems* or *questions*.

**Next:** Considered *methods*, i.e. computational approaches to find solutions.

When is a method *accurate*?

Closeness of method output to true answer for unperturbed input.

When is a method *stable*?

- ▶ “A method is stable if the result it produces is the exact answer for a nearby input.”
- ▶ The above is commonly called *backward stability* and is a stricter requirement than just the temptingly simple:  
If the method’s sensitivity to variation in the input is no (or not much) greater than that of the problem itself.

## Getting into Trouble with Accuracy and Stability

How can I produce inaccurate results?

- ▶ Apply an inaccurate method
- ▶ Apply an unstable method to a well-conditioned problem
- ▶ Apply any type of method to an ill-conditioned problem

## In-Class Activity: Forward/Backward Error

In-class activity: Forward/Backward Error

## Wanted: Real Numbers... in a computer

Computers can represent *integers*, using bits:

$$23 = 1 \cdot 2^4 + 0 \cdot 2^3 + 1 \cdot 2^2 + 1 \cdot 2^1 + 1 \cdot 2^0 = (10111)_2$$

How would we represent fractions?

**Idea:** Keep going down past zero exponent:

$$23.625 = 1 \cdot 2^4 + 0 \cdot 2^3 + 1 \cdot 2^2 + 1 \cdot 2^1 + 1 \cdot 2^0 \\ + 1 \cdot 2^{-1} + 0 \cdot 2^{-2} + 1 \cdot 2^{-3}$$

**So:** Could store

- ▶ a fixed number of bits with exponents  $\geq 0$
- ▶ a fixed number of bits with exponents  $< 0$

This is called *fixed-point arithmetic*.



## Fixed-Point Numbers

Suppose we use units of 64 bits, with 32 bits for exponents  $\geq 0$  and 32 bits for exponents  $< 0$ . What numbers can we represent?

$2^{31}$	$\dots$	$2^0$	$2^{-1}$	$\dots$	$2^{-32}$
----------	---------	-------	----------	---------	-----------

**Smallest:**  $2^{-32} \approx 10^{-10}$

**Largest:**  $2^{31} + \dots + 2^{-32} \approx 10^9$

How many 'digits' of relative accuracy (think relative rounding error) are available for the smallest vs. the largest number?

**For large numbers:** about 19

**For small numbers:** few or none

**Idea:** Instead of *fixing* the location of the 0 exponent, let it *float*.

## Floating Point Numbers

Convert  $13 = (1101)_2$  into floating point representation.

$$13 = 2^3 + 2^2 + 2^0 = (1.101)_2 \cdot 2^3$$

What pieces do you need to store an FP number?

*Significand:*  $(1.101)_2$

- ▶ Generally:  $0 \leq \text{significand} < 2$
- ▶ Can actually require:  $1 \leq \text{significand} < 2$  (**normalization**)

*Exponent:* 3

## Floating Point: Implementation, Normalization

**Previously:** Consider *mathematical* view of FP. (via example:  $(1101)_2$ )

**Next:** Consider *implementation* of FP in hardware.

Do you notice a source of inefficiency in our number representation?

**Idea:** Notice that the leading digit (in binary) of the significand is always one.

Only store '101'. Final storage format:

*Significand:* 101 – a fixed number of bits

*Exponent:* 3 – a (*signed!*) integer allowing a certain range

Exponent is most often stored as a positive 'offset' from a certain negative number. E.g.

$$3 = \underbrace{-1023}_{\text{implicit offset}} + \underbrace{1026}_{\text{stored}}$$

Actually stored: 1026, a positive integer.

## Unrepresentable numbers?

Can you think of a somewhat central number that we cannot represent as

$$x = (1.\text{-----})_2 \cdot 2^{-p}?$$

Zero. Which is somewhat embarrassing.

**Core problem:** The implicit 1. It's a great idea, were it not for this issue.

Have to break the pattern. **Idea:**

- ▶ Declare one exponent 'special', and turn off the leading one for that one.  
(say,  $-1023$ , a.k.a. stored exponent 0)
- ▶ For all larger exponents, the leading one remains in effect.

**Bonus Q:** With this convention, what is the binary representation of a zero?

## Subnormal Numbers

What is the smallest representable number in an FP system with 4 stored bits (5 total) in the significand and a stored exponent range of  $[-7, 8]$ ?

First attempt:

- ▶ Significand as small as possible  $\rightarrow$  all zeros after the implicit leading one
- ▶ Exponent as small as possible:  $-7$

So:

$$(1.0000)_2 \cdot 2^{-7}.$$

Unfortunately: **wrong**.

## Subnormal Numbers, Attempt 2

What is the smallest representable number in an FP system with 4 stored bits in the significand and a (stored) exponent range of  $[-7, 8]$ ?

- ▶ Can go way smaller using the *special exponent* (turns off the leading one)
- ▶ Assume that the special exponent is  $-7$ ; interpreted as  $-6$ .
- ▶ So:  $(0.0001)_2 \cdot 2^{-6}$  (with four digits after the point stored).

Numbers with the special exponent are called *subnormal* (or *denormal*) FP numbers. Technically, zero is also a subnormal.

Why learn about subnormals?

- ▶ Subnormal FP is often slow: not implemented in hardware.
- ▶ Many compilers support options to 'flush subnormals to zero'.

## Underflow

- ▶ FP systems without subnormals will *underflow* (return 0) as soon as the exponent range is exhausted.
- ▶ This smallest representable *normal* number is called the *underflow level*, or *UFL*.
- ▶ Beyond the underflow level, subnormals provide for *gradual underflow* by ‘keeping going’ as long as there are bits in the significand, but it is important to note that subnormals don’t have as many accurate digits as normal numbers.  
[Read a story on the epic battle about gradual underflow](#)
- ▶ Analogously (but much more simply—no ‘supernormals’): the overflow level, *OFL*.

## Rounding Modes

How is rounding performed? (Imagine trying to represent  $\pi$ .)

$$\left( \underbrace{1.1101010}_{\text{representable}} 11 \right)_2$$

- ▶ “Chop” a.k.a. *round-to-zero*:  $(1.1101010)_2$
- ▶ *Round-to-nearest*:  $(1.1101011)_2$  (most accurate)

What is done in case of a tie?  $0.5 = (0.1)_2$  (“Nearest”?)

Up or down? It turns out that picking the same direction every time introduces *bias*. Trick: *round-to-even*.

$$0.5 \rightarrow 0, \quad 1.5 \rightarrow 2$$

[Demo: Density of Floating Point Numbers](#) [cleared]

[Demo: Floating Point vs Program Logic](#) [cleared]



## Smallest Numbers Above...

- ▶ What is smallest FP number  $> 1$ ? Assume 4 stored bits (5 total) in the significand.

$$(1.0001)_2 \cdot 2^0 = x \cdot (1 + 0.0001)_2$$

What's the smallest FP number  $> 1024$  in that same system?

$$(1.0001)_2 \cdot 2^{10} = x \cdot (1 + 0.0001)_2$$

Can we give that number a name?

## Unit Roundoff

*Unit roundoff or machine precision or machine epsilon or  $\epsilon_{\text{mach}}$  is...*

the smallest number such that  $\text{float}(1 + \epsilon) > 1$ .

- ▶ **Technically** that makes  $\epsilon_{\text{mach}}$  depend on the rounding rule.
- ▶ For rules with tie-breaking (e.g. round-to-nearest): assume bias towards 'different'.
- ▶ For example: Assuming round-to-nearest, in a system with five bits in the significand,  $\epsilon_{\text{mach}} = (0.00001)_2$ .
- ▶ Another, related, quantity is *ULP*, or *unit in the last place*.  
For round-to-nearest: ( $\epsilon_{\text{mach}} = 0.5 \text{ ULP}$ )

## FP: Relative Rounding Error

What does this say about the relative error incurred in floating point calculations?

- ▶ The factor to get from one FP number to the next larger one is (mostly) independent of magnitude:  $1 + \epsilon_{\text{mach}}$ .
- ▶ Since we can't represent any results between  $x$  and  $x \cdot (1 + \epsilon_{\text{mach}})$ , that's really the minimum error incurred.
- ▶ In terms of relative error:

$$\left| \frac{\tilde{x} - x}{x} \right| = \left| \frac{x(1 + \epsilon_{\text{mach}}) - x}{x} \right| = \epsilon_{\text{mach}}.$$

At least theoretically,  $\epsilon_{\text{mach}}$  is the maximum relative error in any FP operations. (Practical implementations do fall short of this.)

## FP: Machine Epsilon

What's machine epsilon for double-precision floating point with round-to-nearest? (52 stored bits in the significand, 53 total)

$$2^{-53} \approx 10^{-16}$$

**Bonus Q:** What does  $1 + 2^{-53}$  do on your computer? Why?

We can expect FP math to consistently introduce little relative errors of about  $10^{-16}$ .

Working in double precision gives you about 16 (decimal) accurate digits.

[Demo: Floating point and the Harmonic Series](#) [cleared]

## Implementing Arithmetic

How is floating point addition implemented?

Consider adding  $a = (1.101)_2 \cdot 2^1$  and  $b = (1.001)_2 \cdot 2^{-1}$  in a system with three stored bits (four total) in the significand.

Rough algorithm:

1. Bring both numbers onto a common exponent
2. Do grade-school addition from the front, until you run out of digits in your system.
3. Round result.

$$\begin{aligned}a &= 1. \text{ 101} \cdot 2^1 \\b &= 0. \text{ 01001} \cdot 2^1 \\a + b &\approx 1. \text{ 111} \cdot 2^1\end{aligned}$$

## Problems with FP Addition

What happens if you subtract two numbers of very similar magnitude?  
As an example, consider  $a = (1.1011)_2 \cdot 2^0$  and  $b = (1.1010)_2 \cdot 2^0$ .

$$\begin{aligned}a &= 1.1011 \cdot 2^1 \\b &= 1.1010 \cdot 2^1 \\a - b &\approx 0.0001???? \cdot 2^1\end{aligned}$$

or, once we normalize,

$$1.???? \cdot 2^{-3}.$$

There is no data to indicate what the missing digits should be.

→ Machine fills them with zero.

This phenomenon is called *Catastrophic Cancellation*.

[Demo: Catastrophic Cancellation \[cleared\]](#)

## In-Class Activity: Floating Point

In-class activity: Floating Point

## Supplementary Material

- ▶ Josh Haberman, [Floating Point Demystified, Part 1](#)
- ▶ David Goldberg, [What every computer programmer should know about floating point](#)
- ▶ Evan Wallace, [Float Toy](#)
- ▶ Julia Evans, [Examples of Floating Point Problems](#), 2022



# Outline

Introduction to Scientific Computing

## Systems of Linear Equations

Theory: Conditioning

Methods to Solve Systems

LU: Application and Implementation

Linear Least Squares

Eigenvalue Problems

Nonlinear Equations

Optimization

Interpolation

Numerical Integration and Differentiation

Initial Value Problems for ODEs

Boundary Value Problems for ODEs

Partial Differential Equations and Sparse Linear Algebra

Fast Fourier Transform

Additional Topics

## Solving a Linear System

Given:

- ▶  $m \times n$  matrix  $A$
- ▶  $m$ -vector  $\mathbf{b}$

What are we looking for here, and when are we allowed to ask the question?

**Want:**  $n$ -vector  $\mathbf{x}$  so that  $A\mathbf{x} = \mathbf{b}$ .

- ▶ Linear combination of columns of  $A$  to yield  $\mathbf{b}$ .
- ▶ **Restrict** to square case ( $m = n$ ) for now.
- ▶ Even with that: solution may not exist, or may not be unique.

Unique solution exists iff  $A$  is *nonsingular*.

**Next:** Want to talk about conditioning of this operation. Need to measure distances of matrices.

## Matrix Norms

What norms would we apply to matrices?

**Could use:** “*Flatten*” matrix as vector, use vector norm.

**But we won't:** Not very meaningful.

**Instead:** Choose norms for matrices to interact with an ‘associated’ vector norm  $\|\cdot\|$  so that  $\|A\|$  obeys

$$\|A\mathbf{x}\| \leq \|A\| \|\mathbf{x}\|.$$

For a given vector norm, define **induced matrix norm**  $\|\cdot\|$ ,

$$\|A\| := \max_{\|\mathbf{x}\|=1} \|A\mathbf{x}\|.$$

For each vector norm, we get a different matrix norm, e.g. for the vector 2-norm  $\|\mathbf{x}\|_2$  we get a matrix 2-norm  $\|A\|_2$ .

## Intuition for Matrix Norms

Provide some intuition for the matrix norm.

$$\max_{x \neq 0} \frac{\|Ax\|}{\|x\|} = \max_{x \neq 0} \left\| Ax \cdot \frac{1}{\|x\|} \right\| = \max_{\|y\|=1} \|Ay\| = \|A\|.$$

I.e. the matrix norm gives the maximum (relative) growth of the vector norm after multiplying a vector by  $A$ .

## Identifying Matrix Norms

What is  $\|A\|_1$ ?  $\|A\|_\infty$ ?

$$\|A\|_1 = \max_{\text{col } j} \sum_{\text{row } i} |A_{i,j}|, \quad \|A\|_\infty = \max_{\text{row } i} \sum_{\text{col } j} |A_{i,j}|.$$

2-norm? Actually fairly difficult to evaluate. See in a bit.

How do matrix and vector norms relate for  $n \times 1$  matrices?

They agree. Why? For  $n \times 1$ , the vector  $\mathbf{x}$  in  $A\mathbf{x}$  is just a scalar:

$$\max_{\|\mathbf{x}\|=1} \|A\mathbf{x}\| = \max_{x \in \{-1,1\}} \|Ax\| = \|A[:, 1]\|$$

This can help to remember 1- and  $\infty$ -norm.

[Demo: Matrix norms](#) [\[cleared\]](#)

## Properties of Matrix Norms

Matrix norms inherit the vector norm properties:

- ▶  $\|A\| > 0 \Leftrightarrow A \neq 0$ .
- ▶  $\|\gamma A\| = |\gamma| \|A\|$  for all scalars  $\gamma$ .
- ▶ Obeys triangle inequality  $\|A + B\| \leq \|A\| + \|B\|$

But also some more properties that stem from our definition:

- ▶  $\|Ax\| \leq \|A\| \|x\|$
- ▶  $\|AB\| \leq \|A\| \|B\|$  (easy consequence)

Both of these are called *submultiplicativity* of the matrix norm.

In these notes: If we write  $\|\cdot\|$  (for matrix norms) without any specifics, then the statement is true for any induced norm. If a specific norm is needed, the notation will indicate that.

## Conditioning

What is the condition number of solving a linear system  $A\mathbf{x} = \mathbf{b}$ ?

**Input:**  $\mathbf{b}$  with error  $\Delta\mathbf{b}$ ,

**Output:**  $\mathbf{x}$  with error  $\Delta\mathbf{x}$ .

Observe  $A(\mathbf{x} + \Delta\mathbf{x}) = (\mathbf{b} + \Delta\mathbf{b})$ , so  $A\Delta\mathbf{x} = \Delta\mathbf{b}$ .

$$\begin{aligned} \frac{\text{rel err. in output}}{\text{rel err. in input}} &= \frac{\|\Delta\mathbf{x}\| / \|\mathbf{x}\|}{\|\Delta\mathbf{b}\| / \|\mathbf{b}\|} = \frac{\|\Delta\mathbf{x}\| \|\mathbf{b}\|}{\|\Delta\mathbf{b}\| \|\mathbf{x}\|} \\ &= \frac{\|A^{-1}\Delta\mathbf{b}\| \|A\mathbf{x}\|}{\|\Delta\mathbf{b}\| \|\mathbf{x}\|} \\ &\leq \|A^{-1}\| \|A\| \frac{\|\Delta\mathbf{b}\| \|\mathbf{x}\|}{\|\Delta\mathbf{b}\| \|\mathbf{x}\|} \\ &= \|A^{-1}\| \|A\|. \end{aligned}$$

## Conditioning of Linear Systems: Observations

Showed  $\kappa(\text{Solve } A\mathbf{x} = \mathbf{b}) \leq \|A^{-1}\| \|A\|$ .

I.e. found an *upper bound* on the condition number. With a little bit of fiddling, it's not too hard to find examples that achieve this bound, i.e. that it is *sharp*.

So we've found the *condition number of linear system solving*, also called the **condition number of the matrix  $A$** :

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



## Conditioning of Linear Systems: More properties

- ▶  $\text{cond}$  is relative to a given norm. So, to be precise, use

$$\text{cond}_2 \quad \text{or} \quad \text{cond}_\infty .$$

- ▶ If  $A^{-1}$  does not exist:  $\text{cond}(A) = \infty$  by convention.

What is  $\kappa(A^{-1})$ ?

$$\kappa(A)$$

What is the condition number of matrix-vector multiplication?

$\kappa(A)$  because it is equivalent to solving with  $A^{-1}$ .

[Demo: Condition number visualized](#) [\[cleared\]](#)

[Demo: Conditioning of 2x2 Matrices](#) [\[cleared\]](#)

## Residual Vector

What is the **residual vector** of solving the linear system

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

It's the thing that's 'left over'. Suppose our approximate solution is  $\hat{\mathbf{x}}$ . Then the residual vector is

$$\mathbf{r} = \mathbf{b} - A\hat{\mathbf{x}}.$$

## Residual and Error: Relationship

How do the (norms of the) residual vector  $\mathbf{r}$  and the error  $\Delta\mathbf{x} = \mathbf{x} - \hat{\mathbf{x}}$  relate to one another?

$$\|\Delta\mathbf{x}\| = \|\mathbf{x} - \hat{\mathbf{x}}\| = \|A^{-1}(\mathbf{b} - A\hat{\mathbf{x}})\| = \|A^{-1}\mathbf{r}\|$$

Divide both sides by  $\|\hat{\mathbf{x}}\|$ :

$$\frac{\|\Delta\mathbf{x}\|}{\|\hat{\mathbf{x}}\|} = \frac{\|A^{-1}\mathbf{r}\|}{\|\hat{\mathbf{x}}\|} \leq \frac{\|A^{-1}\| \|\mathbf{r}\|}{\|\hat{\mathbf{x}}\|} = \text{cond}(A) \frac{\|\mathbf{r}\|}{\|A\hat{\mathbf{x}}\|} \leq \text{cond}(A) \frac{\|\mathbf{r}\|}{\|A\hat{\mathbf{x}}\|}$$

- ▶ “rel err.”  $\leq$  cond  $\cdot$  rel. resid
- ▶ Given small (rel.) residual, (rel.) error is only (guaranteed to be) small if the condition number is also small.

## Changing the Matrix

So far, only discussed changing the RHS, i.e.  $A\mathbf{x} = \mathbf{b} \rightarrow A\hat{\mathbf{x}} = \hat{\mathbf{b}}$ .  
The matrix consists of FP numbers, too—it, too, is approximate. I.e.

$$A\mathbf{x} = \mathbf{b} \rightarrow \hat{A}\hat{\mathbf{x}} = \mathbf{b}.$$

What can we say about the error due to an approximate matrix?

Consider

$$\Delta\mathbf{x} = \hat{\mathbf{x}} - \mathbf{x} = A^{-1}(A\hat{\mathbf{x}} - \mathbf{b}) = A^{-1}(A\hat{\mathbf{x}} - \hat{A}\hat{\mathbf{x}}) = -A^{-1}\Delta A\hat{\mathbf{x}}.$$

Thus

$$\|\Delta\mathbf{x}\| \leq \|A^{-1}\| \|\Delta A\| \|\hat{\mathbf{x}}\|.$$

And we get

$$\frac{\|\Delta\mathbf{x}\|}{\|\hat{\mathbf{x}}\|} \leq \text{cond}(A) \frac{\|\Delta A\|}{\|A\|}.$$

## Changing Condition Numbers

Once we have a matrix  $A$  in a linear system  $A\mathbf{x} = \mathbf{b}$ , are we stuck with its condition number? Or could we improve it?

### *Preconditioning*

- ▶ **Left' preconditioning:**  $MA\mathbf{x} = M\mathbf{b}$
- ▶ **Right preconditioning:**  $AM\hat{\mathbf{x}} = \mathbf{b}$   
Different  $\hat{\mathbf{x}}$ : Recover  $\mathbf{x} = M\hat{\mathbf{x}}$ .

A typical case: use diagonal matrix as the preconditioner. What is the effect in each case?

- ▶ Row-wise scaling:  $DA\mathbf{x} = D\mathbf{b}$
- ▶ Column-wise scaling:  $AD\hat{\mathbf{x}} = \mathbf{b}$   
Different  $\hat{\mathbf{x}}$ : Recover  $\mathbf{x} = D\hat{\mathbf{x}}$ .

## In-Class Activity: Matrix Norms and Conditioning

In-class activity: Matrix Norms and Conditioning

## Recap: Orthogonal Matrices

What's an *orthogonal* (=orthonormal) matrix?

One that satisfies  $Q^T Q = I$  and  $Q Q^T = I$ .

How do orthogonal matrices interact with the 2-norm?

$$\|Q\mathbf{v}\|_2^2 = (Q\mathbf{v})^T(Q\mathbf{v}) = \mathbf{v}^T Q^T Q \mathbf{v} = \mathbf{v}^T \mathbf{v} = \|\mathbf{v}\|_2^2.$$

## Singular Value Decomposition (SVD)

What is the *Singular Value Decomposition* of an  $m \times n$  matrix?

$$A = U\Sigma V^T,$$

with

- ▶  $U$  is  $m \times m$  and orthogonal  
Columns called the **left singular vectors**.
- ▶  $\Sigma = \text{diag}(\sigma_i)$  is  $m \times n$  and non-negative  
Typically  $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_{\min(m,n)} \geq 0$ .  
Called the **singular values**.
- ▶  $V$  is  $n \times n$  and orthogonal  
Columns called the **right singular vectors**.

**Existence, Computation:** Not yet, later.



## Computing the 2-Norm

Using the SVD of  $A$ , identify the 2-norm.

$A = U\Sigma V^T$  with  $U, V$  orthogonal.

- ▶ 2-norm satisfies  $\|QB\|_2 = \|B\|_2 = \|BQ\|_2$  for any matrix  $B$  and orthogonal  $Q$ .
- ▶ So  $\|A\|_2 = \|\Sigma\|_2 = \sigma_{\max}$

Express the matrix condition number  $\text{cond}_2(A)$  in terms of the SVD:

- ▶  $A^{-1}$  has singular values  $1/\sigma_i$ .
- ▶  $\text{cond}_2(A) = \|A\|_2 \|A^{-1}\|_2 = \sigma_{\max}/\sigma_{\min}$

## Not a matrix norm: Frobenius

The 2-norm is very costly to compute. Can we make something simpler?

$$\|A\|_F = \sqrt{\sum_{i=1}^m \sum_{j=1}^n |a_{ij}|^2}$$

is called the **Frobenius norm**.

What about its properties?

Satisfies the mat. norm properties.

- ▶ definiteness
- ▶ scaling
- ▶ triangle inequality
- ▶ submultiplicativity (proof via Cauchy-Schwarz)

## Frobenius Norm: Properties

Is the Frobenius norm induced by any vector norm?

Can't be! What's  $\|I\|_F$ ? What's  $\|I\|$  for an induced norm?

How does it relate to the SVD?

$$\|A\|_F = \sqrt{\sum_{i=1}^n \sigma_i^2}$$

(Proof?)

## Solving Systems: Simple cases

Solve  $D\mathbf{x} = \mathbf{b}$  if  $D$  is diagonal. (Computational cost?)

$$x_i = b_i / D_{ii} \text{ with cost } O(n)$$

Solve  $Q\mathbf{x} = \mathbf{b}$  if  $Q$  is orthogonal. (Computational cost?)

$$\mathbf{x} = Q^T \mathbf{b} \text{ with cost } O(n^2).$$

Given SVD  $A = U\Sigma V^T$ , solve  $A\mathbf{x} = \mathbf{b}$ . (Computational cost?)

- ▶ Compute  $\mathbf{z} = U^T \mathbf{b}$
- ▶ Solve  $\Sigma \mathbf{y} = \mathbf{z}$
- ▶ Compute  $\mathbf{x} = V \mathbf{y}$

Cost:  $O(n^2)$  to solve, and  $O(n^3)$  to compute SVD.

## Solving Systems: Triangular matrices

Solve

$$\begin{bmatrix} a_{11} & a_{12} & a_{13} & a_{14} \\ & a_{22} & a_{23} & a_{24} \\ & & a_{33} & a_{34} \\ & & & a_{44} \end{bmatrix} \begin{bmatrix} x \\ y \\ z \\ w \end{bmatrix} = \begin{bmatrix} b_1 \\ b_2 \\ b_3 \\ b_4 \end{bmatrix}.$$

- ▶ Rewrite as individual equations.
- ▶ This process is called **back-substitution**.
- ▶ The analogous process for lower triangular matrices is called **forward substitution**.

**Demo: Coding back-substitution [cleared]**

What about non-triangular matrices?

Can do *Gaussian Elimination*, just like in linear algebra class.

## Gaussian Elimination

### Demo: Vanilla Gaussian Elimination [cleared]

What do we get by doing Gaussian Elimination?

*Row Echelon Form.*

How is that different from being upper triangular?

- ▶ REF reveals the rank of the matrix.
- ▶ REF can take “multiple column-steps” to the right per row.

What if we do not just eliminate downward but also upward?

That's called *Gauss-Jordan elimination*. Turns out to be computationally inefficient. We won't look at it.

# LU Factorization

What is the LU factorization?

A factorization  $A = LU$  with:

- ▶  $L$  lower triangular, unit diagonal
- ▶  $U$  upper triangular

Existence? Good question  $\rightarrow$  will answer later.

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

Does LU help solve  $A\mathbf{x} = \mathbf{b}$ ?

$$\begin{aligned} A\mathbf{x} &= \mathbf{b} \\ L \underbrace{U\mathbf{x}}_y &= \mathbf{b} \\ L\mathbf{y} &= \mathbf{b} \quad \leftarrow \text{solvable by fwd. subst.} \\ U\mathbf{x} &= \mathbf{y} \quad \leftarrow \text{solvable by bwd. subst.} \end{aligned}$$

Now know  $\mathbf{x}$  that solves  $A\mathbf{x} = \mathbf{b}$ .



## Determining an LU factorization

$$\begin{bmatrix} a_{11} & \mathbf{a}_{12}^T \\ \mathbf{a}_{21} & A_{22} \end{bmatrix} = \begin{bmatrix} L_{11} & \\ L_{21} & L_{22} \end{bmatrix} \begin{bmatrix} U_{11} & U_{12} \\ & U_{22} \end{bmatrix}.$$

Or, written differently:

$$\begin{bmatrix} 1 & \\ \ell_{21} & L_{22} \end{bmatrix} \begin{bmatrix} u_{11} & \mathbf{u}_{12}^T \\ & U_{22} \\ \mathbf{a}_{11} & \mathbf{a}_{12} \\ \mathbf{a}_{21} & A_{22} \end{bmatrix}$$

- ▶ Clear:  $u_{11} = a_{11}$ ,  $\mathbf{u}_{12}^T = \mathbf{a}_{12}^T$ .
- ▶  $\mathbf{a}_{21} = u_{11}\ell_{21}$ , or  $\ell_{21} = \mathbf{a}_{21}/u_{11}$ .
- ▶  $A_{22} = \ell_{21}\mathbf{u}_{12}^T + L_{22}U_{22}$ , or  $L_{22}U_{22} = A_{22} - \ell_{21}\mathbf{u}_{12}^T$ .

Demo: LU Factorization [cleared]

## Computational Cost

What is the computational cost of multiplying two  $n \times n$  matrices?

$$O(n^3)$$

- ▶  $u_{11} = a_{11}, \mathbf{u}_{12}^T = \mathbf{a}_{12}^T.$
- ▶  $\ell_{21} = \mathbf{a}_{21}/u_{11}.$
- ▶  $L_{22}U_{22} = A_{22} - \ell_{21}\mathbf{u}_{12}^T.$

What is the computational cost of carrying out LU factorization on an  $n \times n$  matrix?

$$O(n^2) \text{ for each step, } n - 1 \text{ of those steps: } O(n^3).$$

[Demo: Complexity of Mat-Mat multiplication and LU \[cleared\]](#)

## LU: Failure Cases?

Is LU/Gaussian Elimination bulletproof?

Not bulletproof:

$$A = \begin{bmatrix} 0 & 1 \\ 2 & 1 \end{bmatrix}.$$

Q: Is this a problem with the process or with the entire *idea* of LU?

$$\begin{bmatrix} u_{11} & u_{12} \\ & u_{22} \end{bmatrix} \begin{bmatrix} 1 & \\ \ell_{21} & 1 \end{bmatrix} \begin{bmatrix} 0 & 1 \\ 2 & 1 \end{bmatrix} \rightarrow u_{11} = 0$$
$$\underbrace{u_{11} \cdot \ell_{21}}_0 + 1 \cdot 0 = 2$$

It turns out to be that  $A$  doesn't *have* an LU factorization.

LU has exactly one failure mode: the division when  $u_{11} = 0$ .

## Saving the LU Factorization

What can be done to get something *like* an LU factorization?

**Idea from linear algebra class:** In Gaussian elimination, simply swap rows, equivalent linear system.

- ▶ Good idea: Swap rows if there's a zero in the way
- ▶ Even better idea: Find the largest entry (by absolute value), swap it to the top row.

The entry we divide by is called the *pivot*.

- ▶ Swapping rows to get a bigger pivot is called **partial pivoting**.
- ▶ Swapping rows *and columns* to get an even bigger pivot is called **complete pivoting**. (downside: additional  $O(n^2)$  cost *per step* to find the pivot!)

**Demo:** [LU Factorization with Partial Pivoting](#) [cleared]

## Cholesky: LU for Symmetric Positive Definite

LU *can* be used for SPD matrices. But can we do better?

**Cholesky factorization:**  $A = LL^T$ , i.e. *like* LU, but using  $L^T$  for  $U$ .

$$\begin{bmatrix} \ell_{11} & \\ \ell_{21} & L_{22} \end{bmatrix} \begin{bmatrix} \ell_{11} & \ell_{21}^T \\ & L_{22}^T \end{bmatrix} = \begin{bmatrix} a_{11} & \mathbf{a}_{21}^T \\ \mathbf{a}_{21} & A_{22} \end{bmatrix}.$$

$\ell_{11}^2 = a_{11}$ , then  $\ell_{11}\ell_{21} = \mathbf{a}_{21}$ , i.e.  $\ell_{21} = \mathbf{a}_{21}/\ell_{11}$ . Finally,

$$\begin{aligned} \ell_{21}\ell_{21}^T + L_{22}L_{22}^T &= A_{22}, \quad \text{or} \\ L_{22}L_{22}^T &= A_{22} - \ell_{21}\ell_{21}^T. \end{aligned}$$

- ▶ Fails if  $a_{11}$  is negative at any point. ( $\Rightarrow A$  not SPSEmiD)
- ▶ If  $a_{11}$  zero:  $A$  is positive *semidefinite*.
- ▶ *Cheaper than LU*: no pivoting, only one factor to compute!

## More cost concerns

What's the cost of solving  $A\mathbf{x} = \mathbf{b}$ ?

LU:  $O(n^3)$

FW/BW Subst:  $2 \times O(n^2) = O(n^2)$

What's the cost of solving  $A\mathbf{x} = \mathbf{b}_1, \mathbf{b}_2, \dots, \mathbf{b}_n$ ?

LU:  $O(n^3)$

FW/BW Subst:  $2n \times O(n^2) = O(n^3)$

What's the cost of finding  $A^{-1}$ ?

Same as solving

$$AX = I,$$

so still  $O(n^3)$ .

## Cost: Worrying about the Constant, BLAS

$O(n^3)$  really means

$$\alpha \cdot n^3 + \beta \cdot n^2 + \gamma \cdot n + \delta.$$

All the non-leading and constants terms swept under the rug. But: at least the leading constant ultimately matters.

Shrinking the constant: surprisingly hard (even for 'just' matmul)

**Idea:** Rely on library implementation: *BLAS* (Fortran)

Level 1  $\mathbf{z} = \alpha \mathbf{x} + \mathbf{y}$  vector-vector operations

$$O(n)$$

?axpy

Level 2  $\mathbf{z} = \mathbf{A}\mathbf{x} + \mathbf{y}$  matrix-vector operations

$$O(n^2)$$

?gemv

Level 3  $\mathbf{C} = \mathbf{A}\mathbf{B} + \beta \mathbf{C}$  matrix-matrix operations

$$O(n^3)$$

?gemm, ?trsm

**Show (using perf):** numpy matmul calls BLAS dgemm

# LAPACK

LAPACK: Implements 'higher-end' things (such as LU) using BLAS  
Special matrix formats can also help save const significantly, e.g.

- ▶ banded
- ▶ sparse
- ▶ symmetric
- ▶ triangular

Sample routine names:

- ▶ dgesvd, zgesdd
- ▶ dgetrf, dgetrs



## LU on Blocks: The Schur Complement

Given a linear system

$$\left[ \begin{array}{cc|c} A & B & \mathbf{b}_1 \\ C & D & \mathbf{b}_2 \end{array} \right],$$

can we do 'block Gaussian elimination' to get a *block triangular matrix*?

Multiply the top row by  $-CA^{-1}$ , add to second row, gives:

$$\left[ \begin{array}{cc|c} A & B & \mathbf{b}_1 \\ 0 & D - CA^{-1}B & \mathbf{b}_2 - CA^{-1}\mathbf{b}_1 \end{array} \right],$$

- ▶  $D - CA^{-1}B$  is called the **Schur complement**.
- ▶ Block pivoting is also possible if needed.

## LU: Special cases

What happens if we feed a non-invertible matrix to LU?

$$PA = LU$$

(invertible, not invertible) (Why?)

What happens if we feed LU an  $m \times n$  non-square matrices?

Think carefully about sizes of factors and columns/rows that do/don't matter. Two cases:

- ▶  $m > n$  (tall&skinny):  $L : m \times n$ ,  $U : n \times n$
- ▶  $m < n$  (short&fat):  $L : m \times m$ ,  $U : m \times n$

This is called **reduced LU factorization**.

## Round-off Error in LU without Pivoting

Consider factorization of  $\begin{bmatrix} \epsilon & 1 \\ 1 & 1 \end{bmatrix}$  where  $\epsilon < \epsilon_{\text{mach}}$ :

▶ Without pivoting:  $L = \begin{bmatrix} 1 & 0 \\ 1/\epsilon & 1 \end{bmatrix}$ ,  $U = \begin{bmatrix} \epsilon & 1 \\ 0 & 1 - 1/\epsilon \end{bmatrix}$

▶ Rounding:  $\text{fl}(U) = \begin{bmatrix} \epsilon & 1 \\ 0 & -1/\epsilon \end{bmatrix}$

▶ This leads to  $L \text{fl}(U) = \begin{bmatrix} \epsilon & 1 \\ 1 & 0 \end{bmatrix}$ , a backward error of  $\begin{bmatrix} 0 & 0 \\ 0 & 1 \end{bmatrix}$

## Round-off Error in LU with Pivoting

Permuting the rows of  $A$  in partial pivoting gives  $PA = \begin{bmatrix} 1 & 1 \\ \epsilon & 1 \end{bmatrix}$

- ▶ We now compute  $L = \begin{bmatrix} 1 & 0 \\ \epsilon & 1 \end{bmatrix}$ ,  $U = \begin{bmatrix} 1 & 1 \\ 0 & 1 - \epsilon \end{bmatrix}$ , so

$$\text{fl}(U) = \begin{bmatrix} 1 & 1 \\ 0 & 1 \end{bmatrix}$$

- ▶ This leads to  $L \text{fl}(U) = \begin{bmatrix} 1 & 1 \\ \epsilon & 1 + \epsilon \end{bmatrix}$ , a backward error of  $\begin{bmatrix} 0 & 0 \\ 0 & \epsilon \end{bmatrix}$ .

## Changing matrices

Seen: LU cheap to re-solve if RHS changes. (Able to keep the expensive bit, the LU factorization) What if the *matrix* changes?

Special cases allow something to be done (a so-called *rank-one update*):

$$\hat{A} = A + \mathbf{u}\mathbf{v}^T$$

The **Sherman-Morrison formula** gives us

$$(A + \mathbf{u}\mathbf{v}^T)^{-1} = A^{-1} - \frac{A^{-1}\mathbf{u}\mathbf{v}^T A^{-1}}{1 + \mathbf{v}^T A^{-1}\mathbf{u}}.$$

Proof: Multiply the above by  $\hat{A}$  get the identity.

FYI: There is a rank- $k$  analog called the **Sherman-Morrison-Woodbury formula**.

[Demo: Sherman-Morrison \[cleared\]](#)

## In-Class Activity: LU

In-class activity: LU and Cost

# Outline

Introduction to Scientific Computing

Systems of Linear Equations

**Linear Least Squares**

Introduction

Sensitivity and Conditioning

Solving Least Squares

Eigenvalue Problems

Nonlinear Equations

Optimization

Interpolation

Numerical Integration and Differentiation

Initial Value Problems for ODEs

Boundary Value Problems for ODEs

Partial Differential Equations and Sparse Linear Algebra

Fast Fourier Transform

Additional Topics

## What about non-square systems?

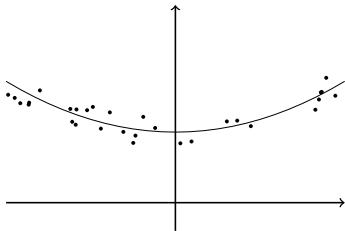
Specifically, what about linear systems with 'tall and skinny' matrices? (A:  $m \times n$  with  $m > n$ ) (aka *overdetermined* linear systems)

Specifically, any hope that we will solve those exactly?

Not really: more equations than unknowns.



## Example: Data Fitting



Have data:  $(x_i, y_i)$  and model:

$$y(x) = \alpha + \beta x + \gamma x^2$$

Find data that (best) fit model!

## Data Fitting Continued

$$\begin{aligned}\alpha + \beta x_1 + \gamma x_1^2 &= y_1 \\ &\vdots \\ \alpha + \beta x_n + \gamma x_n^2 &= y_n\end{aligned}$$

Not going to happen for  $n > 3$ . Instead:

$$\begin{aligned}&|\alpha + \beta x_1 + \gamma x_1^2 - y_1|^2 \\ &\quad + \dots + \\ &|\alpha + \beta x_n + \gamma x_n^2 - y_n|^2 \rightarrow \min!\end{aligned}$$

→ *Least Squares*

This is called *linear least squares* specifically because the coefficients  $\mathbf{x}$  enter linearly into the residual.

## Rewriting Data Fitting

Rewrite in matrix form.

$$\|A\mathbf{x} - \mathbf{b}\|_2^2 \rightarrow \min!$$

with

$$A = \begin{bmatrix} 1 & x_1 & x_1^2 \\ \vdots & \vdots & \vdots \\ 1 & x_n & x_n^2 \end{bmatrix}, \quad \mathbf{x} = \begin{bmatrix} \alpha \\ \beta \\ \gamma \end{bmatrix}, \quad \mathbf{b} = \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix}$$

- ▶ Matrices like  $A$  are called **Vandermonde matrices**.
- ▶ Easy to generalize to higher polynomial degrees.

## Least Squares: The Problem In Matrix Form

$$\|A\mathbf{x} - \mathbf{b}\|_2^2 \rightarrow \min!$$

is cumbersome to write.

Invent new notation, defined to be equivalent:

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

### NOTE:

- ▶ Data Fitting is *one example* where LSQ problems arise.
- ▶ Many other application lead to  $A\mathbf{x} \cong \mathbf{b}$ , with different matrices.

## Data Fitting: Nonlinearity

Give an example of a nonlinear data fitting problem.

$$\begin{aligned} & |\exp(\alpha) + \beta x_1 + \gamma x_1^2 - y_1|^2 \\ & \quad + \dots + \\ & |\exp(\alpha) + \beta x_n + \gamma x_n^2 - y_n|^2 \rightarrow \min! \end{aligned}$$

But that would be easy to remedy: Do linear least squares with  $\exp(\alpha)$  as the unknown. More difficult:

$$\begin{aligned} & |\alpha + \exp(\beta x_1 + \gamma x_1^2) - y_1|^2 \\ & \quad + \dots + \\ & |\alpha + \exp(\beta x_n + \gamma x_n^2) - y_n|^2 \rightarrow \min! \end{aligned}$$

[Demo: Interactive Polynomial Fit \[cleared\]](#)

## Properties of Least-Squares

Consider LSQ problem  $A\mathbf{x} \cong \mathbf{b}$  and its associated *objective function*  $\varphi(\mathbf{x}) = \|\mathbf{b} - A\mathbf{x}\|_2^2$ . Assume  $A$  has full rank. Does this always have a solution?

Yes.  $\varphi \geq 0$ ,  $\varphi \rightarrow \infty$  as  $\|\mathbf{x}\| \rightarrow \infty$ ,  $\varphi$  continuous  $\Rightarrow$  has a minimum.

Is it always unique?

No, for example if  $A$  has a nullspace.

## Least-Squares: Finding a Solution by Minimization

Examine the objective function, find its minimum.

$$\begin{aligned}\varphi(\mathbf{x}) &= (\mathbf{b} - A\mathbf{x})^T (\mathbf{b} - A\mathbf{x}) \\ &= \mathbf{b}^T \mathbf{b} - 2\mathbf{x}^T A^T \mathbf{b} + \mathbf{x}^T A^T A \mathbf{x}\end{aligned}$$

$$\nabla \varphi(\mathbf{x}) = -2A^T \mathbf{b} + 2A^T A \mathbf{x}$$

$\nabla \varphi(\mathbf{x}) = 0$  yields  $A^T A \mathbf{x} = A^T \mathbf{b}$ . Called the *normal equations*.

## Least squares: Demos

[Demo: Polynomial fitting with the normal equations](#) [cleared]

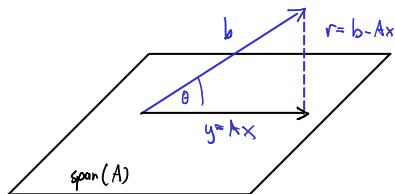
What's the shape of  $A^T A$ ?

Always square.

[Demo: Issues with the normal equations](#) [cleared]



## Least Squares, Viewed Geometrically



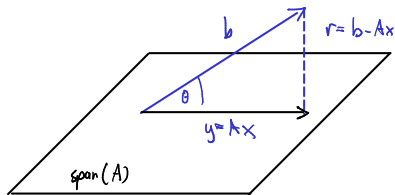
Why is  $\mathbf{r} \perp \text{span}(A)$  a good thing to require?

Because then the distance between  $\mathbf{y} = \mathbf{Ax}$  and  $\mathbf{b}$  is minimal.

Q: Why?

Because of Pythagoras's theorem—another  $\mathbf{y}$  would mean additional distance traveled in  $\text{span}(A)$ .

## Least Squares, Viewed Geometrically (II)



Phrase the Pythagoras observation as an equation.

$$\begin{aligned}\text{span}(A) &\perp \mathbf{b} - A\mathbf{x} \\ A^T \mathbf{b} - A^T A\mathbf{x} &= 0\end{aligned}$$

Congratulations: Just rediscovered the normal equations.

Write that with an orthogonal projection matrix  $P$ .

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

## About Orthogonal Projectors

What is a *projector*?

A matrix satisfying  $P^2 = P$ .

What is an *orthogonal projector*?

A symmetric projector.

How do I make one projecting onto  $\text{span}\{\mathbf{q}_1, \mathbf{q}_2, \dots, \mathbf{q}_\ell\}$  for orthogonal  $\mathbf{q}_i$ ?

First define  $Q = [\mathbf{q}_1 \ \mathbf{q}_2 \ \dots \ \mathbf{q}_\ell]$ . Then

$$QQ^T$$

will project and is obviously symmetric.

## Least Squares and Orthogonal Projection

Check that  $P = A(A^T A)^{-1}A^T$  is an orthogonal projector onto  $\text{colspan}(A)$ .

$$P^2 = A(A^T A)^{-1}A^T A(A^T A)^{-1}A^T = A(A^T A)^{-1}A^T = P.$$

Symmetry: also yes.

Onto  $\text{colspan}(A)$ : Last matrix is  $A \rightarrow$  result of  $P\mathbf{x}$  must be in  $\text{colspan}(A)$ .

**Conclusion:**  $P$  is the projector from the previous slide!

What assumptions do we need to define the  $P$  from the last question?

$A^T A$  has full rank (i.e. is invertible).

## Pseudoinverse

What is the **pseudoinverse** of  $A$ ?

Nonsquare  $m \times n$  matrix  $A$  has no inverse in usual sense.  
If  $\text{rank}(A) = n$ , **pseudoinverse** is  $A^+ = (A^T A)^{-1} A^T$ . (colspan-projector with final  $A$  missing)

What can we say about the condition number in the case of a tall-and-skinny, full-rank matrix?

$$\text{cond}_2(A) = \|A\|_2 \|A^+\|_2$$

If not full rank,  $\text{cond}(A) = \infty$  by convention.

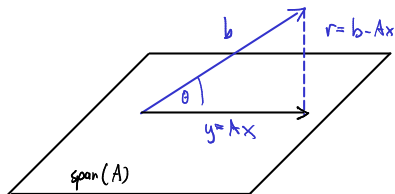
What does all this have to do with solving least squares problems?

$$\mathbf{x} = A^+ \mathbf{b} \text{ solves } A\mathbf{x} \cong \mathbf{b}.$$

## In-Class Activity: Least Squares

In-class activity: Least Squares

## Sensitivity and Conditioning of Least Squares



Relate  $\|\mathbf{A}\mathbf{x}\|$  and  $\mathbf{b}$  with  $\theta$  via trig functions.

$$\cos(\theta) = \frac{\|\mathbf{A}\mathbf{x}\|_2}{\|\mathbf{b}\|_2},$$

## Sensitivity and Conditioning of Least Squares (II)

Derive a conditioning bound for the least squares problem.

Recall  $\mathbf{x} = A^+ \mathbf{b}$ . Also  $\Delta \mathbf{x} = A^+ \Delta \mathbf{b}$ . Take norms, divide by  $\|\mathbf{x}\|_2$ :

$$\begin{aligned} \frac{\|\Delta \mathbf{x}\|_2}{\|\mathbf{x}\|_2} &\leq \|A^+\|_2 \frac{\|\Delta \mathbf{b}\|_2}{\|\mathbf{x}\|_2} \\ &= \frac{\kappa(A)}{\|A\|_2 \|A^+\|_2} \|A^+\|_2 \frac{\|\mathbf{b}\|_2}{\|\mathbf{b}\|_2} \frac{\|\Delta \mathbf{b}\|_2}{\|\mathbf{x}\|_2} \\ &= \kappa(A) \frac{\|\mathbf{b}\|_2}{\|A\|_2 \|\mathbf{x}\|_2} \frac{\|\Delta \mathbf{b}\|_2}{\|\mathbf{b}\|_2} \leq \kappa(A) \underbrace{\frac{\|\mathbf{b}\|_2}{\|A\mathbf{x}\|_2}}_{1/\cos\theta} \frac{\|\Delta \mathbf{b}\|_2}{\|\mathbf{b}\|_2}. \end{aligned}$$

What values of  $\theta$  are bad?

$\mathbf{b} \perp \text{colspan}(A)$ , i.e.  $\theta \approx \pi/2$ .



## Sensitivity and Conditioning of Least Squares (III)

Any comments regarding dependencies?

Unlike for  $A\mathbf{x} = \mathbf{b}$ , the sensitivity of least squares solution depends on both  $A$  and  $\mathbf{b}$ .

What about changes in the matrix?

$$\frac{\|\Delta\mathbf{x}\|_2}{\|\mathbf{x}\|_2} \leq [\text{cond}(A)^2 \tan(\theta) + \text{cond}(A)] \cdot \frac{\|\Delta A\|_2}{\|A\|_2}.$$

Two behaviors:

- ▶ If  $\tan(\theta) \approx 0$ , condition number is  $\text{cond}(A)$ .
- ▶ Otherwise,  $\text{cond}(A)^2$ .

## Transforming Least Squares to Upper Triangular

Suppose we have  $A = QR$ , with  $Q$  square and orthogonal, and  $R$  upper triangular. This is called a **QR factorization**.

How do we transform the least squares problem  $A\mathbf{x} \cong \mathbf{b}$  to one with an upper triangular matrix?

$$\begin{aligned} & \|A\mathbf{x} - \mathbf{b}\|_2 \\ &= \left\| Q^T(QR\mathbf{x} - \mathbf{b}) \right\|_2 \\ &= \left\| R\mathbf{x} - Q^T\mathbf{b} \right\|_2 \end{aligned}$$

## Simpler Problems: Triangular

What do we win from transforming a least-squares system to upper triangular form?

$$[R_{\text{top}}] \mathbf{x} \cong \begin{bmatrix} (Q^T \mathbf{b})_{\text{top}} \\ (Q^T \mathbf{b})_{\text{bottom}} \end{bmatrix}$$

How would we minimize the residual norm?

For the residual vector  $\mathbf{r}$ , we find

$$\|\mathbf{r}\|_2^2 = \left\| (Q^T \mathbf{b})_{\text{top}} - R_{\text{top}} \mathbf{x} \right\|_2^2 + \left\| (Q^T \mathbf{b})_{\text{bottom}} \right\|_2^2.$$

$R_{\text{top}}$  is invertible, so we can find  $\mathbf{x}$  to zero out the first term, leaving

$$\|\mathbf{r}\|_2^2 = \left\| (Q^T \mathbf{b})_{\text{bottom}} \right\|_2^2.$$

## Computing QR

- ▶ Gram-Schmidt
- ▶ Householder Reflectors
- ▶ Givens Rotations

Demo: Gram-Schmidt–The Movie [cleared] (shows *modified G-S*)

Demo: Gram-Schmidt and Modified Gram-Schmidt [cleared]

Demo: Keeping track of coefficients in Gram-Schmidt [cleared]

Seen: Even modified Gram-Schmidt still unsatisfactory in finite precision arithmetic because of roundoff.

**NOTE:** Textbook makes further modification to ‘modified’ Gram-Schmidt:

- ▶ Orthogonalize *subsequent* rather than *preceding* vectors.
- ▶ Numerically: no difference, but sometimes algorithmically helpful.

## Economical/Reduced QR

Is QR with square  $Q$  for  $A \in \mathbb{R}^{m \times n}$  with  $m > n$  efficient?

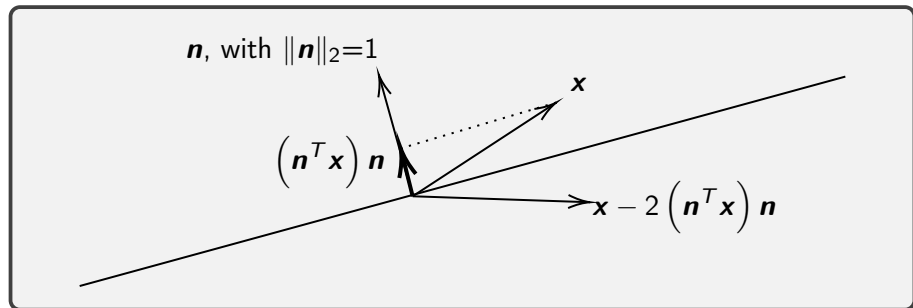
No. Can obtain **economical** or **reduced QR** with  $Q \in \mathbb{R}^{m \times n}$  and  $R \in \mathbb{R}^{n \times n}$ . Least squares solution process works unmodified with the economical form, though the equivalence proof relies on the 'full' form.

In-Class Activity: QR

In-class activity: QR

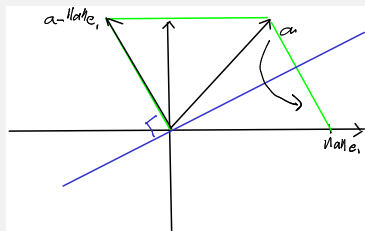
## Constructing Reflections

Given a plane represented by its (unit) normal vector  $\mathbf{n}$ , construct a reflection about that plane.



## Householder Transformations

Find an *orthogonal* matrix  $Q$  to zero out the lower part of a vector  $\mathbf{a}$ .



Orthogonality in figure:  $(\mathbf{a} - \|\mathbf{a}\|_2 \mathbf{e}_1) \cdot (\mathbf{a} + \|\mathbf{a}\|_2 \mathbf{e}_1) = \|\mathbf{a}\|_2^2 - \|\mathbf{a}\|_2^2$ .

Let's call  $\mathbf{v} = \mathbf{a} - \|\mathbf{a}\|_2 \mathbf{e}_1$ . How do we reflect about the plane orthogonal to  $\mathbf{v}$ ? Project-and-keep-going:

$$H := I - 2 \frac{\mathbf{v}\mathbf{v}^T}{\mathbf{v}^T \mathbf{v}}.$$

This is called a **Householder reflector**.



## Householder Reflectors: Properties

Seen from picture (and easy to see with algebra):

$$H\mathbf{a} = \pm \|\mathbf{a}\|_2 \mathbf{e}_1.$$

Remarks:

- ▶ **Q:** What if we want to zero out only the  $i + 1$ th through  $n$ th entry?  
**A:** Use  $\mathbf{e}_i$  above.
- ▶ A product  $H_n \cdots H_1 A = R$  of Householders makes it easy (and quite efficient!) to build a QR factorization.
- ▶ It turns out  $\mathbf{v}' = \mathbf{a} + \|\mathbf{a}\|_2 \mathbf{e}_1$  works out, too—just pick whichever one causes less cancellation.
- ▶  $H$  is symmetric
- ▶  $H$  is orthogonal

[Demo: 3x3 Householder demo](#) [cleared]

## Givens Rotations

If reflections work, can we make rotations work, too?

$$\begin{bmatrix} c & s \\ -s & c \end{bmatrix} \begin{bmatrix} a_1 \\ a_2 \end{bmatrix} = \begin{bmatrix} \sqrt{a_1^2 + a_2^2} \\ 0 \end{bmatrix}.$$

where  $c = a_1 / \|\mathbf{a}\|_2$  and  $s = a_2 / \|\mathbf{a}\|_2$ , with  $\mathbf{a} = [a_1, a_2]^T$ .

**Downside?** Produces only one zero at a time.

[Demo: 3x3 Givens demo \[cleared\]](#)

## Givens Rotations: Elimination Order

Given a matrix

$$A = \begin{bmatrix} a_{11} & a_{12} & a_{13} \\ a_{21} & a_{22} & a_{23} \\ a_{31} & a_{32} & a_{33} \end{bmatrix},$$

in what order can we apply Givens rotations to eliminate the nonzeros below the diagonal?

$$G_{32}G_{21}G_{31}A = R,$$

where  $G_{ij}$  is the Givens rotation that zeros out the entry at row  $i$ , column  $j$  of the matrix preceding it. (NB: not of  $A$ , as the entries change as soon as other Givens rotations are applied.)

## Rank-Deficient Matrices and QR

What happens with QR for rank-deficient matrices?

$A = QR$ , where  $R$  has some zero diagonal entries, in undetermined order.

Practically, it makes sense to ask for all these 'small' columns to be gathered near the 'right' of  $R \rightarrow$  Column pivoting.

Q: What does the resulting factorization look like?

$$AP = QR$$

$$AP = Q \begin{bmatrix} * & * & * \\ & \text{(small)} & \text{(small)} \\ & & \text{(smaller)} \end{bmatrix}$$

Also used as the basis for *rank-revealing QR*.

## Rank-Deficient Matrices and Least-Squares

What happens with Least Squares for rank-deficient matrices?

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

- ▶ QR still finds a solution with minimal residual  
But: need to deal with zeros on diagonal in  $R$
- ▶ By QR it's easy to see that least squares with a short-and-fat matrix is equivalent to a rank-deficient one.
- ▶ **But:** No longer unique.  $\mathbf{x} + \mathbf{n}$  for  $\mathbf{n} \in N(A)$  has the same residual.
- ▶ **In other words:** Have more freedom  
**Or:** Can demand another condition, for example:
  - ▶ Minimize  $\|\mathbf{b} - A\mathbf{x}\|_2^2$ , and
  - ▶ minimize  $\|\mathbf{x}\|_2^2$ , simultaneously.  
Unfortunately, QR does not help much with that  $\rightarrow$  Need better tool, the SVD  $A = U\Sigma V^T$ . Let's learn more about it.

## SVD: Reduced and Full

For a matrix of shape  $m \times n$  with  $m > n$ , what are the shapes of the factors in the SVD?

Again, there is the **full** version of the factorization:

- ▶  $U: m \times m$
- ▶  $\Sigma: m \times n$
- ▶  $V: n \times n$

and the **economical/reduced** version:

- ▶  $U: m \times n$
- ▶  $\Sigma: n \times n$
- ▶  $V: n \times n$

## SVD: What's this thing good for? (I)

- ▶ Recall:  $\|A\|_2 = \sigma_1$
- ▶ Recall:  $\text{cond}_2(A) = \sigma_1/\sigma_n$
- ▶ Nullspace  $N(A) = \text{span}(\{\mathbf{v}_i : \sigma_i = 0\})$ .
- ▶  $\text{rank}(A) = \#\{i : \sigma_i \neq 0\}$

Computing rank in the presence of round-off error is laughably non-robust. More robust:

- ▶ *Numerical rank:*

$$\text{rank}_\varepsilon = \#\{i : \sigma_i > \varepsilon\}$$

## SVD: What's this thing good for? (II)

### ► *Low-rank Approximation*

#### Theorem (Eckart-Young-Mirsky)

If  $k < r = \text{rank}(A)$  and

$$A_k = \sum_{i=1}^k \sigma_i u_i v_i^T, \quad \text{then}$$

$$\min_{\text{rank}(B)=k} \|A - B\|_2 = \|A - A_k\|_2 = \sigma_{k+1},$$

$$\min_{\text{rank}(B)=k} \|A - B\|_F = \|A - A_k\|_F = \sqrt{\sum_{j=k+1}^n \sigma_j^2}.$$



## SVD: What's this thing good for? (III)

- ▶ The minimum norm solution to  $A\mathbf{x} \cong \mathbf{b}$ :

$$\begin{aligned}U\Sigma V^T \mathbf{x} &\cong \mathbf{b} \\ \Leftrightarrow \Sigma \underbrace{V^T \mathbf{x}}_y &\cong U^T \mathbf{b} \\ \Leftrightarrow \Sigma \mathbf{y} &\cong U^T \mathbf{b}\end{aligned}$$

Then define

$$\Sigma^+ = \text{diag}(\sigma_1^+, \dots, \sigma_n^+),$$

where  $\Sigma^+$  is  $n \times m$  if  $A$  is  $m \times n$ , and

$$\sigma_i^+ = \begin{cases} 1/\sigma_i & \sigma_i \neq 0, \\ 0 & \sigma_i = 0. \end{cases}$$

## SVD: Minimum-Norm, Pseudoinverse

What is the minimum 2-norm solution to  $A\mathbf{x} \cong \mathbf{b}$  and why?

Observe  $\|\mathbf{x}\|_2 = \|\mathbf{y}\|_2$ , and recall that  $\|\mathbf{y}\|_2$  was already minimal. (why?)

$$\mathbf{x} = V\Sigma^+U^T\mathbf{b}$$

solves the minimum-norm least-squares problem.

Generalize the pseudoinverse to the case of a rank-deficient matrix.

Define  $A^+ = V\Sigma^+U^T$  and call it the **pseudoinverse** of  $A$ .

Coincides with prior definition in case of full rank.

## Comparing the Methods

Methods to solve least squares with  $A$  an  $m \times n$  matrix:

- ▶ Form:  $A^T A$ :  $n^2 m / 2$  (symmetric—only need to fill half)  
Solve with  $A^T A$ :  $n^3 / 6$  (Cholesky)
- ▶ Solve with Householder:  $mn^2 - n^3 / 3$
- ▶ If  $m \approx n$ , about the same
- ▶ If  $m \gg n$ : Householder QR requires about twice as much work as normal equations
- ▶ SVD:  $mn^2 + n^3$  (with a large constant)

[Demo: Relative cost of matrix factorizations \[cleared\]](#)

## In-Class Activity: Householder, Givens, SVD

In-class activity: Householder, Givens, SVD

# Outline

Introduction to Scientific Computing

Systems of Linear Equations

Linear Least Squares

**Eigenvalue Problems**

Properties and Transformations

Sensitivity

Computing Eigenvalues

Krylov Space Methods

Nonlinear Equations

Optimization

Interpolation

Numerical Integration and Differentiation

Initial Value Problems for ODEs

Boundary Value Problems for ODEs

Partial Differential Equations and Sparse Linear Algebra

Fast Fourier Transform

Additional Topics

## Eigenvalue Problems: Setup/Math Recap

$A$  is an  $n \times n$  matrix.

- ▶  $\mathbf{x} \neq 0$  is called an *eigenvector* of  $A$  if there exists a  $\lambda$  so that

$$A\mathbf{x} = \lambda\mathbf{x}.$$

- ▶ In that case,  $\lambda$  is called an *eigenvalue*.
- ▶ The set of all eigenvalues  $\lambda(A)$  is called the *spectrum*.
- ▶ The *spectral radius* is the magnitude of the biggest eigenvalue:

$$\rho(A) = \max \{|\lambda| : \lambda \in \lambda(A)\}$$

## Eigenvalue Problems: Motivation from Mechanics

Consider mass-spring systems, e.g. as modeled in (e.g.) [myphysicslab.com](http://myphysicslab.com)

What is needed to model?

- ▶ Hooke's law:  $F = \alpha(x - x_0)$ : spring force prop. to extension
- ▶  $\mathbf{x}(t)$ : vector of mass position
- ▶ Model Hooke's law as  $\mathbf{F}(t) = \mathbf{A}\mathbf{x}(t)$ , i.e. spring constants and position differences computed via matrix

$$\mathbf{F} = m\mathbf{a} \Leftrightarrow \mathbf{A}\mathbf{x}(t) = m \frac{\partial^2 \mathbf{x}(t)}{\partial t^2}.$$

Observe these systems 'wobble': assume  $\mathbf{x}(t) = \sin(\omega t)\mathbf{x}_0$ . Then

$$\mathbf{A}\mathbf{x}_0(t) \sin(\omega t) = m(-\omega^2)\mathbf{x}_0 \sin(\omega t).$$

Cancel  $\sin(\omega t)$ : yields an eigenvalue problem for  $\mathbf{x}_0$  and  $-\omega^2$ .

## Finding Eigenvalues

How do you find eigenvalues?

$$\begin{aligned} A\mathbf{x} = \lambda\mathbf{x} &\Leftrightarrow (A - \lambda I)\mathbf{x} = 0 \\ &\Leftrightarrow A - \lambda I \text{ singular} \Leftrightarrow \det(A - \lambda I) = 0 \end{aligned}$$

$\det(A - \lambda I)$  is called the *characteristic polynomial*, which has degree  $n$ , and therefore  $n$  (potentially complex) roots.

**Does that help algorithmically?** Abel-Ruffini theorem: for  $n \geq 5$  is no general formula for roots of polynomial. IOW: no.

- ▶ For LU and QR, we obtain *exact* answers (except rounding).
- ▶ For eigenvalue problems: not possible—must *iterate*.

[Demo: Rounding in characteristic polynomial using SymPy \[cleared\]](#)



# Multiplicity

What is the *multiplicity* of an eigenvalue?

Actually, there are two notions called multiplicity:

- ▶ *Algebraic Multiplicity*: multiplicity of the root of the characteristic polynomial
- ▶ *Geometric Multiplicity*: #of lin. indep. eigenvectors

In general:  $AM \geq GM$ .

If  $AM > GM$ , the matrix is called *defective*.

## An Example

Give characteristic polynomial, eigenvalues, eigenvectors of

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

CP:  $(\lambda - 1)^2$

Eigenvalues: 1 (with algebraic multiplicity 2)

Eigenvectors:

$$\begin{bmatrix} 1 & 1 \\ & 1 \end{bmatrix} \begin{bmatrix} x \\ y \end{bmatrix} = \begin{bmatrix} x \\ y \end{bmatrix}$$

$\Rightarrow x + y = x \Rightarrow y = 0$ . So only a 1D space of eigenvectors.

## Diagonalizability

When is a matrix called *diagonalizable*?

If it is not defective, i.e. if it has a  $n$  linear independent eigenvectors (i.e. a full basis of them). Call those  $(\mathbf{x}_i)_{i=1}^n$ .

In that case, let

$$X = \begin{bmatrix} | & & | \\ \mathbf{x}_1 & \cdots & \mathbf{x}_n \\ | & & | \end{bmatrix},$$

and observe  $AX = XD$  or

$$A = XDX^{-1},$$

where  $D$  is a diagonal matrix with the eigenvalues.

## Similar Matrices

Related **definition**: Two matrices  $A$  and  $B$  are called similar if there exists an invertible matrix  $X$  so that  $A = XBX^{-1}$ .

In that sense: “Diagonalizable” = “Similar to a diagonal matrix”.

Observe: Similar  $A$  and  $B$  have same eigenvalues. (Why?)

Suppose  $A\mathbf{x} = \lambda\mathbf{x}$ . We have  $B = X^{-1}AX$ . Let  $\mathbf{w} = X^{-1}\mathbf{v}$ . Then

$$B\mathbf{w} = X^{-1}A\mathbf{v} = \lambda\mathbf{w}.$$

## Eigenvalue Transformations (I)

What do the following transformations of the eigenvalue problem  $A\mathbf{x} = \lambda\mathbf{x}$  do?

*Shift.*  $A \rightarrow A - \sigma I$

$$(A - \sigma I)\mathbf{x} = (\lambda - \sigma)\mathbf{x}$$

*Inversion.*  $A \rightarrow A^{-1}$

$$A^{-1}\mathbf{x} = \lambda^{-1}\mathbf{x}$$

*Power.*  $A \rightarrow A^k$

$$A^k\mathbf{x} = \lambda^k\mathbf{x}$$

## Eigenvalue Transformations (II)

*Polynomial*  $A \rightarrow aA^2 + bA + cI$

$$(aA^2 + bA + cI)\mathbf{x} = (a\lambda^2 + b\lambda + c)\mathbf{x}$$

*Similarity*  $T^{-1}AT$  with  $T$  invertible

Let  $\mathbf{y} := T^{-1}\mathbf{x}$ . Then

$$T^{-1}AT\mathbf{y} = \lambda\mathbf{y}$$

## Sensitivity (I)

Assume  $A$  not defective. Suppose  $X^{-1}AX = D$ . Perturb  $A \rightarrow A + E$ .  
What happens to the eigenvalues?

$$X^{-1}(A + E)X = D + F$$

- ▶  $A + E$  and  $D + F$  have same eigenvalues
- ▶  $D + F$  is not necessarily diagonal

Suppose  $\mathbf{v}$  is an eigenvector of  $D + F$ .

$$(D + F)\mathbf{v} = \mu\mathbf{v}$$

$$\Leftrightarrow F\mathbf{v} = (\mu I - D)\mathbf{v}$$

$$\Leftrightarrow (\mu I - D)^{-1}F\mathbf{v} = \mathbf{v} \quad (\text{when is that invertible?})$$

$$\Rightarrow \|\mathbf{v}\| \leq \|(\mu I - D)^{-1}\| \|F\| \|\mathbf{v}\|$$

$$\Rightarrow \|(\mu I - D)^{-1}\|^{-1} \leq \|F\|$$

## Sensitivity (II)

$X^{-1}(A + E)X = D + F$ . Have  $\|(\mu I - D)^{-1}\|^{-1} \leq \|F\|$ .

Demo: Bauer-Fike Eigenvalue Sensitivity Bound [cleared]

$$\|(\mu I - D)^{-1}\|^{-1} = |\mu - \lambda_k|$$

where  $\lambda_k$  is the closest eigenvalue of  $D$  to  $\mu$ . (see demo)

Since we made no assumptions on  $\mu$  other than that it doesn't match an existing eigenvalue, this bound holds for *all* eigenvalues of  $A + E$ .

$$|\mu - \lambda_k| = \|(\mu I - D)^{-1}\|^{-1} \leq \|F\| = \|X^{-1}EX\| \leq \text{cond}(X) \|E\|.$$

- ▶ 'Bad' if  $X$  ill-conditioned, i.e. if eigenvectors are nearly lin.dep.
- ▶  $X$  orthogonal (e.g. for symmetric  $A$ )  $\Rightarrow \text{cond}_2(X) = 1$ .
- ▶ This bound is in terms of *all* eigenvalues, so may overestimate for each individual eigenvalue.



## Power Iteration

### Demo: Motivating Power Iteration [cleared]

Let  $A \in \mathbb{R}^{n \times n}$  and  $A\mathbf{v}_j = \lambda_j \mathbf{v}_j$  ( $j \in \{1, 2, \dots, n\}$ ) and  $|\lambda_1| > |\lambda_2| > \dots > |\lambda_n|$ .

Pick some  $\mathbf{x}_0$ , consider  $\mathbf{x}_{i+1} = A\mathbf{x}_i$  ( $i \in \{0, \dots\}$ ). Called **Power Iteration**.

Let  $\mathbf{x}_0 = \sum_{j=1}^n \alpha_j \mathbf{v}_j$ . Observe that  $\mathbf{x}_i = A^i \mathbf{x}_0 = \sum_{j=1}^n \alpha_j \lambda_j^i \mathbf{v}_j$ .

Define  $\mathbf{e}_i = \mathbf{x}_i / \lambda_1^i - \alpha_1 \mathbf{v}_1$ .

$$\begin{aligned} \|\mathbf{e}_{i+1}\| &= \left\| \frac{\mathbf{x}_{i+1}}{\lambda_1^{i+1}} - \alpha_1 \mathbf{v}_1 \right\| = \left\| \frac{\sum_{j=1}^n \alpha_j \lambda_j^{i+1} \mathbf{v}_j}{\lambda_1^{i+1}} - \alpha_1 \mathbf{v}_1 \right\| \\ &= \left\| \sum_{j=2}^n \alpha_j \left( \frac{\lambda_j}{\lambda_1} \right)^{i+1} \mathbf{v}_j \right\| \leq \left| \frac{\lambda_2}{\lambda_1} \right|^{i+1} \left\| \sum_{j=2}^n \alpha_j \mathbf{v}_j \right\| = \left| \frac{\lambda_2}{\lambda_1} \right|^{i+1} \|\mathbf{e}_0\|. \end{aligned}$$

I.e. converges to (a multiple of)  $\mathbf{v}_1$  'linearly' (see later).

## Power Iteration: Issues?

What could go wrong with Power Iteration?

- ▶ Starting vector has no component along  $\mathbf{x}_1$   
Not a problem in practice: Rounding will introduce one.
- ▶ Overflow in computing  $\lambda_1^i$   
→ Normalize after each step
- ▶  $|\lambda_1| = |\lambda_2|$ 
  - ▶ If  $\lambda_1 = \lambda_2$ : multiplicity, defer.
  - ▶ If  $\lambda_1 \neq \lambda_2$ : use shift+invert to separate magnitudes
- ▶ Complex eigenvalues  
→ use complex-valued shift, and invert.

Recall: for real-valued matrices, eigenvalues come in conjugate pairs  $\lambda$  and  $\lambda^*$ , with  $|\lambda| = |\lambda^*|$ .

## Convergence of Power Iteration: Notation

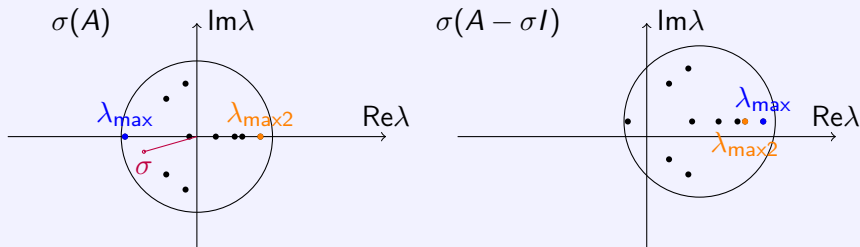
- ▶  $\lambda_{\max}(A)$ : biggest eigenvalue by magnitude
- ▶  $\lambda_{\max 2}(A)$ : second-biggest eigenvalue by magnitude.
- ▶  $\lambda_{\min 2}(A)$ : second-smallest eigenvalue by magnitude
- ▶  $\lambda_{\min}(A)$ : smallest eigenvalue by magnitude

(Not well-defined if there are multiple  $\lambda$  with the same magnitudes.  
Assume that's not the case.)

## Power Iteration: Shift

How does a shift  $(A - \sigma I)$  change power iteration?

- ▶ Converges to eigenvector for  $\lambda_{\max}(A - \sigma I)$  with convergence factor  $\left| \frac{\lambda_{\max 2}(A - \sigma I)}{\lambda_{\max}(A - \sigma I)} \right|$ .
- ▶ Can help guide convergence to eigenvalues 'on boundary' of spectrum.



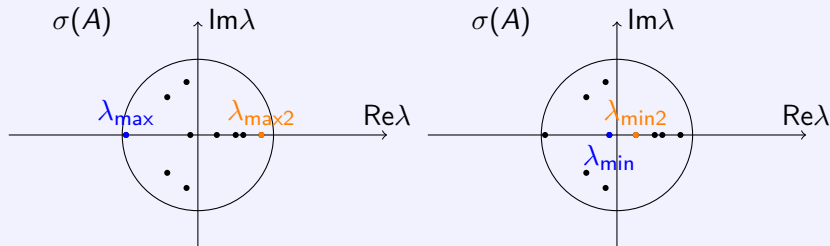
## Power Iteration: Inversion

How does inversion ( $A^{-1}$ ) change power iteration?

- ▶ Converges to eigenvector for  $\lambda_{\max}(A^{-1}) = 1/\lambda_{\min}(A)$  with convergence factor

$$\left| \frac{\lambda_{\max 2}(A^{-1})}{\lambda_{\max}(A^{-1})} \right| = \left| \frac{1/\lambda_{\min 2}(A)}{1/\lambda_{\min}(A)} \right| = \left| \frac{\lambda_{\min}(A)}{\lambda_{\min 2}(A)} \right|.$$

- ▶ Guide convergence to smallest eigenvalues.



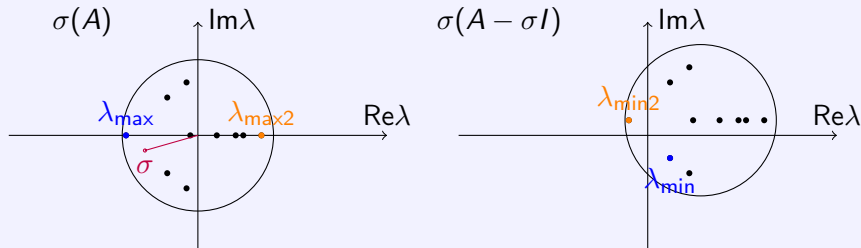
## Power Iteration: Shift and Inversion

How does shift-invert  $((A - \sigma I)^{-1})$  change power iteration?

- ▶ Converges to eigenvector for  $\lambda_{\max}((A - \sigma I)^{-1}) = 1/\lambda_{\min}(A - \sigma I)$  with convergence factor

$$\left| \frac{\lambda_{\max 2}((A - \sigma I)^{-1})}{\lambda_{\max}((A - \sigma I)^{-1})} \right| = \left| \frac{\lambda_{\min}(A - \sigma I)}{\lambda_{\min 2}(A - \sigma I)} \right|.$$

- ▶ Guide convergence to eigenvalue **closest to  $\sigma$** .



## What about Eigenvalues?

Power Iteration generates eigenvectors. What if we would like to know eigenvalues?

Estimate them:

$$\frac{\mathbf{x}^T A \mathbf{x}}{\mathbf{x}^T \mathbf{x}}$$

- ▶  $= \lambda$  if  $\mathbf{x}$  is an eigenvector w/ eigenvalue  $\lambda$
- ▶ Otherwise, an estimate of a 'nearby' eigenvalue

This is called the *Rayleigh quotient*.

Idea: Could use Rayleigh quotient as a shift. Yields **Rayleigh quotient iteration**.

[Demo: Power Iteration and its Variants \[cleared\]](#)

## In-Class Activity: Eigenvalues

In-class activity: Eigenvalues



## Schur form: Motivation

For finding multiple eigenvalues, want factorization that allows access to **all** eigenvalues and eigenvectors.

Suggestions?

- ▶ Diagonalization  $A = XDX^{-1}$  cannot provide what we need: it does not always exist.
- ▶ Even if it did exist, computing/applying  $X^{-1}$  would be subject to rounding concerns.
  - ▶ Idea: use a similarity transform with *orthogonal* matrices.

## Schur form

Show: Every matrix is orthonormally similar to an upper triangular matrix, i.e.  $A = QUQ^T$ . This is called the **Schur form** or **Schur factorization**.

Let  $(\lambda, \mathbf{v})$  be an eigenpair (at least one always exists), i.e.  $A\mathbf{v} = \lambda\mathbf{v}$  ( $\mathbf{v} \neq 0$ ). Let  $V = \text{span}\{\mathbf{v}\}$ . Then

$$\begin{aligned} A: \quad V &\rightarrow V \\ V^\perp &\rightarrow V \oplus V^\perp \end{aligned}$$

$$A = \underbrace{\begin{bmatrix} | & & & & \\ \mathbf{v} & & & & \\ | & & & & \\ | & & & & \end{bmatrix}}_{Q_1} \begin{bmatrix} \lambda & * & * & * & * \\ 0 & * & * & * & * \\ \vdots & * & * & * & * \\ 0 & * & * & * & * \end{bmatrix} Q_1^T.$$

$U_1$

- ▶ Bottom right of  $U_1$ : Same eigenvalues as  $A$  without  $\lambda$ .
- ▶ Repeat  $n$  times to triangular:  $A = Q_n \cdots Q_1 U Q_1^T \cdots Q_n^T$ .

## Schur Form: Comments, Eigenvalues, Eigenvectors

$A = QUQ^T$ . For complex  $\lambda$ :

- ▶ Either complex matrices, or
- ▶  $2 \times 2$  blocks on diag.

If we had a Schur form of  $A$  (no  $2 \times 2$  blocks), can we find the eigenvalues?

The eigenvalues (of  $U$  and  $A$ !) are on the diagonal of  $U$ .

And the eigenvectors?

Find eigenvector of  $U$ : Suppose  $\lambda$  is an eigenvalue.

$$U - \lambda I = \begin{bmatrix} U_{11} & \mathbf{u} & U_{13} \\ 0 & 0 & \mathbf{v}^T \\ 0 & 0 & U_{31} \end{bmatrix}$$

$\mathbf{x} = [U_{11}^{-1}\mathbf{u}; -1; 0]^T$  eigenvector of  $U$ , and  $Q\mathbf{x}$  eigenvector of  $A$ .

## Computing Multiple Eigenvalues

All Power Iteration Methods compute one eigenvalue at a time.  
What if I want *all* eigenvalues?

Two ideas:

1. *Deflation*: similarity transform to

$$\begin{bmatrix} \lambda_1 & * \\ & B \end{bmatrix},$$

i.e. use the argument for the existence of Schur form as a computational procedure.

2. Iterate with multiple vectors simultaneously.

## Simultaneous Iteration

What happens if we carry out power iteration on multiple vectors simultaneously?

*Simultaneous Iteration:*

1. Start with  $X_0 \in \mathbb{R}^{n \times p}$  ( $p \leq n$ ) with (arbitrary) iteration vectors in columns
2.  $X_{k+1} = AX_k$

Problems:

- ▶ Needs rescaling
- ▶  $X$  increasingly ill-conditioned: all columns go towards  $\mathbf{x}_1$

Fix: orthogonalize!

## Orthogonal Iteration

*Orthogonal Iteration:*

1. Start with  $X_0 \in \mathbb{R}^{n \times p}$  ( $p \leq n$ ) with (arbitrary) iteration vectors in columns
2.  $\hat{Q}_k R_k = X_k$  (reduced)
3.  $X_{k+1} = A \hat{Q}_k$

**Good:**  $X_k$  obey convergence theory from power iteration

**Bad:**

- ▶ Slow/linear convergence
- ▶ Expensive iteration

## Toward the QR Algorithm

$$Q_0 R_0 = X_0$$

$$X_1 = A Q_0$$

$$Q_1 R_1 = X_1 = A Q_0 \Rightarrow Q_1 R_1 Q_0^T = A$$

$$X_2 = A Q_1$$

$$Q_2 R_2 = X_2 = A Q_1 \Rightarrow Q_2 R_2 Q_1^T = A$$

Once the  $Q_k$  converge ( $Q_{n+1} \approx Q_n$ ), we have a Schur factorization!

**Problem:**  $Q_{n+1} \approx Q_n$  works poorly as a convergence test.

**Observation 1:** Once  $Q_{n+1} \approx Q_n$ , we also have  $Q_n R_n Q_n^T \approx A$ .

**Observation 2:**  $\hat{X}_n := Q_n^T A Q_n \approx R_n$ .

**Idea:** Use below-diag part of  $\hat{X}_n$  for convergence check.

**Q:** Can we restate the iteration to compute  $\hat{X}_k$  directly?

Demo: Orthogonal Iteration [cleared]

## QR Iteration/QR Algorithm

Orthogonal iteration:      QR iteration:

$$X_0 = A$$

$$Q_k R_k = X_k$$

$$X_{k+1} = A Q_k$$

$$\bar{X}_0 = A$$

$$\bar{Q}_k \bar{R}_k = \bar{X}_k$$

$$\bar{X}_{k+1} = \bar{R}_k \bar{Q}_k$$

$$\blacktriangleright \bar{X}_{k+1} = \bar{R}_k \bar{Q}_k = \bar{Q}_k^H \bar{X}_k \bar{Q}_k = \underbrace{\bar{Q}_k^H \bar{Q}_{k-1}^H \cdots \bar{Q}_0^H}_{\text{orthogonal}} A \underbrace{\bar{Q}_0 \cdots \bar{Q}_k}_{\text{orthogonal}}$$

$\blacktriangleright$  The  $\bar{X}_k$  are all similar to  $A \rightarrow$  have same eigenvalues.

$\blacktriangleright A^2 = \bar{Q}_0 \bar{R}_0 \bar{Q}_0 \bar{R}_0 = \bar{Q}_0 \bar{Q}_1 \bar{R}_1 \bar{R}_0$  (analogous for  $A^k$ )

**Claim:** (see next slide) Orth.it. and QR it. are equivalent, via

$$\blacktriangleright Q_k = \bar{Q}_0 \bar{Q}_1 \cdots \bar{Q}_k.$$

$$\blacktriangleright \hat{X}_k = \bar{X}_{k+1}.$$

From orthogonal iteration: Observed  $\hat{X}_k = \bar{X}_{k+1}$  converge.

$\rightarrow$  QR iteration produces Schur form.



## Proof sketch: Equivalence of QR iteration/Orth. iteration

### Orthogonal Iteration (no bars)

- ▶  $X_0 := A$ 
  - ▶  $Q_0 R_0 := X_0$ ,
  - ▶ where we may choose  
 $Q_0 = \bar{Q}_0$
  - ▶  $\hat{X}_0 = Q_0^H A Q_0 =$   
 $Q_0^H Q_0 R_0 Q_0 = R_0 Q_0$
- ▶  $X_1 := A Q_0$ 
  - ▶  $Q_1 R_1 := X_1$ ,
  - and because of  
 $X_1 = Q_0 Q_0^H A Q_0 = Q_0 \hat{X}_0 =$   
 $Q_0 \bar{Q}_1 \bar{R}_1$   
we may choose  
 $Q_1 = Q_0 \bar{Q}_1 = \bar{Q}_0 \bar{Q}_1$ .
- ▶  $\vdots$

### QR Iteration (with bars)

- ▶  $\bar{X}_0 := A$ 
  - ▶  $\bar{Q}_0 \bar{R}_0 := A$
- ▶  $\bar{X}_1 := \bar{R}_0 \bar{Q}_0 = \hat{X}_0$ 
  - ▶  $\bar{Q}_1 \bar{R}_1 := \bar{X}_1$
- ▶  $\bar{X}_2 := \bar{R}_1 \bar{Q}_1$ 
  - ▶  $\bar{X}_2 = Q_1^H A Q_1 = \hat{X}_1$
- ▶  $\vdots$

Demo: QR Iteration [cleared]

## QR Iteration: Forward *and* Inverse

QR iteration may be viewed as performing **inverse iteration**. How?

Take an inverse (conjugate) transpose of the whole method.

▶  $\bar{X}_0^{-H} = A^{-H}$ .

▶ Recall  $\bar{Q}_k \bar{R}_k = \bar{X}_k$ . Invert and transpose both sides:

$$Q_k R_k^{-H} = \bar{X}_k^{-H}$$

▶ Recall  $\bar{X}_{k+1} = \bar{R}_k \bar{Q}_k$ . Invert and transpose both sides:

$$\bar{X}_{k+1}^{-H} = \bar{R}_k^{-H} \bar{Q}_k$$

I.e. *exact same iterates* as QR iteration (power iteration ‘from the left’) would be produced by “QL iteration” on  $A^{-H}$ , i.e. inverse iteration ‘from the right’. Therefore: would expect *shifts* to be effective.

## QR Iteration: Incorporating a Shift

How can we accelerate convergence of QR iteration using shifts?

$$\begin{aligned}\bar{Q}_k \bar{R}_k &= \bar{X}_k - \sigma_k I \\ \bar{X}_{k+1} &= \bar{R}_k \bar{Q}_k + \sigma_k I\end{aligned}$$

Still a similarity transform:

$$\bar{X}_{k+1} = \bar{R}_k \bar{Q}_k + \sigma_k I = [\bar{Q}_k^T \bar{X}_k - \bar{Q}_k^T \sigma_k] \bar{Q}_k + \sigma_k I$$

**Q:** How should the shifts be chosen?

- ▶ Ideally: Close to existing eigenvalue
- ▶ Heuristics:
  - ▶ Lower right entry of  $\bar{X}_k$
  - ▶ Eigenvalues of lower right  $2 \times 2$  of  $\bar{X}_k$

**Demo:** [QR Iteration \[cleared\]](#) (Shifted)

## QR Iteration: Computational Expense

A full QR factorization at each iteration costs  $O(n^3)$ —can we make that cheaper?

**Idea:** *Upper Hessenberg form*

$$A = Q \begin{bmatrix} * & * & * & * \\ * & * & * & * \\ & * & * & * \\ & & * & * \end{bmatrix} Q^T$$

- ▶ Attainable by *similarity transforms* (!)  $HAH^T$  with Householders that start 1 entry lower than ‘usual’
- ▶ QR factorization of Hessenberg matrices can be achieved in  $O(n^2)$  time using Givens rotations.

[Demo: Householder Similarity Transforms \[cleared\]](#)

## QR/Hessenberg: Overall procedure

Overall procedure:

1. Reduce matrix to Hessenberg form
2. Apply QR iteration using Givens QR to obtain Schur form

Why does QR iteration *stay* in Hessenberg form?

Assume  $\bar{X}_k$  is upper Hessenberg ("UH").

- ▶  $\bar{Q}_k \bar{R}_k = \bar{X}_k$ :  $\bar{Q}_k = \bar{R}_k^{-1} \bar{X}_k$  is UH (UH  $\cdot$  upper  $\Delta$  = UH)
- ▶  $\bar{X}_{k+1} = \bar{R}_k \bar{Q}_k$  is UH (upper  $\Delta \cdot$  UH = UH)

What does this process look like for symmetric matrices?

- ▶ Use Householders to attain tridiagonal form
- ▶ Use QR iteration with Givens to attain diagonal form

## Krylov space methods: Intro

What subspaces can we use to look for eigenvectors?

QR:

$$\text{span}\{A^l \mathbf{y}_1, A^l \mathbf{y}_2, \dots, A^l \mathbf{y}_k\}$$

Krylov space:

$$\text{span}\{\underbrace{\mathbf{x}}_{\mathbf{x}_0}, A\mathbf{x}, \dots, \underbrace{A^{k-1}\mathbf{x}}_{\mathbf{x}_{k-1}}\}$$

Define:

$$K_k := \begin{bmatrix} | & & | \\ \mathbf{x}_0 & \cdots & \mathbf{x}_{k-1} \\ | & & | \end{bmatrix}. \quad (n \times k)$$

## Krylov for Matrix Factorization

What matrix factorization is obtained through Krylov space methods?

$$AK_n = \begin{bmatrix} | & & | \\ \mathbf{x}_1 & \cdots & \mathbf{x}_n \\ | & & | \end{bmatrix} = K_n \underbrace{\begin{bmatrix} | & & | & | \\ \mathbf{e}_2 & \cdots & \mathbf{e}_n & K_n^{-1}\mathbf{x}_n \\ | & & | & | \end{bmatrix}}_{C_n}.$$

- ▶  $K_n^{-1}AK_n = C_n$
- ▶  $C_n$  is upper Hessenberg
- ▶ So Krylov is 'just' another way to get a matrix into upper Hessenberg form.
- ▶ But: works well when only matvec is available (searches in Krylov space, not the space spanned by first columns)

## Conditioning in Krylov Space Methods/Arnoldi Iteration (I)

What is a problem with Krylov space methods? How can we fix it?

$(\mathbf{x}_i)$  converge rapidly to eigenvector for largest eigenvalue  
→  $K_k$  become ill-conditioned

**Idea:** Orthogonalize! (at end... for now)

$$Q_n R_n = K_n \quad \Rightarrow \quad Q_n = K_n R_n^{-1}$$

Then

$$Q_n^T A Q_n = R_n \underbrace{K_n^{-1} A K_n}_{C_n} R_n^{-1}.$$

- ▶  $C_n$  is upper Hessenberg
- ▶  $Q_n^T A Q_n$  is also UH  
(because upper  $\Delta \cdot$  UH = UH and UH  $\cdot$  upper  $\Delta$  = UH).



## Conditioning in Krylov Space Methods/Arnoldi Iteration (II)

We find that  $Q_n^T A Q_n$  is also upper Hessenberg:  $Q_n^T A_n Q_n = H$ .  
Also readable as  $A Q_n = Q_n H$ , which, read column-by-column, is:

$$A \mathbf{q}_k = h_{1k} \mathbf{q}_1 + \cdots + h_{k+1,k} \mathbf{q}_{k+1}$$

We find:  $h_{jk} = \mathbf{q}_j^T A \mathbf{q}_k$ . Use that to rewrite, letting  $\mathbf{v} = A \mathbf{q}_k$ :

$$\mathbf{v} - (\mathbf{q}_1^T \mathbf{v}) \mathbf{q}_1 - \cdots - (\mathbf{q}_k^T \mathbf{v}) \mathbf{q}_k = h_{k+1,k} \mathbf{q}_{k+1}$$

- ▶ Looks just like Gram-Schmidt QR!
- ▶ Can compute  $(k+1)$ st column of  $H$  and  $\mathbf{q}_{k+1}$  from  $\mathbf{q}_1, \dots, \mathbf{q}_k$ .

This is called **Arnoldi iteration**. For symmetric: **Lanczos iteration**.

[Demo: Arnoldi Iteration \[cleared\]](#) (Part 1)

## Krylov: What about eigenvalues?

How can we use Arnoldi/Lanczos to compute eigenvalues?

$$Q = [Q_k \quad U_k]$$

Green: known (i.e. already computed), red: not yet computed.

$$H = Q^T A Q = \begin{bmatrix} Q_k^T \\ U_k^T \end{bmatrix} A [Q_k \quad U_k] = \begin{bmatrix} * & * & * & * & * \\ * & * & * & * & * \\ & * & * & * & * \\ & & * & * & * \\ & & & * & * \end{bmatrix}$$

Use eigenvalues of top-left matrix as approximate eigenvalues.  
(still need to be computed, using QR it.)

Those are called **Ritz values**.

## Computing the SVD (Kiddy Version)

1. Compute (orth.) eigenvectors  $V$  and eigenvalues  $D$  of  $A^T A$ ,

$$A^T A V = V D \quad \Rightarrow \quad V^T A^T A V = D =: \Sigma^2.$$

2. Find  $U$  from  $A = U \Sigma V^T \Leftrightarrow U \Sigma = A V$ .

Observe  $U$  is orthogonal if  $\Sigma^{-1}$  exists: (If not, can choose so.)

$$U^T U = \Sigma^{-1} V^T A^T A V \Sigma^{-1} = \Sigma^{-1} \Sigma^2 \Sigma^{-1} = I.$$

### Demo: Computing the SVD [cleared]

“Actual”/“non-kiddy” computation of the SVD:

- ▶ Bidiagonalize  $A = U \begin{bmatrix} B \\ 0 \end{bmatrix} V^T$ , then diagonalize via variant of QR.
- ▶ References: [Chan '82](#) or Golub/van Loan Sec 8.6.

# Outline

Introduction to Scientific Computing

Systems of Linear Equations

Linear Least Squares

Eigenvalue Problems

**Nonlinear Equations**

Introduction

Iterative Procedures

Methods in One Dimension

Methods in  $n$  Dimensions ("Systems of Equations")

Optimization

Interpolation

Numerical Integration and Differentiation

Initial Value Problems for ODEs

Boundary Value Problems for ODEs

Partial Differential Equations and Sparse Linear Algebra

Fast Fourier Transform

Additional Topics

## Solving Nonlinear Equations

What is the goal here?

Solve  $\mathbf{f}(\mathbf{x}) = 0$  for  $f : \mathbb{R}^n \rightarrow \mathbb{R}^n$ .

If looking for solution to  $\tilde{\mathbf{f}}(\mathbf{x}) = \mathbf{y}$ , simply consider  $\mathbf{f}(\mathbf{x}) = \tilde{\mathbf{f}}(\mathbf{x}) - \mathbf{y}$ .

*Intuition:* Each of the  $n$  equations describes a surface. Looking for intersections.

[Demo: Three quadratic functions](#) [cleared]

## Showing Existence

How can we show existence of a root?

- ▶ Intermediate value theorem (uses continuity, 1D only)
- ▶ Inverse function theorem (relies on invertible Jacobian  $J_f$ )  
Get *local* invertibility, i.e.  $f(\mathbf{x}) = \mathbf{y}$  solvable

- ▶ **Contraction mapping theorem**

A function  $\mathbf{g} : \mathbb{R}^n \rightarrow \mathbb{R}^n$  is called *contractive* if there exists a  $0 < \gamma < 1$  so that  $\|\mathbf{g}(\mathbf{x}) - \mathbf{g}(\mathbf{y})\| \leq \gamma \|\mathbf{x} - \mathbf{y}\|$ . A fixed point of  $\mathbf{g}$  is a point where  $\mathbf{g}(\mathbf{x}) = \mathbf{x}$ .

*Then:* On a closed set  $S \subseteq \mathbb{R}^n$  with  $\mathbf{g}(S) \subseteq S$  there exists a unique fixed point.

*Example:* (real-world) map

In general, *no uniqueness* results available.

## Sensitivity and Multiplicity

What is the sensitivity/conditioning of root finding?

$\text{cond}(\text{root finding}) = \text{cond}(\text{evaluation of the inverse function at } 0)$   
Evaluation (of the inverse) at 0: must use *absolute* condition numbers.

What are multiple roots?

$$0 = f(x) = f'(x) = \dots = f^{(m-1)}(x)$$

This is called a **root of multiplicity  $m$** .

How do multiple roots interact with conditioning?

The inverse function is steep near one, so conditioning is poor.

## In-Class Activity: Krylov and Nonlinear Equations

In-class activity: Krylov and Nonlinear Equations



## Rates of Convergence

What is *linear convergence*? *quadratic convergence*?

$\mathbf{e}_k = \hat{\mathbf{u}}_k - \mathbf{u}$ : error in the  $k$ th iterate  $\hat{\mathbf{u}}_k$ . Assume  $\mathbf{e}_k \rightarrow 0$  as  $k \rightarrow \infty$ .

An iterative method *converges with rate*  $r$  if

$$\lim_{k \rightarrow \infty} \frac{\|\mathbf{e}_{k+1}\|}{\|\mathbf{e}_k\|^r} = C \begin{cases} > 0, \\ < \infty. \end{cases}$$

- ▶  $r = 1$  is called *linear convergence*.
- ▶  $r > 1$  is called *superlinear convergence*.
- ▶  $r = 2$  is called *quadratic convergence*.

Examples:

- ▶ Power iteration is linearly convergent.
- ▶ Rayleigh quotient iteration is quadratically convergent.

## About Convergence Rates

### Demo: Rates of Convergence [cleared]

Characterize linear, quadratic convergence in terms of the 'number of accurate digits'.

- ▶ Linear convergence gains a constant number of digits each step:

$$\|\mathbf{e}_{k+1}\| \leq C \|\mathbf{e}_k\|$$

(and  $C < 1$  matters!)

- ▶ Quadratic convergence

$$\|\mathbf{e}_{k+1}\| \leq C \|\mathbf{e}_k\|^2$$

(Only starts making sense once  $\|\mathbf{e}_k\|$  is small.  $C$  doesn't matter much.)

## Stopping Criteria

Comment on the 'foolproof-ness' of these stopping criteria:

1.  $|f(x)| < \varepsilon$  ('residual is small')
2.  $\|\mathbf{x}_{k+1} - \mathbf{x}_k\| < \varepsilon$
3.  $\|\mathbf{x}_{k+1} - \mathbf{x}_k\| / \|\mathbf{x}_k\| < \varepsilon$

1. Can trigger far away from a root in the case of multiple roots (or a 'flat'  $f$ )
2. Allows different 'relative accuracy' in the root depending on its magnitude.
3. Enforces a relative accuracy in the root, but *does not* actually check that the function value is small.  
So if convergence 'stalls' away from a root, this may trigger without being anywhere near the desired solution.

**Lesson:** No stopping criterion is bulletproof. The 'right' one almost always depends on the application.

# Bisection Method

Demo: Bisection Method [cleared]

What's the rate of convergence? What's the constant?

Linear with constant  $1/2$ .

## Mini Review: Taylor's Theorem

$$f(x+h) = f(x) + f'(x)h + \frac{f''(x)}{2!}h^2 + \dots \quad (\text{requires } f \text{ analytic})$$

Taylor **with explicit remainder term** ( $\theta \in [x, x+h]$ ,  $f \in C^{k+1}$ ):

$$f(x+h) = f(x) + \dots + \frac{f^{(k-1)}(x)}{(k-1)!}h^{k-1} + \frac{f^k(\theta)}{k!}h^k$$

Special case  $k = 1$ : Equivalent to **mean value theorem**:

$$f(x+h) = f(x) + f'(\theta)h \quad (\theta \in [x, x+h], f \in C^1)$$

With **big-O truncation**:

$$f(x+h) = f(x) + \dots + \frac{f^{(k-1)}(x)}{(k-1)!}h^{k-1} + O(h^k) \quad (h \rightarrow 0, f \in C^{k+1})$$

## Fixed Point Iteration

$$\begin{aligned}x_0 &= \langle \text{starting guess} \rangle \\x_{k+1} &= g(x_k)\end{aligned}$$

### Demo: Fixed point iteration [cleared]

When does fixed point iteration converge? Assume  $g$  is smooth.

Let  $x^*$  be the fixed point with  $x^* = g(x^*)$ .

**Claim:** If  $|g'(x^*)| < 1$  at the fixed point, FPI converges.

Error:

$$e_{k+1} = x_{k+1} - x^* = g(x_k) - g(x^*)$$

[cont'd.]

## Fixed Point Iteration: Convergence cont'd.

Error in FPI:  $e_{k+1} = x_{k+1} - x^* = g(x_k) - g(x^*)$

Mean value theorem says: There is a  $\theta_k$  between  $x_k$  and  $x^*$  so that

$$g(x_k) - g(x^*) = g'(\theta_k)(x_k - x^*) = g'(\theta_k)e_k.$$

So:  $e_{k+1} = g'(\theta_k)e_k$  and if  $\|g'\| \leq C < 1$  near  $x^*$ , then we have linear convergence with constant  $C$ .

**Q:** What if  $g'(x^*) = 0$ ?

By Taylor:

$$g(x_k) - g(x^*) = g''(\xi_k)(x_k - x^*)^2/2$$

So we have *quadratic convergence* in this case!

We would obviously like a systematic way of finding  $g$  that produces quadratic convergence.

# Newton's Method

Derive Newton's method.

**Idea:** Approximate  $f$  at  $x_k$  using Taylor:  $f(x_k + h) \approx f(x_k) + f'(x_k)h$   
Now find root of this linear approximation in terms of  $h$ :

$$f(x_k) + f'(x_k)h = 0 \quad \Leftrightarrow \quad h = -\frac{f(x_k)}{f'(x_k)}.$$

$$x_0 = \langle \text{starting guess} \rangle$$

$$x_{k+1} = x_k - \frac{f(x_k)}{f'(x_k)} = g(x_k)$$

[Demo: Newton's method \[cleared\]](#)



## Convergence and Properties of Newton

What's the rate of convergence of Newton's method?

Recall the quotient rule  $(f/g)' = (f'g - g'f)/g^2$ .

$$g'(x) = \frac{f(x)f''(x)}{f'(x)^2}$$

So if  $f(x^*) = 0$  and  $f'(x^*) \neq 0$ , we have  $g'(x^*) = 0$ , i.e. quadratic convergence toward single roots.

*Drawbacks of Newton?*

- ▶ Convergence (argument) only *locally*
- ▶ Have to know  $f'$ !

**Demo:** [Convergence of Newton's Method](#) [cleared]

## Secant Method

What would Newton without the use of the derivative look like?

Approximate

$$f'(x_k) \approx \frac{f(x_k) - f(x_{k-1})}{x_k - x_{k-1}}.$$

So

$$\begin{aligned}x_0 &= \langle \text{starting guess} \rangle \\x_{k+1} &= x_k - \frac{f(x_k)}{\frac{f(x_k) - f(x_{k-1})}{x_k - x_{k-1}}}.\end{aligned}$$

# Convergence of Properties of Secant

Rate of convergence is  $(1 + \sqrt{5}) / 2 \approx 1.618$ . ([proof](#))

*Drawbacks* of Secant?

- ▶ Convergence argument only good *locally*  
Will see: convergence only local (near root)
- ▶ Slower convergence
- ▶ Need two starting guesses

[Demo: Secant Method \[cleared\]](#)

[Demo: Convergence of the Secant Method \[cleared\]](#)

Secant (and similar methods) are called **Quasi-Newton Methods**.

## Improving on Newton?

How would we do “Newton + 1” (i.e. even faster, even better)?

Easy:

- ▶ Use second order Taylor approximation, solve resulting quadratic
- ▶ Get cubic convergence!
- ▶ Get a method that's *extremely* fast and *extremely* brittle
- ▶ Need **second** derivatives
- ▶ What if the quadratic has no solution?

## Root Finding with Interpolants

Secant method uses a linear approximation to  $f$  based on points  $f(x_k)$ ,  $f(x_{k-1})$ , could use more points and higher-order approximation:

- ▶ Can fit polynomial to (subset of)  $(x_0, f(x_0)), \dots, (x_k, f(x_k))$
- ▶ Look for a root of that
- ▶ Fit a quadratic to the last three: **Muller's method**
  - ▶ Also finds complex roots
  - ▶ Convergence rate  $r \approx 1.84$

What about existence of roots in that case?

- ▶ **Inverse quadratic interpolation**
  - ▶ Fit/'interpolate' quadratic polynomial  $q$  so that  $q(f(x_i)) = x_i$  for  $i \in \{k, k-1, k-2\}$ .
  - ▶ Approximate root by evaluating  $x_{k+1} = q(0)$ .

## Achieving Global Convergence

The linear approximations in Newton and Secant are only good locally. How could we use that?

- ▶ Hybrid methods: bisection + Newton
  - ▶ Stop if Newton leaves bracket
- ▶ Fix a region where they're 'trustworthy' (**trust region methods**)
- ▶ Limit step size
- ▶ Sufficient conditions for convergence of Newton (under *strong* assumptions) are available.

## In-Class Activity: Nonlinear Equations

In-class activity: Nonlinear Equations

## Fixed Point Iteration

$$\begin{aligned}\mathbf{x}_0 &= \langle \text{starting guess} \rangle \\ \mathbf{x}_{k+1} &= \mathbf{g}(\mathbf{x}_k)\end{aligned}$$

When does this converge?

Converges (locally) if  $\|J_{\mathbf{g}}(\mathbf{x}^*)\| < 1$  in some norm, where the *Jacobian matrix*

$$J_{\mathbf{g}}(\mathbf{x}^*) = \begin{bmatrix} \partial_{x_1} g_1 & \cdots & \partial_{x_n} g_1 \\ \vdots & & \\ \partial_{x_1} g_n & \cdots & \partial_{x_n} g_n \end{bmatrix}.$$

Similarly: If  $J_{\mathbf{g}}(\mathbf{x}^*) = 0$ , we have at least quadratic convergence.

**Better:** There exists a norm  $\|\cdot\|_A$  such that  $\rho(A) \leq \|A\|_A < \rho(A) + \epsilon$ , so  $\rho(A) < 1$  suffices. ([proof](#))



## Newton's Method

What does Newton's method look like in  $n$  dimensions?

Approximate by linear:  $\mathbf{f}(\mathbf{x} + \mathbf{s}) = \mathbf{f}(\mathbf{x}) + J_{\mathbf{f}}(\mathbf{x})\mathbf{s}$ .

Set to 0:  $J_{\mathbf{f}}(\mathbf{x})\mathbf{s} = -\mathbf{f}(\mathbf{x}) \Rightarrow \mathbf{s} = -(J_{\mathbf{f}}(\mathbf{x}))^{-1}\mathbf{f}(\mathbf{x})$ .

$\mathbf{x}_0 = \langle \text{starting guess} \rangle$

$\mathbf{x}_{k+1} = \mathbf{x}_k - (J_{\mathbf{f}}(\mathbf{x}_k))^{-1}\mathbf{f}(\mathbf{x}_k)$

Downsides of  $n$ -dim. Newton?

- ▶ Still only locally convergent
- ▶ Computing and inverting  $J_{\mathbf{f}}$  is expensive.

[Demo: Newton's method in  \$n\$  dimensions \[cleared\]](#)

## Secant in $n$ dimensions?

What would the secant method look like in  $n$  dimensions?

Need an 'approximate Jacobian' satisfying

$$\tilde{J} \cdot (\mathbf{x}_{k+1} - \mathbf{x}_k) = \mathbf{f}(\mathbf{x}_{k+1}) - \mathbf{f}(\mathbf{x}_k).$$

Suppose we have *already taken* a step to  $\mathbf{x}_{k+1}$ . Could we 'reverse engineer'  $\tilde{J}$  from that equation?

- ▶ No:  $n^2$  unknowns in  $\tilde{J}$ , but only  $n$  equations
- ▶ Can only hope to 'update'  $\tilde{J}$  with information from current guess.

Some choices, all called *Broyden's method*:

- ▶ update  $J_n$ , minimize  $\|J_n - J_{n-1}\|_F$
- ▶ update  $J_n^{-1}$  (via Sherman-Morrison), minimize  $\|J_n^{-1} - J_{n-1}^{-1}\|_F$   
multiple variants ("good" Broyden and "bad" Broyden)

## Numerically Testing Derivatives

Getting derivatives right is important. How can I test/debug them?

Verify convergence of the Taylor remainder by checking that, for a unit vector  $\mathbf{s}$  and an input vector  $\mathbf{x}$ ,

$$\left\| \frac{\mathbf{f}(\mathbf{x} + h\mathbf{s}) - \mathbf{f}(\mathbf{x})}{h} - J_{\mathbf{f}}(\mathbf{x})\mathbf{s} \right\| = O(h) \quad (h \rightarrow 0).$$

- ▶ Same trick can be used to check the Hessian (needed in optimization): It is the Jacobian of the gradient.
- ▶ Norm is not necessarily small. Convergence (i.e. decrease with  $h$ ) matters.
- ▶ Important to divide by  $h$ , so that the norm is  $O(1)$ .
- ▶ Can “bootstrap” the derivatives: do the above one term at a time.

# Outline

Introduction to Scientific Computing

Systems of Linear Equations

Linear Least Squares

Eigenvalue Problems

Nonlinear Equations

## Optimization

Introduction

Methods for unconstrained opt. in one dimension

Methods for unconstrained opt. in  $n$  dimensions

Nonlinear Least Squares

Constrained Optimization

Interpolation

Numerical Integration and Differentiation

Initial Value Problems for ODEs

Boundary Value Problems for ODEs

Partial Differential Equations and Sparse Linear Algebra

Fast Fourier Transform

Additional Topics

## Optimization: Problem Statement

Have: **Objective function**  $f : \mathbb{R}^n \rightarrow \mathbb{R}$

Want: **Minimizer**  $\mathbf{x}^* \in \mathbb{R}^n$  so that

$$f(\mathbf{x}^*) = \min_{\mathbf{x}} f(\mathbf{x}) \quad \text{subject to} \quad \mathbf{g}(\mathbf{x}) = 0 \quad \text{and} \quad \mathbf{h}(\mathbf{x}) \leq 0.$$

- ▶  $\mathbf{g}(\mathbf{x}) = 0$  and  $\mathbf{h}(\mathbf{x}) \leq 0$  are called **constraints**.  
They define the set of **feasible points**  $\mathbf{x} \in S \subseteq \mathbb{R}^n$ .
- ▶ If  $\mathbf{g}$  or  $\mathbf{h}$  are present, this is **constrained optimization**.  
Otherwise **unconstrained optimization**.
- ▶ If  $f$ ,  $\mathbf{g}$ ,  $\mathbf{h}$  are *linear*, this is called **linear programming**.  
Otherwise **nonlinear programming**.

## Optimization: Observations

Q: What if we are looking for a *maximizer* not a minimizer?

Give some examples:

- ▶ What is the fastest/cheapest/shortest... way to do...?
- ▶ Solve a system of equations 'as well as you can' (if no exact solution exists)—similar to what least squares does for linear systems:

$$\min \|F(\mathbf{x})\|$$

What about multiple objectives?

- ▶ In general: Look up **Pareto optimality**.
- ▶ **For 450:** Make up your mind—decide on one (or build a combined objective). Then we'll talk.

## Existence/Uniqueness

Terminology: **global minimum** / **local minimum**

Under what conditions on  $f$  can we say something about existence/uniqueness?

If  $f : S \rightarrow \mathbb{R}$  is continuous on a closed and bounded set  $S \subseteq \mathbb{R}^n$ , then

a minimum exists.

$f : S \rightarrow \mathbb{R}$  is called *coercive* on  $S \subseteq \mathbb{R}^n$  if

$$\lim_{\|x\| \rightarrow \infty} f(x) = +\infty$$

If  $f$  is coercive and continuous and  $S$  is closed, ...

a global minimum exists (but is possibly non-unique).

## Convexity

$S \subseteq \mathbb{R}^n$  is called **convex** if for all  $\mathbf{x}, \mathbf{y} \in S$  and all  $0 \leq \alpha \leq 1$

$$\alpha \mathbf{x} + (1 - \alpha) \mathbf{y} \in S.$$

$f : S \rightarrow \mathbb{R}$  is called **convex on**  $S \subseteq \mathbb{R}^n$  if for  $\mathbf{x}, \mathbf{y} \in S$  and all  $0 \leq \alpha \leq 1$

$$f(\alpha \mathbf{x} + (1 - \alpha) \mathbf{y}) \leq \alpha f(\mathbf{x}) + (1 - \alpha) f(\mathbf{y}).$$

With ' $<$ ': *strictly convex*.

**Q:** Give an example of a convex, but not strictly convex function.



## Convexity: Consequences

If  $f$  is convex, ...

- ▶ then  $f$  is continuous at interior points.  
(Why? What would happen if it had a jump?)
- ▶ a local minimum is a *global* minimum.

If  $f$  is strictly convex, ...

- ▶ a local minimum is a *unique global* minimum.

## Optimality Conditions

If we have found a candidate  $\mathbf{x}^*$  for a minimum, how do we know it actually is one? Assume  $f$  is smooth, i.e. has all needed derivatives.

- ▶ In one dimension:
  - ▶ Necessary:  $f'(x^*) = 0$  (i.e.  $x^*$  is an extremal point)
  - ▶ Sufficient:  $f'(x^*) = 0$  and  $f''(x^*) \geq 0$   
(implies  $x^*$  is a local minimum)
- ▶ In  $n$  dimensions:
  - ▶ Necessary:  $\nabla f(\mathbf{x}^*) = 0$  (i.e.  $\mathbf{x}^*$  is an extremal point)
  - ▶ Sufficient:  $\nabla f(\mathbf{x}^*) = 0$  and  $H_f(\mathbf{x}^*)$  positive semidefinite  
(implies  $\mathbf{x}^*$  is a local minimum)

where the **Hessian**

$$H_f(\mathbf{x}^*) = \begin{bmatrix} \frac{\partial^2}{\partial x_1^2} & \cdots & \frac{\partial^2}{\partial x_1 \partial x_n} \\ \vdots & & \vdots \\ \frac{\partial^2}{\partial x_n \partial x_1} & \cdots & \frac{\partial^2}{\partial x_n^2} \end{bmatrix} f(\mathbf{x}^*).$$

## Optimization: Observations

Q: Come up with a hypothetical approach for finding minima.

A: Solve  $\nabla f = 0$ .

Q: Is the Hessian symmetric?

A: Yes. (Schwarz's theorem)

Q: How can we practically test for positive definiteness?

A: Attempt a Cholesky factorization. If it succeeds, the matrix is PD.

## Sensitivity and Conditioning (1D)

How does optimization react to a slight perturbation of the minimum?

Suppose we still have  $|f(\tilde{x}) - f(x^*)| < \text{tol}$  (where  $x^*$  is true min.).  
Apply Taylor's theorem:

$$f(x^* + h) = f(x^*) + \underbrace{f'(x^*)}_0 h + f''(x^*) \frac{h^2}{2} + O(h^3)$$

Ignore higher-order terms, solve for  $h$ :  $|\tilde{x} - x^*| \lesssim \sqrt{2 \text{tol} / f''(x^*)}$ .

**In other words:** Can expect *half as many digits* in  $\tilde{x}$  as in  $f(\tilde{x})$ .

This is important to keep in mind when setting tolerances.

It's only possible to do better when derivatives are explicitly known and convergence is not based on function values alone. (then: can solve  $\nabla f = 0$ )

## Sensitivity and Conditioning (nD)

How does optimization react to a slight perturbation of the minimum?

Suppose we still have  $|f(\tilde{\mathbf{x}}) - f(\mathbf{x}^*)| < \text{tol}$ , where  $\mathbf{x}^*$  is true min. and  $\mathbf{x} = \mathbf{x}^* + h\mathbf{s}$ . Assume  $\|\mathbf{s}\| = 1$ .

$$f(\mathbf{x}^* + h\mathbf{s}) = f(\mathbf{x}^*) + \underbrace{h \nabla f(\mathbf{x}^*)^T}_{0} \mathbf{s} + \frac{h^2}{2} \mathbf{s}^T H_f(\mathbf{x}^*) \mathbf{s} + O(h^3)$$

After ignoring high-order terms, yields:

$$|h|^2 \lesssim \frac{2 \text{tol}}{\lambda_{\min}(H_f(\mathbf{x}^*))}.$$

**In other words:** Conditioning of  $H_f$  determines sensitivity of  $\mathbf{x}^*$ .

# Unimodality

Would like a method like bisection, but for optimization.

In general: No invariant that can be preserved.

Need *extra assumption*.

$f$  is called **unimodal** on an open interval if there exists an  $x^*$  in the interval such that for all  $x_1 < x_2$  in the interval

▶  $x_2 < x^* \Rightarrow f(x_1) > f(x_2)$

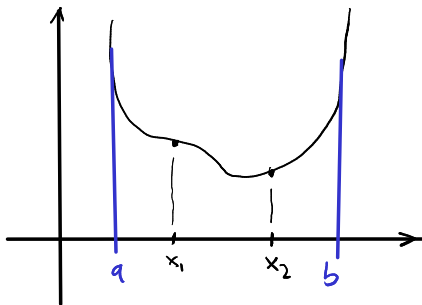
▶  $x^* < x_1 \Rightarrow f(x_1) < f(x_2)$

## In-Class Activity: Optimization Theory

In-class activity: Optimization Theory

## Golden Section Search

Suppose we have an interval with  $f$  unimodal:



Would like to maintain unimodality.

1. Pick  $x_1, x_2$
2. If  $f(x_1) > f(x_2)$ , reduce to  $(x_1, b)$
3. If  $f(x_1) \leq f(x_2)$ , reduce to  $(a, x_2)$



## Golden Section Search: Efficiency

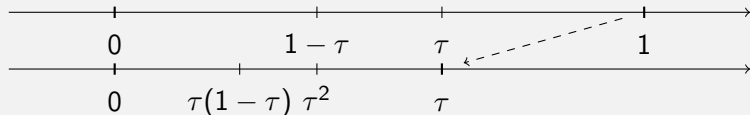
Where to put  $x_1, x_2$ ?

- ▶ Want symmetry:

$$x_1 = a + (1 - \tau)(b - a)$$

$$x_2 = a + \tau(b - a)$$

- ▶ Want to reuse function evaluations



Need:  $\tau^2 = 1 - \tau$ . Find:  $\tau = (\sqrt{5} - 1) / 2$ .

Also known as the *golden section*. Hence *golden section search*.

Convergence rate?

Linearly convergent. Can we do better?

## Newton's Method

Reuse the Taylor approximation idea, but for optimization.

$$f(x+h) \approx f(x) + f'(x)h + f''(x)\frac{h^2}{2} =: \hat{f}(h)$$

Solve  $0 = \hat{f}'(h) = f'(x) + f''(x)h$ :  $h = -f'(x)/f''(x)$ .

1.  $x_0 = \langle \text{some starting guess} \rangle$

2.  $x_{k+1} = x_k - \frac{f'(x_k)}{f''(x_k)}$

**Q:** Notice something? Identical to Newton for solving  $f'(x) = 0$ .  
Because of that: locally quadratically convergent.

**Good idea:** Combine slow-and-safe (bracketing) strategy with fast-and-risky (Newton).

[Demo: Newton's Method in 1D \[cleared\]](#)

## Steepest Descent/Gradient Descent

Given a scalar function  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  at a point  $\mathbf{x}$ , which way is down?

Direction of steepest descent:  $-\nabla f$

**Q:** How far along the gradient should we go?

Unclear—do a line search. For example using Golden Section Search.

1.  $\mathbf{x}_0 = \langle \text{some starting guess} \rangle$
2.  $\mathbf{s}_k = -\nabla f(\mathbf{x}_k)$
3. Use line search to choose  $\alpha_k$  to minimize  $f(\mathbf{x}_k + \alpha_k \mathbf{s}_k)$
4.  $\mathbf{x}_{k+1} = \mathbf{x}_k + \alpha_k \mathbf{s}_k$
5. Go to 2.

**Observation:** (from demo)

- ▶ Linear convergence

## Steepest Descent: Convergence

Consider quadratic model problem:

$$f(\mathbf{x}) = \frac{1}{2} \mathbf{x}^T A \mathbf{x} + \mathbf{c}^T \mathbf{x}$$

where  $A$  is SPD. (A good model of  $f$  near a minimum.)

Define error  $\mathbf{e}_k = \mathbf{x}_k - \mathbf{x}^*$ . Then can show:

$$\|\mathbf{e}_{k+1}\|_A = \sqrt{\mathbf{e}_{k+1}^T A \mathbf{e}_{k+1}} = \frac{\sigma_{\max}(A) - \sigma_{\min}(A)}{\sigma_{\max}(A) + \sigma_{\min}(A)} \|\mathbf{e}_k\|_A$$

where  $\|\mathbf{x}\|_A = \sqrt{\mathbf{x}^T A \mathbf{x}}$ .  $\rightarrow$  confirms linear convergence.

Convergence constant related to conditioning:

$$\frac{\sigma_{\max}(A) - \sigma_{\min}(A)}{\sigma_{\max}(A) + \sigma_{\min}(A)} = \frac{\kappa(A) - 1}{\kappa(A) + 1}.$$

# Hacking Steepest Descent for Better Convergence

## Extrapolation methods:

Look back a step, maintain '*momentum*'.

$$\mathbf{x}_{k+1} = \mathbf{x}_k - \alpha_k \nabla f(\mathbf{x}_k) + \beta_k (\mathbf{x}_k - \mathbf{x}_{k-1})$$

## Heavy ball method:

For specific constant  $\alpha_k = \alpha$  and  $\beta_k = \beta$ , can attain:

$$\|\mathbf{e}_{k+1}\|_A = \frac{\sqrt{\kappa(A)} - 1}{\sqrt{\kappa(A)} + 1} \|\mathbf{e}_k\|_A$$

[Demo: Steepest Descent \[cleared\]](#) (Part 2)

# Optimization in Machine Learning

What is *stochastic gradient descent (SGD)*?

Common in ML: Objective functions of the form

$$f(\mathbf{x}) = \frac{1}{n} \sum_{i=1}^n f_i(\mathbf{x}),$$

where each  $f_i$  comes from an *observation* (“data point”) in a (training) data set. Then “*batch*” (i.e. normal) gradient descent is

$$\mathbf{x}_{k+1} = \mathbf{x}_k - \alpha \frac{1}{n} \sum_{i=1}^n \nabla f_i(\mathbf{x}_k).$$

*Stochastic* GD uses one (or few, “*minibatch*”) observation at a time:

$$\mathbf{x}_{k+1} = \mathbf{x}_k - \alpha \nabla f_{\phi(k)}(\mathbf{x}_k).$$

[ADAM optimizer](#): GD with exp. moving avgs. of  $\nabla$  and its square.

## Conjugate Gradient Methods

Can we optimize in *the space spanned* by the last two step directions?

$$(\alpha_k, \beta_k) = \operatorname{argmin}_{\alpha_k, \beta_k} \left[ f\left(\mathbf{x}_k - \alpha_k \nabla f(\mathbf{x}_k) + \beta_k (\mathbf{x}_k - \mathbf{x}_{k-1})\right) \right]$$

- ▶ Will see in more detail later (for solving linear systems)
- ▶ Provably optimal first-order method for the quadratic model problem
- ▶ Turns out to be closely related to Lanczos ( $A$ -orthogonal search directions)

[Demo: Conjugate Gradient Method](#) [cleared]

# Nelder-Mead Method

Idea:

Form a  $n$ -point polytope in  $n$ -dimensional space and adjust worst point (highest function value) by moving it along a line passing through the centroid of the remaining points.

[Demo: Nelder-Mead Method \[cleared\]](#)



## Newton's method ( $n$ D)

What does Newton's method look like in  $n$  dimensions?

Build a Taylor approximation:

$$f(\mathbf{x} + \mathbf{s}) \approx f(\mathbf{x}) + \nabla f(\mathbf{x})^T \mathbf{s} + \frac{1}{2} \mathbf{s}^T H_f(\mathbf{x}) \mathbf{s} =: \hat{f}(\mathbf{s})$$

Then solve  $\nabla \hat{f}(\mathbf{s}) = 0$  for  $\mathbf{s}$  to find

$$H_f(\mathbf{x}) \mathbf{s} = -\nabla f(\mathbf{x}).$$

1.  $\mathbf{x}_0 = \langle \text{some starting guess} \rangle$
2. Solve  $H_f(\mathbf{x}_k) \mathbf{s}_k = -\nabla f(\mathbf{x}_k)$  for  $\mathbf{s}_k$
3.  $\mathbf{x}_{k+1} = \mathbf{x}_k + \mathbf{s}_k$

## Newton's method ( $n$ D): Observations

Drawbacks?

- ▶ Need second (!) derivatives  
(addressed by Conjugate Gradients, later in the class)
- ▶ local convergence
- ▶ Works poorly when  $H_f$  is nearly indefinite

[Demo: Newton's Method in n dimensions](#) [\[cleared\]](#)

## Quasi-Newton Methods

Secant/Broyden-type ideas carry over to optimization. How?

Come up with a way to update to update the approximate Hessian.

$$\mathbf{x}_{k+1} = \mathbf{x}_k - \alpha_k B_k^{-1} \nabla f(\mathbf{x}_k)$$

- ▶  $\alpha_k$ : a line search/damping parameter.
- ▶  $\mathbf{s}_k = \mathbf{x}_{k+1} - \mathbf{x}_k$
- ▶  $\mathbf{y}_k = \nabla f(\mathbf{x}_{k+1}) - \nabla f(\mathbf{x}_k)$
- ▶ Secant condition:  $B_{k+1} \mathbf{s}_k = \mathbf{y}_k$
- ▶ Ansatz for Hessian update:  $B_{k+1} = B_k + a \mathbf{u} \mathbf{u}^T + b \mathbf{v} \mathbf{v}^T$

**BFGS**: Secant-type method, similar to Broyden:

$$B_{k+1} = B_k + \frac{\mathbf{y}_k \mathbf{y}_k^T}{\mathbf{y}_k^T \mathbf{s}_k} - \frac{B_k \mathbf{s}_k \mathbf{s}_k^T B_k}{\mathbf{s}_k^T B_k \mathbf{s}_k}$$

# In-Class Activity: Optimization Methods

In-class activity: Optimization Methods

## Nonlinear Least Squares: Setup

What if the  $f$  to be minimized is actually a 2-norm?

$$f(\mathbf{x}) = \|\mathbf{r}(\mathbf{x})\|_2, \quad \mathbf{r}(\mathbf{x}) = \mathbf{y} - \mathbf{a}(\mathbf{x})$$

Define 'helper function'

$$\varphi(\mathbf{x}) = \frac{1}{2} \mathbf{r}(\mathbf{x})^T \mathbf{r}(\mathbf{x}) = \frac{1}{2} f^2(\mathbf{x})$$

and minimize that instead.

$$\frac{\partial}{\partial x_i} \varphi = \frac{1}{2} \sum_{j=1}^n \frac{\partial}{\partial x_i} [r_j(\mathbf{x})^2] = \sum_j \left( \frac{\partial}{\partial x_i} r_j \right) r_j,$$

or, in matrix form:

$$\nabla \varphi = J_r(\mathbf{x})^T \mathbf{r}(\mathbf{x}).$$

## Gauss-Newton

For brevity:  $J := J_r(\mathbf{x})$ .

Can show similarly:

$$H_\varphi(\mathbf{x}) = J^T J + \sum_i r_i H_{r_i}(\mathbf{x}).$$

Newton step  $\mathbf{s}$  can be found by solving  $H_\varphi(\mathbf{x})\mathbf{s} = -\nabla\varphi$ .

**Observation:**  $\sum_i r_i H_{r_i}(\mathbf{x})$  is inconvenient to compute *and* unlikely to be large (since it's multiplied by components of the residual, which is supposed to be small)  $\rightarrow$  forget about it.

**Gauss-Newton method:** Find step  $\mathbf{s}$  by  $J^T J \mathbf{s} = -\nabla\varphi = -J^T \mathbf{r}(\mathbf{x})$ . Does that remind you of the *normal equations*?  $J\mathbf{s} \cong -\mathbf{r}(\mathbf{x})$ . Solve that using our existing methods for least-squares problems.

# Gauss-Newton: Observations?

## Demo: Gauss-Newton [cleared]

Observations?

- ▶ Newton on its own is still only locally convergent
- ▶ Gauss-Newton is clearly similar
- ▶ It's worse because the step is only approximate  
→ Much depends on the starting guess.

## Levenberg-Marquardt

If Gauss-Newton on its own is poorly conditioned, can try **Levenberg-Marquardt**:

$$(J_r(\mathbf{x}_k)^T J_r(\mathbf{x}_k) + \mu_k I) \mathbf{s}_k = -J_r(\mathbf{x}_k)^T \mathbf{r}(\mathbf{x}_k)$$

for a 'carefully chosen'  $\mu_k$ . This makes the system matrix 'more invertible' but also less accurate/faithful to the problem.

What Levenberg-Marquardt does is generically called **regularization**:  
Make  $H$  more positive definite.

Easy to rewrite to least-squares problem:

$$\begin{bmatrix} J_r(\mathbf{x}_k) \\ \sqrt{\mu_k} I \end{bmatrix} \mathbf{s}_k \cong \begin{bmatrix} -\mathbf{r}(\mathbf{x}_k) \\ 0 \end{bmatrix}.$$



## Constrained Optimization: Problem Setup

Want  $\mathbf{x}^*$  so that

$$f(\mathbf{x}^*) = \min_{\mathbf{x}} f(\mathbf{x}) \quad \text{subject to} \quad \mathbf{g}(\mathbf{x}) = 0$$

No inequality constraints just yet. This is *equality-constrained optimization*. Develop a (local) necessary condition for a minimum.

Necessary cond.: “no feasible descent possible”. Assume  $\mathbf{g}(\mathbf{x}) = 0$ .

Recall unconstrained necessary condition, “no descent possible”:

$$\nabla f(\mathbf{x}) = 0$$

Look for **feasible descent directions** from  $\mathbf{x}$ . (Necessary cond.:  $\nexists$ )

$\mathbf{s}$  is a **feasible direction** at  $\mathbf{x}$  if

$$\mathbf{x} + \alpha \mathbf{s} \quad \text{feasible for } \alpha \in [0, r] \quad (\text{for some } r)$$

## Constrained Optimization: Necessary Condition

Need:  $\nabla f(\mathbf{x}) \cdot \mathbf{s} \geq 0$  (“uphill that way”) for any feasible direction  $\mathbf{s}$ .

- ▶ **Not at boundary:**  $\mathbf{s}$  and  $-\mathbf{s}$  are feasible directions

$$\Rightarrow \nabla f(\mathbf{x}) = 0$$

$\Rightarrow$  Only the boundary of the feasible set is different from the unconstrained case (i.e. interesting)

- ▶ **At boundary:** (the common case)  $\mathbf{g}(\mathbf{x}) = 0$ . Need:

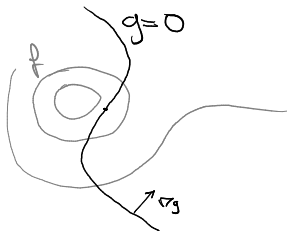
$$-\nabla f(\mathbf{x}) \in \text{rowspan}(J_{\mathbf{g}})$$

a.k.a. “all descent directions would cause a change ( $\rightarrow$ violation) of the constraints.”

**Q:** Why ‘rowspan’? Think about shape of  $J_{\mathbf{g}}$ .

$$\Leftrightarrow -\nabla f(\mathbf{x}) = J_{\mathbf{g}}^T \boldsymbol{\lambda} \quad \text{for some } \boldsymbol{\lambda}.$$

## Lagrange Multipliers



Seen: Need  $-\nabla f(\mathbf{x}) = J_{\mathbf{g}}^T \boldsymbol{\lambda}$  at the (constrained) optimum.

*Idea:* Turn constrained optimization problem for  $\mathbf{x}$  into an *unconstrained* optimization problem for  $(\mathbf{x}, \boldsymbol{\lambda})$ . How?

Need a new function  $\mathcal{L}(\mathbf{x}, \boldsymbol{\lambda})$  to minimize:

$$\mathcal{L}(\mathbf{x}, \boldsymbol{\lambda}) := f(\mathbf{x}) + \boldsymbol{\lambda}^T \mathbf{g}(\mathbf{x}).$$

## Lagrange Multipliers: Development

$$\mathcal{L}(\mathbf{x}, \boldsymbol{\lambda}) := f(\mathbf{x}) + \boldsymbol{\lambda}^T \mathbf{g}(\mathbf{x}).$$

Then:  $\nabla \mathcal{L} = 0$  at unconstrained minimum, i.e.

$$0 = \nabla \mathcal{L} = \begin{bmatrix} \nabla_{\mathbf{x}} \mathcal{L} \\ \nabla_{\boldsymbol{\lambda}} \mathcal{L} \end{bmatrix} = \begin{bmatrix} \nabla f + J_{\mathbf{g}}(\mathbf{x})^T \boldsymbol{\lambda} \\ \mathbf{g}(\mathbf{x}) \end{bmatrix}.$$

Convenient: This matches our necessary condition!

So we could use any unconstrained method to minimize  $\mathcal{L}$ .

**For example:** Using Newton to minimize  $\mathcal{L}$  is called *Sequential Quadratic Programming*. ('SQP')

[Demo: Sequential Quadratic Programming \[cleared\]](#)

## Inequality-Constrained Optimization

Want  $\mathbf{x}^*$  so that

$$f(\mathbf{x}^*) = \min_{\mathbf{x}} f(\mathbf{x}) \quad \text{subject to} \quad \mathbf{g}(\mathbf{x}) = 0 \quad \text{and} \quad \mathbf{h}(\mathbf{x}) \leq 0.$$

Develop a necessary condition for a minimum.

Again: Assume we're at a feasible point, on the boundary of the feasible region. Must ensure descent directions are *infeasible*.

**Motivation:**  $\mathbf{g} = 0 \Leftrightarrow$  two inequality constraints:  $\mathbf{g} \leq 0 \wedge \mathbf{g} \geq 0$ .

Consider the condition  $-\nabla f(\mathbf{x}) = J_{\mathbf{h}}^T \boldsymbol{\lambda}_2$ .

- ▶ Descent direction must start violating constraint.  
But only one direction is dangerous here!
- ▶  $-\nabla f$ : **descent** direction of  $f$ ,  $\nabla h_i$ : **ascent** direction of  $h_i$
- ▶ If we assume  $\lambda_2 > 0$ , going towards  $-\nabla f$  would increase  $\mathbf{h}$   
(and start violating  $\mathbf{h} \leq 0$ )

## Lagrangian, Active/Inactive

Put together the overall Lagrangian.

$$\mathcal{L}(\mathbf{x}, \lambda_1, \lambda_2) := f(\mathbf{x}) + \lambda_1^T \mathbf{g}(\mathbf{x}) + \lambda_2^T \mathbf{h}(\mathbf{x})$$

What are **active** and **inactive** constraints?

- ▶ **Active:** active  $\Leftrightarrow h_i(\mathbf{x}^*) = 0 \Leftrightarrow$  at 'boundary' of ineq. constraint  
(Equality constraints are always 'active')
- ▶ **Inactive:** If  $h_i$  inactive ( $h_i(\mathbf{x}^*) < 0$ ), must force  $\lambda_{2,i} = 0$ .  
Otherwise: Behavior of  $h$  could change location of minimum of  $\mathcal{L}$ . Use **complementarity condition**  $h_i(\mathbf{x}^*)\lambda_{2,i} = 0$ .  
 $\Leftrightarrow$  *at least one* of  $h_i(\mathbf{x}^*)$  and  $\lambda_{2,i}$  is zero.

## Karush-Kuhn-Tucker (KKT) Conditions

Develop a set of necessary conditions for a minimum.

Assuming  $J_g$  and  $J_{h,\text{active}}$  have full rank, this set of conditions is *necessary*:

$$(*) \quad \nabla_{\mathbf{x}} \mathcal{L}(\mathbf{x}^*, \lambda_1^*, \lambda_2^*) = 0$$

$$(*) \quad \mathbf{g}(\mathbf{x}^*) = 0$$

$$\mathbf{h}(\mathbf{x}^*) \leq 0$$

$$\lambda_2 \geq 0$$

$$(*) \quad \mathbf{h}(\mathbf{x}^*) \cdot \lambda_2 = 0$$

These are called the **Karush-Kuhn-Tucker ('KKT') conditions**.

**Computational approach:** Solve (\*) equations by Newton.

# Outline

Introduction to Scientific Computing

Systems of Linear Equations

Linear Least Squares

Eigenvalue Problems

Nonlinear Equations

Optimization

**Interpolation**

Introduction

Methods

Error Estimation

Piecewise interpolation, Splines

Numerical Integration and Differentiation

Initial Value Problems for ODEs

Boundary Value Problems for ODEs

Partial Differential Equations and Sparse Linear Algebra

Fast Fourier Transform

Additional Topics



## Interpolation: Setup

**Given:**  $(x_i)_{i=1}^N, (y_i)_{i=1}^N$

**Wanted:** Function  $f$  so that  $f(x_i) = y_i$

How is this not the same as function fitting? (from least squares)

It's very similar—but the key difference is that we are asking for *exact equality*, not just minimization of a residual norm.

→ Better error control, error not dominated by residual

**Idea:** There is an *underlying function* that we are approximating from the known point values.

**Error here:** Distance from that underlying function

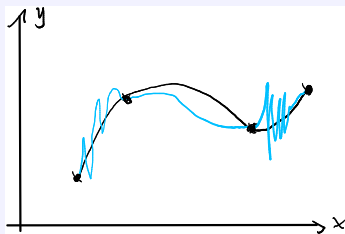
## Interpolation: Setup (II)

**Given:**  $(x_i)_{i=1}^N, (y_i)_{i=1}^N$

**Wanted:** Function  $f$  so that  $f(x_i) = y_i$

Does this problem have a unique answer?

No—there are infinitely many functions that satisfy the problem as stated:



# Interpolation: Importance

Why is interpolation important?

It brings all of calculus within range of numerical operations.

▶ Why?

Because calculus works on functions.

▶ How?

1. Interpolate (go from discrete to continuous)
2. Apply calculus
3. Re-discretize (evaluate at points)

## Making the Interpolation Problem Unique

Limit the set of functions to span of an *interpolation basis*  $\{\varphi_i\}_{i=1}^{N_{\text{Func}}}$ :

$$p_{n-1}(x) = \sum_{j=1}^{N_{\text{func}}} \alpha_j \varphi_j(x)$$

Interpolation becomes solving the linear system:

$$y_i = p_{n-1}(x_i) = \sum_{j=1}^{N_{\text{func}}} \alpha_j \underbrace{\varphi_j(x_i)}_{V_{ij}} \quad \Leftrightarrow \quad V\alpha = \mathbf{y}.$$

- ▶ Want unique answer: Pick  $N_{\text{func}} = N \rightarrow V$  square.
- ▶  $V$  is called the (*generalized*) *Vandermonde matrix*.
- ▶  $V(\text{coefficients}) = (\text{values at nodes})$ .
- ▶ Can prescribe derivatives: Use  $\varphi'_j, f'$  in a row. (**Hermite interp.**)

## Existence/Sensitivity

Solution to the interpolation problem: Existence? Uniqueness?

Equivalent to existence/uniqueness of the linear system

Sensitivity?

- ▶ Shallow answer: Simply consider the condition number of the linear system
- ▶  $\|\text{coefficients}\|$  does not suffice as measure of stability.  
 $f(x)$  can be evaluated in many places. ( $f$  is interpolant.)
- ▶ Want:  $\max_{x \in [a,b]} |f(x)| \leq \Lambda \|\mathbf{y}\|_\infty$
- ▶  $\Lambda$ : **Lebesgue constant**
- ▶  $\Lambda$  depends on  $n$  and  $\{x_i\}_i$ 
  - ▶ Technically also depends on  $\{\phi_i\}_i$
  - ▶ But: same for all polynomial bases

## Modes and Nodes (aka Functions and Points)

Both function basis and point set are under our control. What do we pick?

Ideas for basis functions:

- ▶ Monomials  $1, x, x^2, x^3, x^4, \dots$
- ▶ Functions that make  $V = I \rightarrow$  'Lagrange basis'
- ▶ Functions that make  $V$  triangular  $\rightarrow$  'Newton basis'
- ▶ *Splines* (piecewise polynomials)
- ▶ *Orthogonal polynomials*
- ▶ Sines and cosines
- ▶ 'Bumps' ('*Radial Basis Functions*')

Ideas for points:

- ▶ Equispaced
- ▶ '*Edge-Clustered*' (so-called Chebyshev/Gauss/... nodes)

Specific issues:

- ▶ Why *not* monomials on equispaced points?  
Demo: Monomial interpolation  
[cleared]
- ▶ Why not equispaced?  
Demo: Choice of Nodes for Polynomial Interpolation  
[cleared]

## Lagrange Interpolation

Find a basis so that  $V = I$ , i.e.

$$\varphi_j(x_i) = \begin{cases} 1 & i = j, \\ 0 & \text{otherwise.} \end{cases}$$

Start with simple example. Three nodes:  $x_1, x_2, x_3$

$$\varphi_1(x) = \frac{(x - x_2)(x - x_3)}{(x_1 - x_2)(x_1 - x_3)}$$

$$\varphi_2(x) = \frac{(x - x_1)(x - x_3)}{(x_2 - x_1)(x_2 - x_3)}$$

$$\varphi_3(x) = \frac{(x - x_1)(x - x_2)}{(x_3 - x_1)(x_3 - x_2)}$$

Numerator: Ensures  $\varphi_i$  zero at other nodes.

Denominator: Ensures  $\varphi_i(x_i) = 1$ .

## Lagrange Polynomials: General Form

$$\varphi_j(x) = \frac{\prod_{k=1, k \neq j}^m (x - x_k)}{\prod_{k=1, k \neq j}^m (x_j - x_k)}$$

Write down the Lagrange interpolant for nodes  $(x_i)_{i=1}^m$  and values  $(y_i)_{i=1}^m$ .

$$p_{m-1}(x) = \sum_{j=1}^m y_j \varphi_j(x)$$



## Newton Interpolation

Find a basis so that  $V$  is triangular.

Easier to build than Lagrange, but: coefficient finding costs  $O(n^2)$ .

$$\varphi_j(x) = \prod_{k=1}^{j-1} (x - x_k).$$

(At least) two possibilities for coefficient finding:

- ▶ Set up  $V$ , run forward substitution.
- ▶ [Divided Differences](#) (Wikipedia link)

Why not Lagrange/Newton?

Cheap to form, expensive to evaluate, expensive to do calculus on.

## Better conditioning: Orthogonal polynomials

What caused monomials to have a terribly conditioned Vandermonde?

Being close to linearly dependent.

What's a way to make sure two vectors are *not* like that?

Orthogonality

But polynomials are functions!

## Orthogonality of Functions

How can functions be orthogonal?

Need an **inner product**. Orthogonal then just means  $\langle f, g \rangle = 0$ .

$$\mathbf{f} \cdot \mathbf{g} = \sum_{i=1}^n f_i g_i = \langle \mathbf{f}, \mathbf{g} \rangle$$
$$\langle f, g \rangle = \int_{-1}^1 f(x)g(x)dx$$

## Constructing Orthogonal Polynomials

How can we find an orthogonal basis?

Apply Gram-Schmidt to the monomials.

**Demo: Orthogonal Polynomials [cleared]** — Got: **Legendre polynomials**.  
But how can I practically compute the Legendre polynomials?

→ [DLMF: Chapter on orthogonal polynomials](#)

- ▶ Three-term recurrence:  $(n + 1)P_{n+1} = (2n + 1)xP_n - nP_{n-1}$ ,  
 $P_0 = 1, P_1 = x$
- ▶ A whole zoo of polynomials there, depending on the weight function  $w$  in the (generalized) inner product:

$$\langle f, g \rangle = \int w(x)f(x)g(x)dx.$$

## Chebyshev Polynomials: Definitions

Three equivalent definitions:

- ▶ Result of Gram-Schmidt with weight  $1/\sqrt{1-x^2}$ . What is that weight?

$$1/(\text{Half circle}), \text{ i.e. } x^2 + y^2 = 1, \text{ with } y = \sqrt{1-x^2}$$

(Like for Legendre, you won't exactly get the standard normalization if you do this.)

- ▶  $T_k(x) = \cos(k \cos^{-1}(x))$
- ▶  $T_k(x) = 2xT_{k-1}(x) - T_{k-2}(x)$  plus  $T_0 = 1, T_1 = x$

## Chebyshev Interpolation

What is the Vandermonde matrix for Chebyshev polynomials?

- ▶ Need to know the nodes to answer that
- ▶ The answer would be very simple if the nodes were  $\cos(*)$ .
- ▶ So why not  $\cos(\text{equispaced})$ ? Maybe

$$x_i = \cos\left(\frac{i}{k}\pi\right) \quad (i = 0, 1, \dots, k)$$

These are just the *extrema* (minima/maxima) of  $T_k$ .

$$V_{ij} = \cos\left(j \cos^{-1}\left(\cos\left(\frac{i}{k}\pi\right)\right)\right) = \cos\left(j \frac{i}{k}\pi\right).$$

- ▶ Called the *Discrete Cosine Transform* (DCT)
- ▶ Matvec (and inverse!) available with  $O(N \log N)$  cost ( $\rightarrow$  FFT)

## Chebyshev Nodes

Might also consider roots (instead of extrema) of  $T_k$ :

$$x_i = \cos\left(\frac{2i-1}{2k}\pi\right) \quad (i = 1 \dots, k).$$

Vandermonde for these (with  $T_k$ ) can be applied in  $O(N \log N)$  time, too.

Edge-clustering seemed like a good thing in interpolation nodes. Do these do that?

Yes.

[Demo: Chebyshev Interpolation \[cleared\]](#) (Part I-IV)

## Chebyshev Interpolation: Summary

- ▶ Chebyshev interpolation is fast and works extremely well
- ▶ <http://www.chebfun.org/> and: [ATAP](#)
- ▶ In 1D, they're a very good answer to the interpolation question
- ▶ But sometimes a piecewise approximation (with a specifiable level of smoothness) is more suited to the application



## In-Class Activity: Interpolation

In-class activity: Interpolation

## Truncation Error in Interpolation

If  $f$  is  $n$  times continuously differentiable on a closed interval  $I$  and  $p_{n-1}(x)$  is a polynomial of degree at most  $n$  that interpolates  $f$  at  $n$  distinct points  $\{x_i\}$  ( $i = 1, \dots, n$ ) in that interval, then for each  $x$  in the interval there exists  $\xi$  in that interval such that

$$f(x) - p_{n-1}(x) = \frac{f^{(n)}(\xi(x))}{n!} (x - x_1)(x - x_2) \cdots (x - x_n).$$

Define the error  $R(x) := f(x) - p_{n-1}(x)$  and an auxiliary function:

$$Y_x(t) = R(t) - \frac{R(x)}{W(x)} W(t) \quad \text{where} \quad W(t) = \prod_{i=1}^n (t - x_i).$$

- ▶ Let  $x \in I$  (where we prove the identity) with  $x \neq x_i$  for  $i \in \{1, \dots, n\}$  (WLOG).
- ▶ Note the introduction of the *additional* variable  $t$ .

## Truncation Error in Interpolation: cont'd.

$$Y_x(t) = R(t) - \frac{R(x)}{W(x)}W(t) \quad \text{where} \quad W(t) = \prod_{i=1}^n (t - x_i)$$

- ▶ Since  $x_i$  are roots of  $R(t)$  and  $W(t)$ , we have  $Y_x(x) = Y_x(x_i) = 0$ , which means  $Y_x$  has at least  $n + 1$  roots.
- ▶ From Rolle's theorem,  $Y_x'(t)$  has at least  $n$  roots, then  $Y_x^{(n)}$  has at least one root  $\xi$ , where  $\xi \in I$ .
- ▶ Since  $p_{n-1}(x)$  is a polynomial of degree at most  $n - 1$ ,  $R^{(n)}(t) = f^{(n)}(t)$ . Thus

$$Y_x^{(n)}(t) = f^{(n)}(t) - \frac{R(x)}{W(x)}n!.$$

- ▶ Plugging  $Y_x^{(n)}(\xi) = 0$  into the above yields the result.

## Error Result: Connection to Chebyshev

What is the connection between the error result and Chebyshev interpolation?

- ▶ The error bound suggests choosing the interpolation nodes such that the product  $|\prod_{i=1}^n (x - x_i)|$  is as small as possible. The Chebyshev nodes achieve this.
- ▶ If nodes are edge-clustered,  $\prod_{i=1}^n (x - x_i)$  clamps down the (otherwise quickly-growing) error there.
- ▶ Confusing: *Chebyshev approximating polynomial* (or “polynomial best-approximation”). **Not** the Chebyshev interpolant.
- ▶ Chebyshev nodes also do **not** minimize the Lebesgue constant.

[Demo: Chebyshev Interpolation \[cleared\]](#) (Part V)

## Error Result: Simplified Form

Boil the error result down to a simpler form.

Assume  $x_1 < \dots < x_n$ .

- ▶  $|f^{(n)}(x)| \leq M$  for  $x \in [x_1, x_n]$ ,
- ▶ Set the interval length  $h = x_n - x_1$ .  
Then  $|x - x_i| \leq h$ .

Altogether—there is a constant  $C$  independent of  $h$  so that:

$$\max_x |f(x) - p_{n-1}(x)| \leq CMh^n.$$

For the grid spacing  $h \rightarrow 0$ , we have  $E(h) = O(h^n)$ . This is called *convergence of order  $n$* .

- ▶ [Demo: Interpolation Error \[cleared\]](#)
- ▶ [Demo: Jump with Chebyshev Nodes \[cleared\]](#)

## Going piecewise: Simplest Case

Construct a piecewise linear interpolant at four points.

$x_0, y_0$		$x_1, y_1$		$x_2, y_2$		$x_3, y_3$
	$f_1 = a_1x + b_1$		$f_2 = a_2x + b_2$		$f_3 = a_3x + b_3$	
	2 unk.		2 unk.		2 unk.	
	$f_1(x_0) = y_0$		$f_2(x_1) = y_1$		$f_3(x_2) = y_2$	
	$f_1(x_1) = y_1$		$f_2(x_2) = y_2$		$f_3(x_3) = y_3$	
	2 eqn.		2 eqn.		2 eqn.	

Why three intervals?

General situation  $\rightarrow$  two end intervals and one middle interval. Can just add more middle intervals if needed.

# Piecewise Cubic ('Splines')

$x_0, y_0$		$x_1, y_1$		$x_2, y_2$		$x_3, y_3$
	$f_1$		$f_2$		$f_3$	
	$a_1x^3 + b_1x^2 + c_1x + d_1$		$a_2x^3 + b_2x^2 + c_2x + d_2$		$a_3x^3 + b_3x^2 + c_3x + d_3$	

4 unknowns    4 unknowns    4 unknowns

$$f_1(x_0) = y_0 \quad f_2(x_1) = y_1 \quad f_3(x_2) = y_2$$

$$f_1(x_1) = y_1 \quad f_2(x_2) = y_2 \quad f_3(x_3) = y_3$$

Not enough: need more conditions. Ask for more smoothness.

$$f_1'(x_1) = f_2'(x_1) \quad f_2'(x_2) = f_3'(x_2)$$

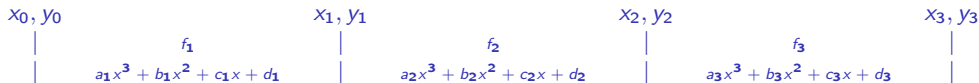
$$f_1''(x_1) = f_2''(x_1) \quad f_2''(x_2) = f_3''(x_2)$$

Not enough: need yet more conditions.

$$f_1''(x_0) = 0 \quad f_3''(x_3) = 0$$

Now: have a square system.

## Piecewise Cubic ('Splines'): Accounting



Number of conditions:  $2N_{\text{intervals}} + 2N_{\text{middle nodes}} + 2$  where

$$N_{\text{intervals}} - 1 = N_{\text{middle nodes}}$$

so

$$2N_{\text{intervals}} + 2(N_{\text{intervals}} - 1) + 2 = 4N_{\text{intervals}},$$

which is exactly the number of unknown coefficients.

These conditions are fairly arbitrary: Can choose different ones basically at will. The above choice: '*natural spline*'.

Can also come up with a basis of spline functions (with the chosen smoothness conditions). These are called **B-Splines**.



# Outline

Introduction to Scientific Computing

Systems of Linear Equations

Linear Least Squares

Eigenvalue Problems

Nonlinear Equations

Optimization

Interpolation

**Numerical Integration and Differentiation**

Numerical Integration

Quadrature Methods

Accuracy and Stability

Gaussian Quadrature

Composite Quadrature

Numerical Differentiation

Richardson Extrapolation

Initial Value Problems for ODEs

Boundary Value Problems for ODEs

Partial Differential Equations and Sparse Linear Algebra

Fast Fourier Transform

Additional Topics

# Numerical Integration: About the Problem

What is numerical integration? (Or **quadrature**?)

Given  $a$ ,  $b$ ,  $f$ , approximate

$$\int_a^b f(x)dx.$$

What about existence and uniqueness?

- ▶ Answer exists e.g. if  $f$  is *integrable* in the Riemann or Lebesgue senses.
- ▶ Answer is unique if  $f$  is e.g. piecewise continuous and bounded. (this also implies existence)

## Conditioning

Derive the (absolute) condition number for numerical integration.

Let  $\hat{f}(x) := f(x) + e(x)$ , where  $e(x)$  is a perturbation.

$$\begin{aligned} & \left| \int_a^b f(x) dx - \int_a^b \hat{f}(x) dx \right| \\ &= \left| \int_a^b e(x) dx \right| \leq \int_a^b |e(x)| dx \leq (b-a) \max_{x \in [a,b]} |e(x)|. \end{aligned}$$

## Interpolatory Quadrature: Examples

**Example:** Fix  $(x_i)$ . Then  $f(x) \approx p_{n-1}(x) = \sum_i f(x_i)l_i(x)$ , where  $l_i(x)$  is the Lagrange polynomial for the node  $x_i$ . Then

$$\int_a^b f(x)dx \approx \sum_i f(x_i) \underbrace{\int_a^b l_i(x)dx}_{\omega_i}.$$

→ a computational recipe (“**quadrature rule**”) applicable to any  $f$ :

$$\int_a^b f(x)dx \approx \sum_i f(x_i)\omega_i.$$

$x_i$ : **quadrature nodes**,  $\omega_i$ : **quadrature weights**.

- ▶ Equispaced nodes: **Newton-Cotes quadrature**.
- ▶ Chebyshev nodes: **Clenshaw-Curtis quadrature**.

## Interpolatory Quadrature: Computing Weights

How do the weights in interpolatory quadrature get computed?

Done by solving linear system.

**Know:** This quadrature should at least integrate monomials exactly.

$$\begin{aligned} b - a &= \int_a^b 1 dx = \omega_1 \cdot 1 + \cdots + \omega_n \cdot 1 \\ &\vdots \\ \frac{1}{k+1} (b^{k+1} - a^{k+1}) &= \int_a^b x^k dx = \omega_1 \cdot x_1^k + \cdots + \omega_n \cdot x_n^k \end{aligned}$$

Write down  $n$  equations for  $n$  unknowns, solve linear system, done.

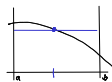
This is called the *method of undetermined coefficients*.

[Demo: Newton-Cotes weight finder](#) [cleared]

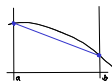
## Examples and Exactness

To what polynomial degree are the following rules exact?

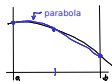
*Midpoint rule*       $(b - a)f\left(\frac{a+b}{2}\right)$



*Trapezoidal rule*       $\frac{b-a}{2}(f(a) + f(b))$



*Simpson's rule*       $\frac{b-a}{6}\left(f(a) + 4f\left(\frac{a+b}{2}\right) + f(b)\right)$



**Midpoint:** technically 0 (constants), actually 1 (linears)

**Trapezoidal:** 1 (linears)

**Simpson's:** technically 2 (parabolas), actually 3 (cubics)

- ▶ Cancellation of odd-order error requires symmetric nodes.
- ▶ (trapez. – midpt.) usable as (“*a-posteriori*”) error estimate.

## Interpolatory Quadrature: Accuracy

Let  $p_{n-1}$  be an interpolant of  $f$  at nodes  $x_1, \dots, x_n$  (of degree  $n - 1$ )

Recall

$$\sum_i \omega_i f(x_i) = \int_a^b p_{n-1}(x) dx.$$

What can you say about the accuracy of the method?

Notation:  $\|f\|_\infty = \max_{x \in [a,b]} |f(x)|$

$$\begin{aligned} & \left| \int_a^b f(x) dx - \int_a^b p_{n-1}(x) dx \right| \\ & \leq \int_a^b |f(x) - p_{n-1}(x)| dx \\ & \leq (b-a) \|f - p_{n-1}\|_\infty \\ \text{(using interp. error)} & \leq C(b-a)h^n \|f^{(n)}\|_\infty \leq Ch^{n+1} \|f^{(n)}\|_\infty \end{aligned}$$

## Quadrature: Overview of Rules

	$n$	Deg.	Ex.Int.Deg. (w/odd)	Intp.Ord.	Quad.Ord. (regular)	Quad.Ord. (w/odd)
		$n - 1$	$(n - 1) + 1_{\text{odd}}$	$n$	$n + 1$	$(n + 1) + 1_{\text{odd}}$
Midp.	1	0	1	1	2	3
Trapz.	2	1	1	2	3	3
Simps.	3	2	3	3	4	5
S. 3/8	4	3	3	4	5	5

- ▶  $n$ : number of points
- ▶ “Deg.”: Degree of polynomial used in interpolation ( $= n - 1$ )
- ▶ “Ex.Int.Deg.”: Polynomials of up to (and including) this degree *actually* get integrated exactly. (including the odd-order bump)
- ▶ “Intp.Ord.”: Order of Accuracy of Interpolation:  $O(h^n)$
- ▶ “Quad.Ord. (regular)”: Order of accuracy for quadrature predicted by the error result above:  $O(h^{n+1})$
- ▶ “Quad.Ord. (w/odd)”: Actual order of accuracy for quadrature given ‘bonus’ degrees for rules with odd point count

**Observation:** Quadrature gets (at least) ‘one order higher’ than interpolation—even more for odd-order rules. (i.e. more accurate)



## Interpolatory Quadrature: Stability

Let  $p_{n-1}$  be an interpolant of  $f$  at nodes  $x_1, \dots, x_n$  (of degree  $n - 1$ )

Recall

$$\sum_i \omega_i f(x_i) = \int_a^b p_{n-1}(x) dx$$

What can you say about the stability of this method?

Again consider  $\hat{f}(x) = f(x) + e(x)$ .

$$\left| \sum_i \omega_i f(x_i) - \sum_i \omega_i \hat{f}(x_i) \right| \leq \sum_i |\omega_i e(x_i)| \leq \left( \sum_i |\omega_i| \right) \|e\|_\infty$$

So, what quadrature weights make for bad stability bounds?

Quadratures with large negative weights. (Recall:  $\sum_i \omega_i$  is fixed.)

## About Newton-Cotes

What's not to like about Newton-Cotes quadrature?

[Demo: Newton-Cotes weight finder \[cleared\]](#) (again, with many nodes)

In fact, Newton-Cotes must have at least one negative weight as soon as  $n \geq 11$ .

More drawbacks:

- ▶ All the fun of high-order interpolation with monomials and equispaced nodes (i.e. convergence not guaranteed)
- ▶ Weights possibly non-negative ( $\rightarrow$  stability issues)
- ▶ Coefficients determined by (possibly ill-conditioned) Vandermonde matrix
- ▶ Thus hard to extend to arbitrary number of points.

## Gaussian Quadrature

So far: nodes chosen from outside.

Can we gain something if we let the quadrature rule choose the nodes, too? **Hope:** More design freedom  $\rightarrow$  Exact to higher degree.

**Idea:** method of undetermined coefficients

**But:** Resulting system would be nonlinear.

Can use orthogonal polynomials to get a leg up.

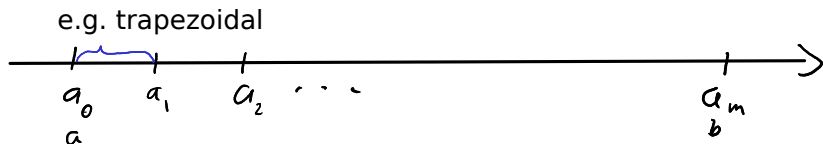
**Gaussian quadrature** with  $n$  points: Exactly integrates polynomials up to degree  $2n - 1$ .

[Demo: Gaussian quadrature weight finder \[cleared\]](#)

## Composite Quadrature

High-order polynomial interpolation requires a high degree of smoothness of the function.

**Idea:** Stitch together multiple lower-order quadrature rules to alleviate smoothness requirement.



## Error in Composite Quadrature

What can we say about the error in the case of composite quadrature?

Error for one panel of length  $h$ :  $|\int (f - p_{n-1})| \leq C \cdot h^{n+1} \|f^{(n)}\|_{\infty}$

$$\left| \int_a^b f(x) dx - \sum_{j=1}^m \sum_{i=1}^n \omega_{j,i} f(x_{j,i}) \right| \leq C \|f^{(n)}\|_{\infty} \sum_{j=1}^m (a_j - a_{j-1})^{n+1}$$

$$= C \|f^{(n)}\|_{\infty} \sum_{j=1}^m \underbrace{(a_j - a_{j-1})^n}_{\leq h^n} (a_j - a_{j-1}) \leq C \|f^{(n)}\|_{\infty} h^n (b - a),$$

where  $h$  is the length of a single panel.

## Composite Quadrature: Notes

**Observation:** Composite quadrature loses an order compared to non-composite.

**Idea:** If we can estimate errors on each subinterval, we can shrink (e.g. by splitting in half) only those contributing the most to the error.  
(**adaptivity**)

## Taking Derivatives Numerically

Why *shouldn't* you take derivatives numerically?

- ▶ 'Unbounded'

A function with small  $\|f\|_\infty$  can have arbitrarily large  $\|f'\|_\infty$

- ▶ Amplifies noise

Imagine a smooth function perturbed by small, high-frequency wiggles

- ▶ Subject to cancellation error

- ▶ Inherently less accurate than integration

- ▶ Interpolation:  $h^n$

- ▶ Quadrature:  $h^{n+1}$

- ▶ Differentiation:  $h^{n-1}$

(where  $n$  is the number of points)

## Numerical Differentiation: How?

How can we take derivatives numerically?

Let  $\mathbf{x} = (x_i)_{i=1}^n$  be nodes and  $(\varphi_i)_{i=1}^n$  an interpolation basis. Find interpolation coefficients  $\boldsymbol{\alpha} = (\alpha_i)_{i=1}^n = V^{-1}f(\mathbf{x})$ . Then

$$f(\xi) \approx p_{n-1}(\xi) = \sum_{i=1}^n \alpha_i \varphi_i(\xi).$$

Then, simply take a derivative:

$$f'(\xi) \approx p'_{n-1}(\xi) = \sum_{i=1}^n \alpha_i \varphi'_i(\xi).$$

$\varphi'_i$  are known because the interpolation basis  $\varphi_i$  is known!

[Demo: Taking Derivatives with Vandermonde Matrices](#) [cleared] (Basics)



## Numerical Differentiation: Accuracy

Interpolation error: If  $\mathbf{x} = (x_i)_{i=1}^n$  and  $f(\mathbf{x}) = p_{n-1}(\mathbf{x})$ , then

$$f(x) - p_{n-1}(x) = \frac{f^{(n)}(\xi)}{n!} \prod_{i=1}^n (x - x_i).$$
$$f'(x) - p'_{n-1}(x) \approx \frac{f^{(n)}(\xi)}{n!} \left( \prod_{i=1}^n (x - x_i) \right)'$$

(ignoring dependency of  $\xi$  on  $x$ ) Note that the derivative of  $\prod(x - x_i)$  has  $n - 1$  roots (extrema of  $\prod(x - x_i)$ ), interspersed between  $(x_i)$ , so can be bounded by  $h^{n-1}$ :

$$|f'(x) - p'_{n-1}(x)| \leq C \left\| f^{(n)} \right\|_{\infty} h^{n-1}.$$

**Demo: Taking Derivatives with Vandermonde Matrices** [cleared]

(Accuracy)

## Differentiation Matrices

How can numerical differentiation be cast as a matrix-vector operation?

Let

$$V' = \begin{bmatrix} \varphi_1'(x_1) & \cdots & \varphi_n'(x_1) \\ \vdots & \ddots & \vdots \\ \varphi_1'(x_n) & \cdots & \varphi_n'(x_n) \end{bmatrix}.$$

Then altogether:

$$f'(\mathbf{x}) \approx p_{n-1}(\mathbf{x}) = V'\boldsymbol{\alpha} = V'V^{-1}f(\mathbf{x}).$$

So  $D = V'V^{-1}$  acts as a differentiation matrix.

[Demo: Taking Derivatives with Vandermonde Matrices \[cleared\]](#) (Build  $D$ )

## Properties of Differentiation Matrices

How do I find second derivatives?

$$D^2.$$

Does  $D$  have a nullspace?

- ▶ Yes, constant vectors.  
(At least for polynomial interpolation bases.)
- ▶ I.e. rows of differentiation matrices always sum to 0.

**Demo:** [Taking Derivatives with Vandermonde Matrices \[cleared\]](#) (Shifting and scaling the nodes)

## Numerical Differentiation: Shift and Scale

Does  $D$  change if we shift the nodes  $(x_i)_{i=1}^n \rightarrow (x_i + c)_{i=1}^n$ ?

Let  $\tilde{f}(x) = f(x - c)$ . Define  $\tilde{p}_{n-1}$  via  $\tilde{p}_{n-1}(\mathbf{x} + c) = \tilde{f}(\mathbf{x} + c)$ . Then  $\tilde{p}_{n-1}(x) = p_{n-1}(x - c)$  for all  $x$  because polynomial bases of degree  $\leq n - 1$  are closed under translation, i.e. a shifted basis again consists of polynomials of degree  $\leq n - 1$ . Thus  $\tilde{p}'_{n-1}(\mathbf{x} + c) = p'_{n-1}(\mathbf{x})$ .

In other words,  $D_{\mathbf{x}} = D_{\mathbf{x}+c}$ .

Does  $D$  change if we scale the nodes  $(x_i)_{i=1}^n \rightarrow (\alpha x_i)_{i=1}^n$ ?

Let  $\tilde{f}(x) = f(x/\alpha)$ . Define  $\tilde{p}_{n-1}$  via  $\tilde{p}_{n-1}(\alpha \mathbf{x}) = \tilde{f}(\alpha \mathbf{x})$ . Then  $\tilde{p}_{n-1}(x) = p_{n-1}(x/\alpha)$  for all  $x$  because polynomial bases of degree  $\leq n - 1$  are closed under dilation, i.e. a dilated basis again consists of polynomials of degree  $\leq n - 1$ . Thus  $\tilde{p}'_{n-1}(\alpha \mathbf{x}) = p'_{n-1}(\mathbf{x})/\alpha$ .

In other words,  $D_{\alpha \mathbf{x}} = D_{\mathbf{x}}/\alpha$ .

## Finite Difference Formulas from Diff. Matrices

How do the rows of a differentiation matrix relate to FD formulas?

Let  $D = (d_{i,j})_{i,j=1}^n$ . Then  $f'(x_i) \approx \sum_{j=1}^n d_{i,j} f(x_j)$ .

For example, if  $D$  is  $3 \times 3$ , then ...

- ▶ first row:  $f'(x_1) \approx d_{1,1}f(x_1) + d_{1,2}f(x_2) + d_{1,3}f(x_3)$ ,
- ▶ second row:  $f'(x_2) \approx d_{2,1}f(x_1) + d_{2,2}f(x_2) + d_{2,3}f(x_3)$ , ...

Assume a large equispaced grid and 3 nodes w/same spacing. How to use?

- ▶ First-row formula for left boundary,
- ▶ second-row formula for interior grid points,
- ▶ third-row formula for right boundary.

## Finite Differences: via Taylor

Idea: Start from definition of derivative. Called a **forward difference**.

$$f'(x) \approx \frac{f(x+h) - f(x)}{h}$$

**Q:** What accuracy does this achieve?

Using Taylor:

$$f(x+h) = f(x) + f'(x)h + f''(x)\frac{h^2}{2} + \dots$$

Plug in:

$$\frac{f(x) + f'(x)h + f''(x)\frac{h^2}{2} + \dots - f(x)}{h} = f'(x) + O(h)$$

→ first order accurate.

## More Finite Difference Rules

Similarly:

$$f'(x) = \frac{f(x+h) - f(x-h)}{2h} + O(h^2)$$

(Centered differences)

Can also take higher order derivatives:

$$f''(x) = \frac{f(x+h) - 2f(x) + f(x-h)}{h^2} + O(h^2)$$

Can find these by trying to match Taylor terms.

Alternative: Use linear algebra with interpolate-then-differentiate to find FD formulas.

[Demo: Finite Differences vs Noise \[cleared\]](#)

[Demo: Floating point vs Finite Differences \[cleared\]](#)

## Richardson Extrapolation

Deriving high-order methods is hard work. Can I just do multiple low-order approximations (with different  $h$  and get a high-order one out?

Suppose we have  $F = \tilde{F}(h) + O(h^p)$  and  $\tilde{F}(h_1)$  and  $\tilde{F}(h_2)$ .

Grab one more term of the Taylor series:  $F = \tilde{F}(h) + ah^p + O(h^q)$   
Typically:  $q = p + 1$  (but not necessarily). Do **not** know  $a$ .

**Idea:** Construct new approximation with the goal of  $O(h^q)$  accuracy:

$$F = \alpha\tilde{F}(h_1) + \beta\tilde{F}(h_2) + O(h^q)$$

- ▶ Need  $\alpha ah_1^p + \beta ah_2^p = 0$
- ▶ Need  $\alpha + \beta = 1$  ( $\Leftrightarrow \beta = 1 - \alpha$ ) (maintain low-order terms!)

$$\alpha(h_1^p - h_2^p) + 1h_2^p = 0 \quad \Leftrightarrow \quad \alpha = \frac{-h_2^p}{h_1^p - h_2^p}$$



## Richardson Extrapolation: Observations,

What are  $\alpha$  and  $\beta$  for a first-order (e.g. finite-difference) method if we choose  $h_2 = h_1/2$ ?

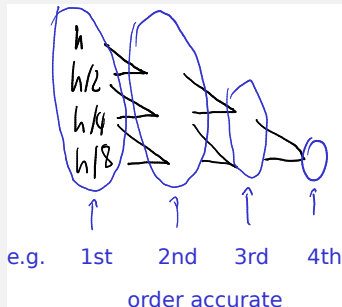
$$p = 1.$$

$$\alpha = \frac{-h_2^p}{h_1^p - h_2^p} = \frac{-\frac{1}{2}}{1 - \frac{1}{2}} = -1, \quad \beta = 1 - \alpha = 2.$$

[Demo: Richardson with Finite Differences \[cleared\]](#)

## Romberg Integration

Can this be used to get *even higher order* accuracy?



Carrying out this process for quadrature is called **Romberg integration**.

## In-Class Activity: Differentiation and Quadrature

In-class activity: Differentiation and Quadrature

# Outline

Introduction to Scientific Computing

Systems of Linear Equations

Linear Least Squares

Eigenvalue Problems

Nonlinear Equations

Optimization

Interpolation

Numerical Integration and Differentiation

**Initial Value Problems for ODEs**

Existence, Uniqueness, Conditioning

Numerical Methods (I)

Accuracy and Stability

Stiffness

Numerical Methods (II)

Boundary Value Problems for ODEs

Partial Differential Equations and Sparse Linear Algebra

Fast Fourier Transform

Additional Topics

## What can we solve already?

- ▶ Linear Systems: **yes**
- ▶ Nonlinear systems: **yes**
- ▶ Systems with derivatives: **no**

## Some Applications

IVPs	BVPs
<ul style="list-style-type: none"><li>▶ Population dynamics <math>y_1' = y_1(\alpha_1 - \beta_1 y_2)</math> (prey) <math>y_2' = y_2(-\alpha_2 + \beta_2 y_1)</math> (predator)</li><li>▶ chemical reactions</li><li>▶ equations of motion</li></ul>	<ul style="list-style-type: none"><li>▶ bridge load</li><li>▶ pollutant concentration (steady state)</li><li>▶ temperature (steady state)</li><li>▶ waves (time-harmonic)</li></ul>

Demo: Predator-Prey System [cleared]

## Initial Value Problems: Problem Statement

Want: Function  $\mathbf{y} : [0, T] \rightarrow \mathbb{R}^n$  so that

- ▶  $\mathbf{y}^{(k)}(t) = \mathbf{f}(t, \mathbf{y}, \mathbf{y}', \mathbf{y}'', \dots, \mathbf{y}^{(k-1)})$  (*explicit*), or
- ▶  $\mathbf{f}(t, \mathbf{y}, \mathbf{y}', \mathbf{y}'', \dots, \mathbf{y}^{(k)}) = 0$  (*implicit*)

are called explicit/implicit  $k$ th-order ordinary differential equations (ODEs).

Give a simple example.

$$y'(t) = \alpha y$$

Not uniquely solvable on its own. What else is needed?

Initial conditions. (Q: How many?)

$$\mathbf{y}(0) = \mathbf{g}_0, \quad \mathbf{y}'(0) = \mathbf{g}_1, \dots, \quad \mathbf{y}^{(k-1)}(0) = \mathbf{g}_{k-1}.$$

Boundary Value Problems (BVPs) trade some derivatives for conditions at the 'other end'.

## Reducing ODEs to First-Order Form

A  $k$ th order ODE can always be reduced to first order. Do this in this example:

$$y''(t) = f(y)$$

In first-order form:

$$\begin{bmatrix} y_1 \\ y_2 \end{bmatrix}'(t) = \begin{bmatrix} y_2(t) \\ f(y_1(t)) \end{bmatrix}$$

Because:

$$y_1''(t) = (y_1'(t))' = y_2'(t) = f(y_1(t)).$$

So we can design our methods to only handle first-order problems.



## Properties of ODEs

What is a **linear** ODE?

$$\mathbf{f}(t, \mathbf{x}) = A(t)\mathbf{x} + \mathbf{b}(t)$$

What is a **linear and homogeneous** ODE?

$$\mathbf{f}(t, \mathbf{x}) = A(t)\mathbf{x}$$

What is a **constant-coefficient** ODE?

$$\mathbf{f}(t, \mathbf{x}) = A\mathbf{x} + \mathbf{b}$$

## Properties of ODEs (II)

What is an **autonomous** ODE?

One in which the function  $f$  does not depend on time  $t$ .

An ODE can be made autonomous by introducing an extra variable:

$$y_0'(t) = 1, \quad y_0(0) = 0.$$

→ Without loss of generality: Get rid of explicit  $t$  dependency.

## Existence and Uniqueness

Consider the perturbed problem

$$\begin{cases} \mathbf{y}'(t) = \mathbf{f}(\mathbf{y}) \\ \mathbf{y}(t_0) = \mathbf{y}_0 \end{cases} \quad \begin{cases} \widehat{\mathbf{y}}'(t) = \mathbf{f}(\widehat{\mathbf{y}}) \\ \widehat{\mathbf{y}}(t_0) = \widehat{\mathbf{y}}_0 \end{cases}$$

Then if  $\mathbf{f}$  is *Lipschitz continuous* (has 'bounded slope'), i.e.

$$\|\mathbf{f}(\mathbf{y}) - \mathbf{f}(\widehat{\mathbf{y}})\| \leq L \|\mathbf{y} - \widehat{\mathbf{y}}\|,$$

- ▶ there exists a solution  $\mathbf{y}$  in a neighborhood of  $t_0$ , and...
- ▶  $\|\mathbf{y}(t) - \widehat{\mathbf{y}}(t)\| \leq e^{L(t-t_0)} \|\mathbf{y}_0 - \widehat{\mathbf{y}}_0\|$
- ▶ This is the **Picard-Lindelöf theorem**.

What does this mean for uniqueness?

It *implies* uniqueness. If there were two separate solutions with identical initial values, they are not allowed to be different.

## Conditioning

Unfortunate terminology accident: “Stability” in ODE-speak

To adapt to conventional terminology, we will use ‘Stability’ for

- ▶ the conditioning of the IVP, *and*
- ▶ the stability of the methods we cook up.

Some terminology:

An IVP is **stable** if and only if...

The solution is continuously dependent on the initial condition, i.e.  
For all  $\varepsilon > 0$  there exists a  $\delta > 0$  so that

$$\|\hat{\mathbf{y}}_0 - \mathbf{y}_0\| < \delta \quad \Rightarrow \quad \|\hat{\mathbf{y}}(t) - \mathbf{y}(t)\| < \varepsilon \quad \text{for all } t \geq t_0.$$

An IVP is **asymptotically stable** if and only if

$$\|\hat{\mathbf{y}}(t) - \mathbf{y}(t)\| \rightarrow 0 \quad (t \rightarrow \infty).$$

## Example I: Scalar, Constant-Coefficient

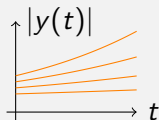
$$\begin{cases} y'(t) = \lambda y \\ y(0) = y_0 \end{cases} \quad \text{where } \lambda = a + ib$$

Solution?

$$y(t) = y_0 e^{\lambda t} = y_0 (e^{at} \cdot e^{ibt})$$

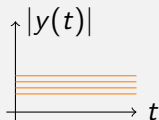
When is this stable?

$\text{Re}(\lambda) > 0$ :



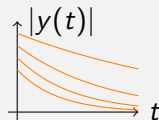
Unstable

$\text{Re}(\lambda) = 0$ :



Stable, not asymptotically stable

$\text{Re}(\lambda) < 0$ :



Asymptotically stable

## Example II: Constant-Coefficient System

$$\begin{cases} \mathbf{y}'(t) = A\mathbf{y}(t) \\ \mathbf{y}(t_0) = \mathbf{y}_0 \end{cases}$$

Assume  $V^{-1}AV = D = \text{diag}(\lambda_1, \dots, \lambda_n)$  diagonal. Find a solution.

Define  $\mathbf{w}(t) := V^{-1}\mathbf{y}(t)$ . Then

$$\mathbf{w}'(t) = V^{-1}\mathbf{y}'(t) = V^{-1}A\mathbf{y}(t) = V^{-1}AV\mathbf{w}(t) = D\mathbf{w}(t).$$

Now:  $n$  decoupled IVPs (with  $\mathbf{w}_0 = V^{-1}\mathbf{y}_0$ )  $\rightarrow$  Solve as scalar.  
Find  $\mathbf{y}(t) = V\mathbf{w}(t)$ .

When is this stable?

When  $\text{Re } \lambda_i \leq 0$  for all eigenvalues  $\lambda_i$ .

## Euler's Method

Discretize the IVP

$$\begin{cases} \mathbf{y}'(t) = \mathbf{f}(\mathbf{y}) \\ \mathbf{y}(t_0) = \mathbf{y}_0 \end{cases}$$

- ▶ Discrete times:  $t_1, t_2, \dots$ , with  $t_{i+1} = t_i + h$
- ▶ Discrete function values:  $\mathbf{y}_k \approx \mathbf{y}(t_k)$ .

**Idea:** Rewrite the IVP in integral form:

$$\mathbf{y}(t) = \mathbf{y}_0 + \int_{t_0}^t \mathbf{f}(\mathbf{y}(\tau)) d\tau,$$

then throw a simple quadrature rule at that. With the rectangle rule, we obtain **Euler's method**.

## Euler's method: Forward and Backward

$$\mathbf{y}(t) = \mathbf{y}_0 + \int_{t_0}^t \mathbf{f}(\mathbf{y}(\tau))d\tau,$$

Use 'left rectangle rule' on integral:

$$\mathbf{y}_{k+1} = \mathbf{y}_k + hf(\mathbf{y}_k)$$

Requires *evaluating the RHS*. Called an **explicit** method. **Forward Euler**.

Use 'right rectangle rule' on integral:

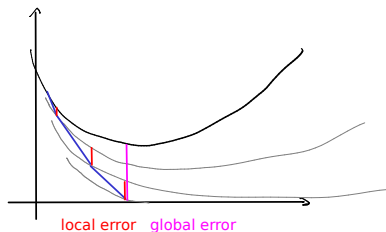
$$\mathbf{y}_{k+1} = \mathbf{y}_k + hf(\mathbf{y}_{k+1})$$

Requires *solving a system of equations*. Called an **implicit** method. **Backward Euler**.

[Demo: Forward Euler stability \[cleared\]](#)



## Global and Local Error



Let  $u_k(t)$  be the function that solves the ODE with the initial condition  $u_k(t_k) = y_k$ . Define the **local error** at step  $k$  as...

$$l_k = y_k - u_{k-1}(t_k)$$

Define the **global error** at step  $k$  as...

$$g_k = y(t_k) - y_k$$

## About Local and Global Error

Is global error =  $\sum$  local errors?

No.

Consider an analogy with interest rates—at any given moment, you receive 5% interest ( $\sim$  incur 5% error) on your current balance.

But your current balance *includes* prior interest (error from prior steps), which yields more interest (in turn contributes to the error).

This contribution to the error is called *propagated error*.

The local error is much easier to estimate  $\rightarrow$  will focus on that.

A time integrator is said to be *accurate of order  $p$*  if...

$$l_k = O(h^{p+1})$$

## ODE IVP Solvers: Order of Accuracy

A time integrator is said to be *accurate of order  $p$*  if  $\ell_k = O(h^{p+1})$

This requirement is one order higher than one might expect—why?

**A:** To get to time 1, at least  $1/h$  steps need to be taken, so that the global error is roughly

$$\underbrace{\frac{1}{h}}_{\text{\#steps}} \cdot O(h^{p+1}) = O(h^p).$$

(Note that this ignores ‘accrual’ of propagated error.)

## Stability of a Method

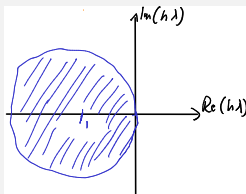
Find out when forward Euler is stable when applied to  $y'(t) = \lambda y(t)$ .

$$\begin{aligned}y_k &= y_{k-1} + h\lambda y_{k-1} \\ &= (1 + h\lambda)y_{k-1} \\ &= (1 + h\lambda)^k y_0\end{aligned}$$

So: stable  $\Leftrightarrow |1 + h\lambda| \leq 1$ .

$|1 + h\lambda|$  is also called the **amplification factor**.

Gives rise to the **stability region** in the complex plane:



## Stability: Systems

What about stability for systems, i.e.

$$\mathbf{y}'(t) = A\mathbf{y}(t)?$$

1. Diagonalize system as before
  2. Notice that same  $V$  also diagonalizes the time stepper
  3. apply scalar analysis to components.
- Stable if  $|1 + h\lambda_i| \leq 1$  for all eigenvalues  $\lambda_i$ .

## Stability: Nonlinear ODEs

What about stability for nonlinear systems, i.e.

$$\mathbf{y}'(t) = \mathbf{f}(\mathbf{y}(t))?$$

Consider perturbation  $\mathbf{e}(t) = \mathbf{y}(t) - \hat{\mathbf{y}}(t)$ . Linearize:

$$\mathbf{e}'(t) = \mathbf{f}(\mathbf{y}(t)) - \mathbf{f}(\hat{\mathbf{y}}(t)) \approx J_{\mathbf{f}}(\mathbf{y}(t))\mathbf{e}(t)$$

i.e. can (at least locally) apply analysis for linear systems to the nonlinear case.

## Stability for Backward Euler

Find out when backward Euler is stable when applied to  $y'(t) = \lambda y(t)$ .

$$\begin{aligned}y_k &= y_{k-1} + h\lambda y_k \\y_k(1 - h\lambda) &= y_{k-1} \\y_k &= \frac{1}{1 - h\lambda} y_{k-1} = \left(\frac{1}{1 - h\lambda}\right)^k y_0.\end{aligned}$$

So: stable  $\Leftrightarrow |1 - h\lambda| \geq 1$ .

In particular: stable for any  $h$  if  $\Re \lambda \leq 0$  (“**unconditionally stable**”).

BE can be stable even when ODE is *unstable*. ( $\Re \lambda > 0$ ). Accuracy?

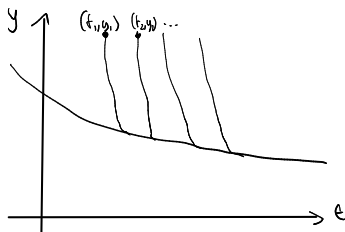
- ▶ *Explicit* methods: main concern in choosing  $h$  is *stability* (but *also* accuracy).
- ▶ *Implicit* methods: main concern in choosing  $h$  is *accuracy*.

## Stiff ODEs: Demo

Demo: Stiffness [cleared]



## 'Stiff' ODEs



- ▶ Stiff problems have *multiple time scales*.  
(In the example above: Fast decay, slow evolution.)
- ▶ In the case of a stable ODE system

$$\mathbf{y}'(t) = \mathbf{f}(\mathbf{y}(t)),$$

stiffness can arise if  $J_f$  has eigenvalues of very different magnitude.

## Stiffness: Observations

Why not just 'small' or 'large' magnitude?

Because the discrepancy between time scales is the root of the problem. If all time scales are similar, then time integration must simply 'deal with' that one time scale.

If there are two, then some (usually the fast ones) may be considered uninteresting.

What is the problem with applying explicit methods to stiff problems?

Fastest time scale governs time step  $\rightarrow$  tiny time step  $\rightarrow$  inefficient.

## Stiffness vs. Methods

Phrase this as a conflict between accuracy and stability.

- ▶ Accuracy (here: capturing the slow time scale) *could* be achieved with large time steps.
- ▶ Stability (in explicit methods) demands a small time step.

Can an implicit method take arbitrarily large time steps?

In terms of stability: sure.  
In terms of accuracy: no.

## Predictor-Corrector Methods

**Idea:** Obtain intermediate result, improve it (with same or different method).

For example:

1. *Predict* with forward Euler:  $\tilde{y}_{k+1} = y_k + hf(y_k)$
2. *Correct* with the trapezoidal rule:  
$$y_{k+1} = y_k + \frac{h}{2}(f(y_k) + f(\tilde{y}_{k+1})).$$

This is called **Heun's method**.

## Runge-Kutta / 'Single-step' / 'Multi-Stage' Methods

**Idea:** Compute intermediate 'stage values', compute new state from those:

$$r_1 = f(t_k + c_1 h, y_k + (a_{11} \cdot r_1 + \dots + a_{1s} \cdot r_s)h)$$

$$\vdots$$

$$r_s = f(t_k + c_s h, y_k + (a_{s1} \cdot r_1 + \dots + a_{ss} \cdot r_s)h)$$

$$y_{k+1} = y_k + (b_1 \cdot r_1 + \dots + b_s \cdot r_s)h$$

Can summarize in a *Butcher tableau*:

$c_1$	$a_{11}$	$\dots$	$a_{1s}$
$\vdots$	$\vdots$		$\vdots$
$c_s$	$a_{s1}$	$\dots$	$a_{ss}$
	$b_1$	$\dots$	$b_s$

## Runge-Kutta: Properties

When is an RK method explicit?

If the diagonal entries in the Butcher tableau and everything above it are zero.

When is it implicit?

(Otherwise)

When is it *diagonally implicit*? (And what does that mean?)

If the everything above the diagonal entries in the Butcher tableau is zero.

This means that one can solve for one stage value at a time (and not multiple).

## Runge-Kutta: Embedded Pairs

How can error in RK integration be controlled?

Most of the cost is in computing stage values  $r_1, \dots, r_s$ . Reuse for a second (order  $p^*$  accurate) state estimate:

$$y_{k+1} = y_k + (b_1 \cdot r_1 + \dots + b_s \cdot r_s)h$$

$$y_{k+1}^* = y_k + (b_1^* \cdot r_1 + \dots + b_s^* \cdot r_s)h$$

$|y_{k+1} - y_{k+1}^*|$  can serve as an estimate of local error  $\ell_{k+1}$ , e.g. for time step control. Called an **embedded pair**.

$c_1$	$a_{11}$	$\dots$	$a_{1s}$
$\vdots$	$\vdots$		$\vdots$
$c_s$	$a_{s1}$	$\dots$	$a_{ss}$
$p$	$b_1$	$\dots$	$b_s$
$p^*$	$b_1^*$	$\dots$	$b_s^*$

## Heun and Butcher

Stuff Heun's method into a Butcher tableau:

1.  $\tilde{y}_{k+1} = y_k + hf(y_k)$
2.  $y_{k+1} = y_k + \frac{h}{2}(f(y_k) + f(\tilde{y}_{k+1}))$ .

0			
1		1	
-----		$\frac{1}{2}$	$\frac{1}{2}$



# RK4

What is RK4?

0				
$\frac{1}{2}$	$\frac{1}{2}$			
$\frac{1}{2}$	0	$\frac{1}{2}$		
1	0	0	1	
<hr/>				
	$\frac{1}{6}$	$\frac{1}{3}$	$\frac{1}{3}$	$\frac{1}{6}$

Note similarity to  
Simpson's rule!

$$k_1 = f(t_n, y_n),$$

$$k_2 = f\left(t_n + \frac{h}{2}, y_n + h\frac{k_1}{2}\right),$$

$$k_3 = f\left(t_n + \frac{h}{2}, y_n + h\frac{k_2}{2}\right),$$

$$k_4 = f(t_n + h, y_n + hk_3)$$

$$y_{n+1} = y_n + \frac{1}{6}h(k_1 + 2k_2 + 2k_3 + k_4),$$

[Demo: Dissipation in Runge-Kutta Methods \[cleared\]](#)

## Multi-step/Single-stage/Adams Methods/Backward Differencing Formulas (BDFs)

**Idea:** Instead of computing stage values, use *history* (of either values of  $f$  or  $y$ —or both):

$$y_{k+1} = \sum_{i=1}^M \alpha_i y_{k+1-i} + h \sum_{i=1}^N \beta_i f(y_{k+1-i})$$

Extensions to implicit possible.

Method relies on existence of history. What if there isn't any? (Such as at the start of time integration?)

These methods are *not self-starting*.

Need another method to produce enough history.

## Stability Regions

Why does the idea of stability regions still apply to more complex time integrators (e.g. RK?)

As long as the method doesn't "treat individual vector entries specially", a matrix that diagonalizes the ODE also diagonalizes the time integrator.

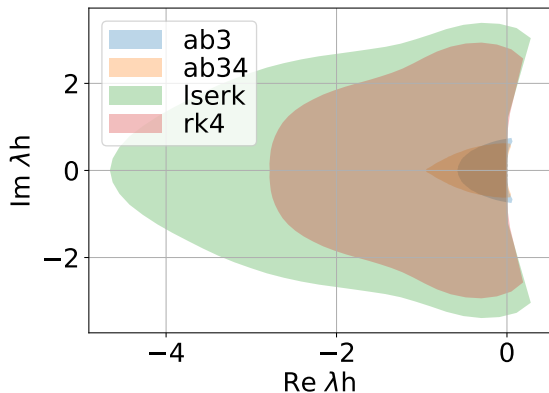
⇒ Can consider stability one eigenvalue at a time.

[Demo: Stability regions \[cleared\]](#)

## More Advanced Methods

Discuss:

- ▶ What is a good cost metric for time integrators?
- ▶ AB3 vs RK4
- ▶ Runge-Kutta-Chebyshev
- ▶ [LSERK](#) and [AB34](#)
- ▶ IMEX and multi-rate
- ▶ Parallel-in-time (["Parareal"](#))



## In-Class Activity: Initial Value Problems

In-class activity: Initial Value Problems

# Outline

Introduction to Scientific Computing

Systems of Linear Equations

Linear Least Squares

Eigenvalue Problems

Nonlinear Equations

Optimization

Interpolation

Numerical Integration and Differentiation

Initial Value Problems for ODEs

**Boundary Value Problems for ODEs**

**Existence, Uniqueness, Conditioning**  
**Numerical Methods**

Partial Differential Equations and Sparse Linear Algebra

Fast Fourier Transform

Additional Topics

## BVP Problem Setup: Second Order

Example: Second-order linear ODE

$$u''(x) + p(x)u'(x) + q(x)u(x) = r(x)$$

with *boundary conditions* ('BCs') at  $a$ :

- ▶ *Dirichlet*  $u(a) = u_a$
- ▶ or *Neumann*  $u'(a) = v_a$
- ▶ or *Robin*  $\alpha u(a) + \beta u'(a) = w_a$

and the same choices for the BC at  $b$ .

*Note:* BVPs in time are rare in applications, hence  $x$  (not  $t$ ) is typically used for the independent variable.

## BVP Problem Setup: General Case

ODE:

$$\mathbf{y}'(x) = \mathbf{f}(\mathbf{y}(x)) \quad \mathbf{f} : \mathbb{R}^n \rightarrow \mathbb{R}^n$$

BCs:

$$\mathbf{g}(\mathbf{y}(a), \mathbf{y}(b)) = 0 \quad \mathbf{g} : \mathbb{R}^{2n} \rightarrow \mathbb{R}^n$$

(Recall the rewriting procedure to first-order for any-order ODEs.)

Does a first-order, scalar BVP make sense?

No—need second order (or  $n \geq 2$ ) to allow two boundary conditions.

**Example:** Linear BCs  $B_a \mathbf{y}(a) + B_b \mathbf{y}(b) = \mathbf{c}$ .

Is this Dirichlet/Neumann/...?

Could be any—we're in the system case, and  $B_a$  and  $B_b$  are matrices—so conditions could be any *any* component.



## Do solutions even exist? How sensitive are they?

General case is harder than root finding, and we couldn't say much there.

→ Only consider linear BVP.

$$(*) \begin{cases} \mathbf{y}'(x) = A(x)\mathbf{y}(x) + \mathbf{b}(x) \\ B_a\mathbf{y}(a) + B_b\mathbf{y}(b) = \mathbf{c} \end{cases}$$

Exploit linearity: split into multiple problems.

Split into boundary (B) and volume (V) parts.

$$(B) \begin{cases} \mathbf{y}'_B(x) = A(x)\mathbf{y}_B(x) \\ B_a\mathbf{y}_B(a) + B_b\mathbf{y}_B(b) = \mathbf{c} \end{cases}$$
$$(V) \begin{cases} \mathbf{y}'_V(x) = A(x)\mathbf{y}_V(x) + \mathbf{b}(x) \\ B_a\mathbf{y}_V(a) + B_b\mathbf{y}_V(b) = 0 \end{cases}$$

Then  $\mathbf{y} = \mathbf{y}_B + \mathbf{y}_V$ .

## Solving the “Boundary” BVP

$$(B) \begin{cases} \mathbf{y}'_B(x) = A(x)\mathbf{y}_B(x) \\ B_a\mathbf{y}_B(a) + B_b\mathbf{y}_B(b) = \mathbf{c} \end{cases}$$

$$\mathbf{y}'_{B,i}(x) = A(x)\mathbf{y}_{B,i}(x), \quad \mathbf{y}_{B,i}(a) = \mathbf{e}_i. \quad (i = 1, \dots, n)$$

$\mathbf{e}_i$  is the  $i$ th unit vector. Define **fundamental sol. matrix**:

$$Y(x) = \begin{bmatrix} | & & | \\ \mathbf{y}_{B,1} & \cdots & \mathbf{y}_{B,n} \\ | & & | \end{bmatrix}$$

Let  $Q := B_a Y(a) + B_b Y(b)$ . (\*) has a unique solution  $\Leftrightarrow Q$  is invertible. Solve  $Q\alpha = \mathbf{c}$  to find coefficients:  $Y(x)\alpha$  solves (B). Define  $\Phi(x) := Y(x)Q^{-1}$ .  $\Phi(x)\mathbf{c}$  also solves (B).

## Solving the “Volume” BVP

$$(V) \begin{cases} \mathbf{y}'_V(x) = A(x)\mathbf{y}_V(x) + \mathbf{b}(x) \\ B_a\mathbf{y}_V(a) + B_b\mathbf{y}_V(b) = 0 \end{cases}$$

Define **Green's function**

$$G(x, z) := \begin{cases} \Phi(x)B_a\Phi(a)\Phi^{-1}(z) & z \leq x, \\ -\Phi(x)B_b\Phi(b)\Phi^{-1}(z) & z > x. \end{cases}$$

Then

$$\mathbf{y}_V(x) = \int_a^b G(x, y)\mathbf{b}(z)dz$$

solves (V).

## ODE Systems: Conditioning

Altogether:

$$\mathbf{y}(x) = \mathbf{y}_B + \mathbf{y}_V = \Phi(x)\mathbf{c} + \int_a^b G(x, y)\mathbf{b}(y)dy.$$

For perturbed problem with  $\mathbf{b}(x) + \Delta\mathbf{b}(x)$  and  $\mathbf{c} + \Delta\mathbf{c}$ , derive a bound on  $\|\Delta\mathbf{y}\|_\infty$ .

$$\|\Delta\mathbf{y}\|_\infty \leq \max(\|\Phi\|_\infty, \|G\|_\infty) \left( \|\Delta\mathbf{c}\|_1 + \int \|\Delta\mathbf{b}(y)\|_1 dy \right).$$

- ▶ Conditioning bound implies uniqueness.
- ▶ Also get continuous dependence on data.

## Shooting Method

**Idea:** Want to make use of the fact that we can already solve IVPs.

**Problem:** Don't know *all* left BCs.

**Demo:** Shooting method [cleared]

What about systems?

No problem—cannons are aimed in 2D as well. :)

What are some downsides of this method?

- ▶ Can fail
- ▶ Can be unstable even if ODE is stable

What's an alternative approach?

Set up a big linear system.

## Finite Difference Method

**Idea:** Replace  $u'$  and  $u''$  with finite differences.

**For example:** second-order centered

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

**Demo:** [Finite differences](#) [cleared]

What happens for a nonlinear ODE?

Get a nonlinear system → Use Newton.

**Demo:** [Sparse matrices](#) [cleared]

## Collocation Method

$$(*) \begin{cases} y'(x) = f(y(x)), \\ g(y(a), y(b)) = 0. \end{cases}$$

1. Pick a basis (for example: Chebyshev polynomials)

$$\hat{y}(x) = \sum_{i=1}^n \alpha_i T_i(x)$$

Want  $\hat{y}$  to be close to solution  $y$ . So: plug into  $(*)$ .

**Problem:**  $\hat{y}$  won't satisfy the ODE at all points at least.  
We do not have enough unknowns for that.

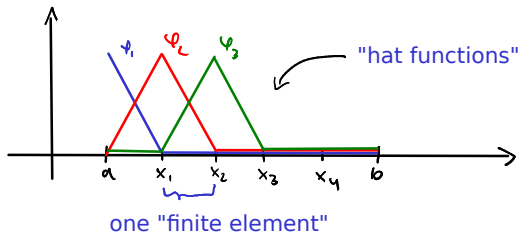
2. **Idea:** Pick  $n$  points where we would like  $(*)$  to be satisfied.  
→ Get a big (non-)linear system
3. Solve that (LU/Newton) → done.

# Galerkin/Finite Element Method

$$u''(x) = f(x), \quad u(a) = u(b) = 0.$$

**Problem** with collocation: Big dense matrix.

**Idea:** Use piecewise basis. Maybe it'll be sparse.



What's the problem with that?

$u'$  does not exist. (at least at a few points where it's discontinuous)  
 $u''$  really does not exist.



## Weak solutions/Weighted Residual Method

**Idea:** Enforce a 'weaker' version of the ODE.

Compute 'moments':

$$\int_a^b u''(x)\psi(x)dx = \int_a^b f(x)\psi(x)dx$$

Require that this holds for some *test functions*  $\psi$  from some set  $W$ .  
Now possible to get rid of (undefined) second derivative using integration by parts:

$$\int_a^b u''(x)\psi(x)dx = [u'(x)\psi(x)]_a^b - \int_a^b u'(x)\psi'(x)dx.$$

- ▶ Also called **weighted residual** methods.
- ▶ Can view collocation as a WR method with  $\psi_j(x) = \delta(x - x_j)$

## Galerkin: Choices in Weak Solutions

Make some choices:

- ▶ Solve for  $u \in \text{span} \{\text{hat functions } \varphi_i\}$
- ▶ Choose  $\psi \in W = \text{span} \{\text{hat functions } \varphi_i\}$  with  $\psi(a) = \psi(b) = 0$ .  
→ Kills boundary term  $[u'(x)\psi(x)]_a^b$ .

These choices are called the **Galerkin method**. Also works with other bases.

## Discrete Galerkin

Assemble a matrix for the Galerkin method.

$$\begin{aligned} - \int_a^b u'(x) \psi'(x) dx &= \int_a^b f(x) \psi(x) dx \\ - \int_a^b \left[ \sum_{j=1}^n \alpha_j \varphi_j'(x) \right] \psi'(x) dx &= \int_a^b f(x) \psi(x) dx \\ - \sum_{j=1}^n \alpha_j \underbrace{\int_a^b \varphi_j'(x) \varphi_i'(x) dx}_{S_{ij}} &= \underbrace{\int_a^b f(x) \varphi_i(x) dx}_{r_i} \\ S \alpha &= r. \end{aligned}$$

**Now:** Compute  $S$ , solve sparse (!) linear system.

# Outline

Introduction to Scientific Computing

Systems of Linear Equations

Linear Least Squares

Eigenvalue Problems

Nonlinear Equations

Optimization

Interpolation

Numerical Integration and Differentiation

Initial Value Problems for ODEs

Boundary Value Problems for ODEs

**Partial Differential Equations and Sparse Linear Algebra**

**Sparse Linear Algebra**  
**PDEs**

Fast Fourier Transform

Additional Topics

## Advertisement

**Remark:** Both PDEs and Large Scale Linear Algebra are big topics. Will only scratch the surface here. Want to know more?

- ▶ CS555 → Numerical Methods for PDEs
- ▶ CS556 → Iterative and Multigrid Methods
- ▶ CS554 → Parallel Numerical Algorithms

We would love to see you there! :)

# Solving Sparse Linear Systems

Solving  $A\mathbf{x} = \mathbf{b}$  has been our bread and butter.

Typical approach: Use factorization (like LU or Cholesky)

Why is this problematic?

**Idea:** Don't factorize, iterate.

**Demo: Sparse Matrix Factorizations and "Fill-In" [cleared]**

## 'Stationary' Iterative Methods

**Idea:** Invert only part of the matrix in each iteration. Split

$$A = M - N,$$

where  $M$  is the part that we are actually inverting. Convergence?

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

$$M\mathbf{x} = N\mathbf{x} + \mathbf{b}$$

$$M\mathbf{x}_{k+1} = N\mathbf{x}_k + \mathbf{b}$$

$$\mathbf{x}_{k+1} = M^{-1}(N\mathbf{x}_k + \mathbf{b})$$

- ▶ These methods are called *stationary* because they do the same thing in every iteration.
- ▶ They carry out fixed point iteration.  
→ Converge if contractive, i.e.  $\rho(M^{-1}N) < 1$ .
- ▶ Choose  $M$  so that it's easy to invert.

## Choices in Stationary Iterative Methods

What could we choose for  $M$  (so that it's easy to invert)?

Name	$M$	$N$
Jacobi	$D$	$-(L + U)$
Gauss-Seidel	$D + L$	$-U$
SOR	$\frac{1}{\omega}D + L$	$(\frac{1}{\omega} - 1)D - U$

where  $L$  is the below-diagonal part of  $A$ , and  $U$  the above-diagonal.

[Demo: Stationary Methods \[cleared\]](#)



## Conjugate Gradient Method

Assume  $A$  is symmetric positive definite.

**Idea:** View solving  $A\mathbf{x} = \mathbf{b}$  as an optimization problem.

$$\text{Minimize } \varphi(\mathbf{x}) = \frac{1}{2}\mathbf{x}^T A\mathbf{x} - \mathbf{x}^T \mathbf{b} \quad \Leftrightarrow \quad \text{Solve } A\mathbf{x} = \mathbf{b}.$$

Observe  $-\nabla\varphi(\mathbf{x}) = \mathbf{b} - A\mathbf{x} = \mathbf{r}$  (residual).

Use an iterative procedure ( $\mathbf{s}_k$  is the search direction):

$$\begin{aligned}\mathbf{x}_0 &= \langle \text{starting vector} \rangle \\ \mathbf{x}_{k+1} &= \mathbf{x}_k + \alpha_k \mathbf{s}_k,\end{aligned}$$

## CG: Choosing the Step Size

What should we choose for  $\alpha_k$  (assuming we know  $\mathbf{s}_k$ )?

$$\begin{aligned} 0 &\stackrel{!}{=} \frac{\partial}{\partial \alpha} \varphi(\mathbf{x}_k + \alpha_k \mathbf{s}_k) \\ &= \nabla \varphi(\mathbf{x}_{k+1}) \cdot \mathbf{s}_k = -\mathbf{r}_{k+1} \cdot \mathbf{s}_k. \end{aligned}$$

**Learned:** Choose  $\alpha$  so that next residual is  $\perp$  to current search direction.

$$\begin{aligned} \mathbf{r}_{k+1} &= \mathbf{r}_k - \alpha_k A \mathbf{s}_k \\ 0 &\stackrel{!}{=} \mathbf{s}_k^T \mathbf{r}_{k+1} = \mathbf{s}_k^T \mathbf{r}_k - \alpha_k \mathbf{s}_k^T A \mathbf{s}_k \end{aligned}$$

Solve:

$$\alpha_k = \frac{\mathbf{s}_k^T \mathbf{r}_k}{\mathbf{s}_k^T A \mathbf{s}_k} = -\frac{\mathbf{s}_k^T A \mathbf{e}_k}{\mathbf{s}_k^T A \mathbf{s}_k}, \quad (*)$$

where  $\mathbf{e}_k = \mathbf{x}_k - \mathbf{x}^*$  and  $\mathbf{r}_k = -A \mathbf{e}_k$ .

## CG: Choosing the Search Direction

What should we choose for  $\mathbf{s}_k$ ?

**Idea:**  $\mathbf{s}_k = \mathbf{r}_k = -\nabla\varphi(\mathbf{x}_k)$ , i.e. steepest descent. No—still a bad idea.

$\mathbf{x}$ ,  $\mathbf{y}$  are called *A-orthogonal* or *conjugate* if and only if  $\mathbf{x}^T \mathbf{A} \mathbf{y} = 0$ .

**Better Idea:** Require  $\mathbf{s}_i^T \mathbf{A} \mathbf{s}_j = 0$  if  $i \neq j$ .

View error as linear combination of search directions, with some (thus far unknown) coefficients:

$$\mathbf{e}_0 = \mathbf{x}_0 - \mathbf{x}^* = \sum_i \delta_i \mathbf{s}_i.$$

- ▶ We run out of *A-orthogonal* directions after  $n$  iterations.
- ▶ Is the error going to be zero then? If  $\delta_k = -\alpha_k$ , then yes.

## CG: Further Development

$$\mathbf{s}_k^T \mathbf{A} \mathbf{e}_0 = \sum_i \delta_i \mathbf{s}_k^T \mathbf{A} \mathbf{s}_i = \delta_k \mathbf{s}_k^T \mathbf{A} \mathbf{s}_k.$$

Solve for  $\delta_k$  and expand:

$$\delta_k = \frac{\mathbf{s}_k^T \mathbf{A} \mathbf{e}_0}{\mathbf{s}_k^T \mathbf{A} \mathbf{s}_k} = \frac{\mathbf{s}_k^T \mathbf{A} \left( \mathbf{e}_0 + \sum_{i=1}^{k-1} \alpha_i \mathbf{s}_i \right)}{\mathbf{s}_k^T \mathbf{A} \mathbf{s}_k} = \frac{\mathbf{s}_k^T \mathbf{A} \mathbf{e}_k}{\mathbf{s}_k^T \mathbf{A} \mathbf{s}_k} = -\alpha_k.$$

How do we generate the  $\mathbf{s}_k$ ?

- ▶ Pick a random one to start with. Perhaps  $\mathbf{r}_0$ ?
- ▶ Generate next one by orthogonalizing from Krylov space procedure  $\mathbf{z}, \mathbf{A}\mathbf{z}, \mathbf{A}^2\mathbf{z}$

**Insight:** Use three-term Lanczos it. to generate.  $\rightarrow$  cheap!

# Introduction

Notation:

$$\frac{\partial}{\partial x} u = \partial_x u = u_x.$$

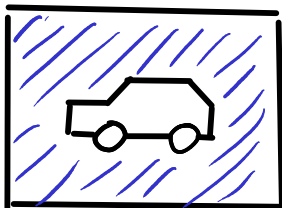
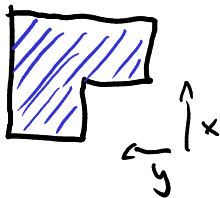
A *PDE* (*partial differential equation*) is an equation with multiple partial derivatives:

$$u_{xx} + u_{yy} = 0$$

Here: solution is a function  $u(x, y)$  of two variables.

**Examples:** Wave propagation, fluid flow, heat diffusion

- ▶ Typical: Solve on domain with complicated geometry.



# Initial and Boundary Conditions

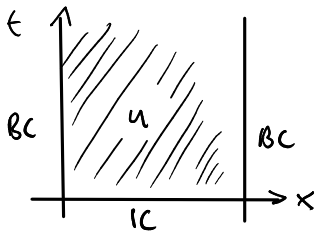
- ▶ Sometimes one variable is time-like.

What makes a variable time-like?

- ▶ Causality
- ▶ No geometry

Have:

- ▶ PDE
- ▶ Boundary conditions
- ▶ Initial conditions (in  $t$ )



## Time-Dependent PDEs

Time-dependent PDEs give rise to a *steady-state* PDE:

$$u_t = f(u_x, u_y, u_{xx}, u_{yy}) \quad \rightarrow \quad 0 = f(u_x, u_y, u_{xx}, u_{yy})$$

Idea for time-dep problems (**Method of Lines**):

- ▶ Discretize spatial derivatives first
- ▶ Obtain large (**semidiscrete**) system of ODEs
- ▶ Use ODE solver from Chapter 9

**Demo: Time-dependent PDEs [cleared]**

## Notation: Laplacian

Laplacian (dimension-independent)

$$\Delta u = \operatorname{div} \operatorname{grad} u = \nabla \cdot (\nabla u) = u_{xx} + u_{yy}$$



# Classifying PDEs

Three main types of PDEs:

- ▶ **hyperbolic** (wave-like, conserve energy)
  - ▶ first-order **conservation laws**:  $u_t + f(u)_x = 0$
  - ▶ second-order **wave equation**:  $u_{tt} = \Delta u$
- ▶ **parabolic** (heat-like, dissipate energy)
  - ▶ **heat equation**:  $u_t = \Delta u$
- ▶ **elliptic** (steady-state, of heat and wave eq. for example)
  - ▶ **Laplace equation**  $\Delta u = 0$
  - ▶ **Poisson equation**  $\Delta u = f$   
(Pure BVP, similar to 1D BVPs, same methods apply—FD, Galerkin, etc.)

# Outline

Introduction to Scientific Computing

Systems of Linear Equations

Linear Least Squares

Eigenvalue Problems

Nonlinear Equations

Optimization

Interpolation

Numerical Integration and Differentiation

Initial Value Problems for ODEs

Boundary Value Problems for ODEs

Partial Differential Equations and Sparse Linear Algebra

**Fast Fourier Transform**

Additional Topics

# Outline

Introduction to Scientific Computing

Systems of Linear Equations

Linear Least Squares

Eigenvalue Problems

Nonlinear Equations

Optimization

Interpolation

Numerical Integration and Differentiation

Initial Value Problems for ODEs

Boundary Value Problems for ODEs

Partial Differential Equations and Sparse Linear Algebra

Fast Fourier Transform

**Additional Topics**