

## Review for Midterm 2

# Contents

- ① Preliminaries (3 lec.)
- F-P numbers
  - Conditioning (problem)
  - Stability (algorithm)
- ② Square Linear Systems (5 lect.)
- $A \vec{x} = \vec{b}$  where  $A \in \mathbb{R}^{n \times n}$
  - polynomial interpolation. (vandermonde)
  - Gaussian Elim.  $\rightarrow$  LU factorization
- ③ Overdetermined Linear Systems (5 lect.) (with partial pivoting) (P)
- $A \vec{x} = \vec{b}$  where  $A \in \mathbb{R}^{m \times n}$ ,  $m > n$
  - polynom. approx. (LLS)
  - Normal eqn ( $A^T A \vec{x} = A^T \vec{b}$ )  $\rightarrow$  QR factorization.

# Preliminaries

# Two Types of Errors

- absolute error
- relative error

# Floating-Point Numbers

- binary scientific notation:

$$\pm \left( 1 + \frac{b_1}{2} + \frac{b_2}{2^2} + \cdots + \frac{b_d}{2^d} \right) 2^E,$$

where  $b_i$  is 0 or 1 and  $E$  is an integer.

- $d$  determines the *resolution*
- the range of  $E$  determines the *scope* or *extent*
- IEEE Standard (double-precision; 64 bits)
  - $d = 52$  and  $-1022 \leq E \leq 1023$
  - $\boxed{\text{eps}} = 2^{-52} \approx 2 \times 10^{-16}$
  - `realmin`, `realmax`

# Floating-Point Numbers (cont')

- Key features
  - On any interval of the form  $[2^E, 2^{E+1})$ , there are  $2^d$  evenly-spaced f-p numbers.
  - The spacing between two adjacent f-p numbers in  $[2^E, 2^{E+1})$  is  $2^{E-d} = 2^E \boxed{\text{eps}}$ .
  - The gap between 1 and the next f-p number is  $\boxed{\text{eps}}$ , the machine epsilon.
  - Representation error (in relative sense) is bounded by  $\frac{1}{2} \boxed{\text{eps}}$ .

# Conditioning (of a problem)

- The condition number measures the ratio of error in the result (or output) to error in the data (or input).
- Recall the definition of condition number  $\kappa_f(x)$
- A large condition number implies that the error in a result may be much greater than the round-off error used to compute it.
- Catastrophic cancellation is one of the most common sources of loss of precision.

## Stability (of an algorithm)

- When an algorithm produces much more error than can be explained by the condition number, the algorithm is unstable.



# Square Linear Systems

# Polynomial Interpolation

- Polynomial interpolation leads to a square linear system of equations with a Vandermonde matrix.

# Gaussian Elimination and (P)LU Factorization

- A triangular linear system is solved by backward substitution or forward elimination.
- A general linear system is solved by Gaussian elimination.
- Gaussian elimination (with partial pivoting) is equivalent to (P)LU factorization.
- Solving a triangular linear system of size  $n \times n$  takes  $\sim n^2$  flops.
- PLU factorization takes  $\sim \frac{2}{3}n^3$  flops.

# Norms

A *norm* generalizes the notion of length for vectors and matrices.

- **Vector  $p$ -norm**

$$\|\mathbf{v}\|_p = \left( \sum_{i=1}^n |b_i|^p \right)^{1/p}, \quad p \in [1, \infty)$$

and

$$\|\mathbf{v}\|_\infty = \max_i |v_i|$$

cf. HW 6 #5

- **Matrix  $p$ -norm (induced)**



$$\|A\|_p = \max_{\|\mathbf{x}\|_p=1} \|A\mathbf{x}\|_p, \quad p \in [1, \infty]$$

- **Frobenius norm (non-induced)**

$$\|A\|_F = \left( \sum_i \sum_j |a_{i,j}|^2 \right)^{1/2}$$

- **MATLAB:** `norm` can calculate both vector and matrix norms

# Row and Column Operations

Various row and column operations can be emulated by matrix multiplications.  
("Left-multiplication for row actions, right-multiplication for column actions")

- row/column extraction (unit vector)
- row/column swap (elementary permutation matrix)
- row/column rearrangement (permutation matrix)
- row replacement  $R_i \rightarrow R_i + cR_j$  (Gaussian transformation matrix)

# Conditioning/Stability

- Partial pivoting is needed for numerical stability.
- The matrix condition number is equal to the condition number of solving a linear system of equations.

# Programming Notes

- Built-in functionalities

- backslash (\)

$A \vec{x} = \vec{b}$  is solved via  $x = A \setminus b$ .

- lu

- norm

- cond, condest, linsolve

$\text{cond}(A, p)$  computes  $\kappa_p(A) = \|A\|_p \|A^{-1}\|_p$ , matrix condition number.

- Demonstration/Instructional codes

- backsub and forelim

- GENp and GEpp

- mylu and myplu

# Overdetermined Linear Systems



# Polynomial Approximation

- The most common solution to overdetermined systems is obtained by *least squares*, which minimizes the 2-norm of the residual vector.
- Least squares is used to find fitting functions that depend linearly on the unknown parameters.
- Equivalence of the LLS problem and the normal equation
  - linear algebra proof
  - calculus proof

# QR Factorization

- Orthogonal sets of vectors are preferred to nonorthogonal ones in computing. (no catastrophic cancellation)
- Matrices with orthonormal columns and orthogonal matrices enjoy many *nice* analytical properties.
- QR factorization plays a role in LLS similar to that of LU factorization in square linear systems.

# Two Types of QR Factorization

For  $A \in \mathbb{R}^{m \times n}$ ,  $m \geq n$ :

orthonormal columns

- Thick QR factorization:  $A = QR$

- $Q \in \mathbb{R}^{m \times m}$  orthogonal (square matrix,  $Q^T Q = I$ )
- $R \in \mathbb{R}^{m \times n}$  upper triangular
- obtained by using successive Householder transformation matrices for triangularization

- Thin:  $A = \hat{Q}\hat{R}$

- $\hat{Q} \in \mathbb{R}^{m \times n}$  orthonormal columns ( $\hat{Q}^T \hat{Q} = I$ )
- $\hat{R} \in \mathbb{R}^{n \times n}$  upper triangular
- obtained by Gram-Schmidt orthonormalization procedure

$\hat{Q}$  is a matrix w/  
ONC, but is not an  
orthogonal matrix!

# Householder Transformation Matrices

- A Householder transformation matrix  $H$  (associated with a vector  $\mathbf{z}$ ) is a reflection matrix which is
  - symmetric,
  - orthogonal, and
  - transforms  $\mathbf{z}$  to  $\pm \|\mathbf{z}\|_2 \mathbf{e}_1$ .     i.e.,  $H\vec{z} = \|\vec{z}\|_2 \vec{e}_1$

HW7 #3 (c) Let  $\vec{v} = \|\vec{z}\|_2 \vec{e}_1 - \vec{z}$  and

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

Show that  $H\vec{z} = \|\vec{z}\| \vec{e}_1$ .

---

Soln Note that

$$\begin{aligned} \vec{v}^T \vec{z} &= (\|\vec{z}\| \vec{e}_1 - \vec{z})^T \vec{z} \\ &= (\|\vec{z}\| \vec{e}_1^T - \vec{z}^T) \vec{z} \\ &= \|\vec{z}\| \vec{e}_1^T \vec{z} - \underbrace{\vec{z}^T \vec{z}}_{\|\vec{z}\|^2} \\ &= \|\vec{z}\| \vec{e}_1^T \vec{z} - \|\vec{z}\|^2 \end{aligned}$$

orth. proj. onto  $\langle \vec{v} \rangle$

Scratch

$$\begin{aligned} H\vec{z} &= \left( I - 2 \frac{\vec{v} \vec{v}^T}{\vec{v}^T \vec{v}} \right) \vec{z} \\ &= \vec{z} - 2 \frac{\vec{v} \vec{v}^T \vec{z}}{\vec{v}^T \vec{v}} \end{aligned}$$

# Programming Notes

- Built-in functionalities
  - backslash (\)
  - qr
- Demonstration/Instructional codes
  - `lsqrfact`: solving least squares using QR
  - `gs`: Gram-Schmidt (for homework)

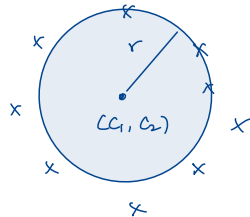
Fitting a circle Given data  $\{(x_k, y_k) : k=1, \dots, m\}$ .

Goal Find  $a_1, a_2$ , and  $r$  as in

$$(x - a_1)^2 + (y - a_2)^2 = r^2$$

---

$$x^2 - 2x a_1 + a_1^2 + y^2 - 2y a_2 + a_2^2 = r^2$$



$$x^2 + y^2 = 2x a_1 + 2y a_2 + \underbrace{r^2 - a_1^2 - a_2^2}_{c_3} \quad \left. \vphantom{\frac{r^2 - a_1^2 - a_2^2}{c_3}} \right\} r = \sqrt{c_3 + a_1^2 + a_2^2}$$

↗  
LLS

$$\begin{bmatrix} 2x_1 & 2y_1 & 1 \\ 2x_2 & 2y_2 & 1 \\ \vdots & \vdots & \vdots \\ 2x_m & 2y_m & 1 \end{bmatrix} \begin{bmatrix} c_1 \\ c_2 \\ c_3 \end{bmatrix} = \begin{bmatrix} x_1^2 + y_1^2 \\ x_2^2 + y_2^2 \\ \vdots \\ x_m^2 + y_m^2 \end{bmatrix}$$

Parametric rep'n

$$\begin{cases} x(\theta) = c_1 + r \cos(\theta) \\ y(\theta) = c_2 + r \sin(\theta) \end{cases}$$

for  $\theta \in [0, 2\pi]$