

Recap: ① Optimization → Newton, exact Hessian  
→ Quasi-Newton, approximate Hessian  
↳ BFGS rank-2 updates to approx Hessian

② Num. Lin. Alg.

- Operation counts
- Condition number  
(sens. of output to input)
- Vector / matrix norms.

Today Matrix factorizations, dimension reduction

- $A = LU$  (or  $LL^T$ ) → solving linear systems  
→ generating random variables (later)
- $= QR$  → least squares problems
- $= USV^T$  → matrix approximation, dimension reduction

Generally, when done directly all matrix factorizations of an  $N \times N$  matrix require  $\mathcal{O}(N^3)$  operations.  
↳ "big oh" notation

Heuristically, each element of the matrix is "touched"  $N$  times (e.g. interacts with all columns or rows).  
↑  
 $N^2$  of them

# The LU Factorization

This factorization is equivalent to Gaussian elimination.

Writing  $A = LU$  } Finding  $L, U$  costs  $\mathcal{O}(N^3)$  ops  
 ↑  
 lower triangular,  $(1, \dots, 1)$  on diagonal }  
 upper triangular 

Then to solve  $A\vec{x} = \vec{b}$ , instead first solve  
 $LU\vec{x} = \vec{b}$  } Forward substitution,  $\mathcal{O}(N^2)$   
 then  $U\vec{x} = \vec{y}$  } Backward substitution,  $\mathcal{O}(N^2)$

Failure mode of LU factorization is the same as for Gaussian elimination - a zero in the pivot spot.

To avoid this, permute rows this is called partial pivoting, and is numerically stable.  $\Rightarrow PA = LU$  (mention fill pivoting)  
 ↑ permutation matrix

Ex:  $A = \begin{pmatrix} a & b \\ c & d \end{pmatrix} / \left( \begin{pmatrix} \epsilon & 1 \\ 1 & 1 \end{pmatrix} \vec{x} = \begin{pmatrix} 1 \\ 2 \end{pmatrix} \right)$

$$\begin{pmatrix} 1 & 0 \\ -c/a & 1 \end{pmatrix} A = \begin{pmatrix} a & b \\ 0 & d - bc/a \end{pmatrix}$$

$$\Rightarrow A = \underbrace{\begin{pmatrix} 1 & 0 \\ c/a & 1 \end{pmatrix}}_L \underbrace{\begin{pmatrix} a & b \\ 0 & d - bc/a \end{pmatrix}}_U$$

## Cholesky Factorization

- special case when  $A$  is symmetric positive definite

$$A = LL^T$$

↑  
not 1's on diagonal

- pivoting never necessary

Furthermore,  $\det A = \det LU = \det L \cdot \det U$   
 $= 1 \cdot \det U$   
 $= \prod_{i=1}^N u_{ii}$  ← product of diagonals

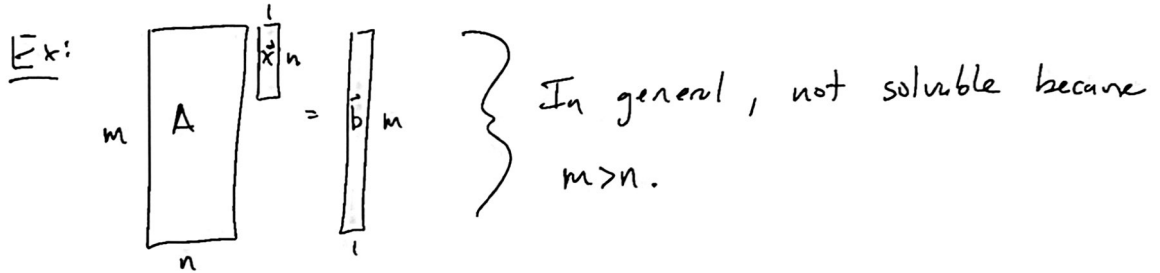
# Least Squares

~~Matrix~~

Two canonical problems in linear algebra:

(1) Solve  $A\vec{x} = \vec{b}$   $\Rightarrow$   $A$  is a square  $n \times n$  matrix.

(2) Find "the best" solution to a system  $A\vec{x} = \vec{b}$  when  $A$  is an  $m \times n$  matrix.



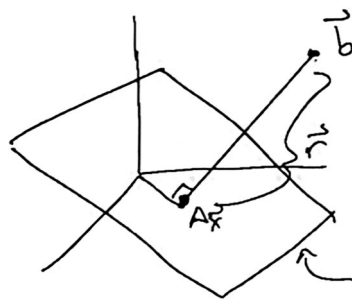
One version of "the best" solution is the least squares solution:

Least squares: For  $A \in \mathbb{R}^{m \times n}$ , with  $m > n$ , find  $\vec{x}$  such that  $\|A\vec{x} - \vec{b}\|_2$  is as small as possible.

$\uparrow$  The fact that this is the 2-norm is important

Ex:  $A \in \mathbb{R}^{3 \times 2}$

Geometrically



$$\min \|A\vec{x} - \vec{b}\|_2$$

minimize the distance between  $\vec{b}$  and  $A\vec{x}$ .

column space of  $A$   
 $\text{col}(A)$

We also have that

$$\|A\vec{x} - \vec{b}\|_2^2 = \sum_{i=1}^m \left( \sum_{j=1}^n a_{ij} x_j - b_i \right)^2 = f(\vec{x})$$

If all you knew was calculus, then you could solve this

by find the solution to  $\left. \begin{array}{l} \frac{\partial f}{\partial x_1} = 0 \\ \vdots \\ \frac{\partial f}{\partial x_n} = 0 \end{array} \right\}$  Finding a zero of  $f'$  [3]  
 $\Rightarrow$  minimizer or maximizer

Let  $\vec{r} = \vec{b} - A\vec{x}$  = residual vector.

$\Rightarrow \|\vec{r}\|$  = distance from  $\text{col}(A)$  to  $\vec{b}$ .

perpendicular.

If  $\vec{x}$  is the least squares solution, then  $\vec{r} \perp \text{col } A$ ,

in particular,  $A^T \vec{r} = \vec{0}$ .

$$\begin{aligned}\Rightarrow A^T \vec{r} = \vec{0} &= A^T (\vec{b} - A\vec{x}) \\ &= A^T \vec{b} - A^T A \vec{x}\end{aligned}$$

$\Rightarrow$  The least squares solution  $\vec{x}$  solves  $\underbrace{A^T A \vec{x} = A^T \vec{b}}$ .

"normal equations"

From a numerical point of view, how best to find  $\vec{x}$  to  $\min \|A\vec{x} - \vec{b}\|_2$ ?

These equations have at least one solution.

Options: (1) Solve the normal equations:

$$\begin{array}{c} A^T A \vec{x} = A^T \vec{b} \\ \uparrow \\ m \times n \end{array}$$

Drawback: Solving  $A^T A \vec{x} = A^T \vec{b}$  has a condition number that is the square of  $A\vec{x} = \vec{b}$ .

Ex: If  $A$  is  $m \times n$  with rank  $n$ , then

$$A = U_{m \times n} S_{n \times n} V_{n \times n}^T \Rightarrow \text{cond}(A) = \frac{\sigma_1}{\sigma_n}$$

Give quick proof

$$A^T A = V S^2 V^T \Rightarrow \text{cond}(A^T A) = \frac{\sigma_1^2}{\sigma_n^2}$$

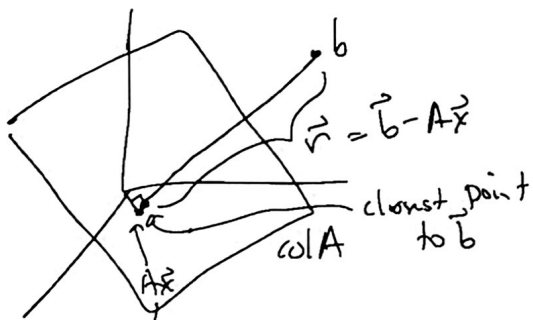
Return to SVD later

So if  $\text{cond}(A) = 10^4$ , then  $\text{cond}(A^T A) = 10^8$ .

Option (2) Use calculus to  $\min \|A\vec{x} - \vec{b}\|$

Option (3): Reformulate the problem to be consistent.

Goal:  $\min \|A\vec{x} - \vec{b}\|$  without solving the normal equations.



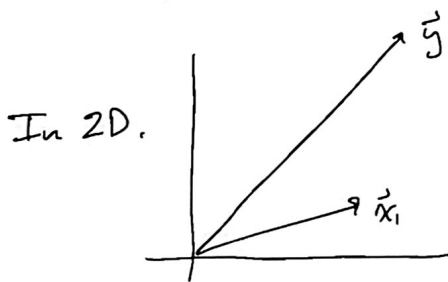
Idea: Construct a linear system that is consistent (i.e. has a solution), without computing  $A^T A$ .

Instead of solving  $A\vec{x} = \vec{b}$ , (which has no solution in general), solve  $A\vec{x} = \vec{b}' \leftarrow \vec{b}' = \text{orthogonal projection of } \vec{b} \text{ onto the column space of } A$ .

$$A\vec{x} = \text{proj}_{\text{col } A} \vec{b} = \vec{b}'$$

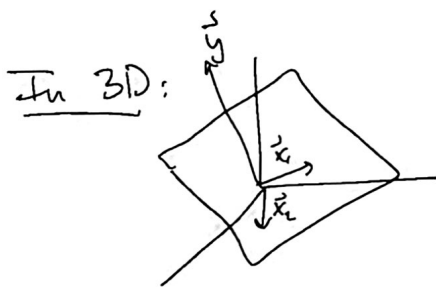
How do we compute  $\text{proj}_{\text{col } A} \vec{b}$ ?

How do we compute a projection of one vector onto another?



The orthogonal projection of  $\vec{y}$  onto  $\vec{x}_1$  is the component of  $\vec{y}$  pointing in the  $\vec{x}_1$  direction.

$$\text{proj}_{\vec{x}_1} \vec{y} = \underbrace{\begin{pmatrix} \vec{y} \cdot \vec{x}_1 \\ \|\vec{x}_1\|^2 \end{pmatrix}}_{\text{inner product}} \underbrace{\begin{pmatrix} \vec{x}_1 \\ \|\vec{x}_1\| \end{pmatrix}}_{\text{unit vector pointing in the direction of } \vec{x}_1}$$



$$\begin{aligned} \text{proj}_{\vec{x}} \vec{y} &= \text{the component of } \vec{y} \text{ in the } \vec{x}_1 - \vec{x}_2 \text{ plane.} \\ &= \underbrace{\left( \vec{y} \cdot \hat{x}_1 \right) \hat{x}_1}_{\hat{x}_1 = \frac{\vec{x}_1}{\|\vec{x}_1\|}} + \underbrace{\left( \vec{y} - (\vec{y} \cdot \hat{x}_1) \hat{x}_1, \hat{x}_2 \right) \hat{x}_2}_{\text{the remaining part of } \vec{y} \text{ after projecting onto } \vec{x}_1} \end{aligned}$$

If  $(\hat{x}_1, \hat{x}_2) = 0$ , (orthogonal), then  $\text{proj}_{\vec{x}} \vec{y} = (\vec{y} \cdot \hat{x}_1) \hat{x}_1 + (\vec{y} \cdot \hat{x}_2) \hat{x}_2$  | Gram Schmidt process.

So this suggests that we want to find an orthonormal basis for  $\text{col}(A)$ ,  $U = \text{span}\{\hat{u}_1, \dots, \hat{u}_n\}$  and then project  $\vec{b}$  onto  $U$ , forming  $\vec{b}'$ .

$$\begin{aligned}\vec{b}' &= (\vec{b}, \hat{u}_1)\hat{u}_1 + (\vec{b}, \hat{u}_2)\hat{u}_2 + \dots + (\vec{b}, \hat{u}_n)\hat{u}_n \\ &= (\hat{u}_1 \ \hat{u}_2 \ \dots \ \hat{u}_n) \begin{pmatrix} (\vec{b}, \hat{u}_1) \\ (\vec{b}, \hat{u}_2) \\ \vdots \\ (\vec{b}, \hat{u}_n) \end{pmatrix}\end{aligned}$$

$$= \underbrace{(\hat{u}_1 \ \hat{u}_2 \ \dots \ \hat{u}_n)}_{U} \begin{pmatrix} \hat{u}_1^T \\ \hat{u}_2^T \\ \vdots \\ \hat{u}_n^T \end{pmatrix} \vec{b}$$

$$= \underbrace{U U^T}_{\text{orthogonal projection of } \vec{b} \text{ onto the columnspace of } A} \vec{b}$$

orthogonal projection of  $\vec{b}$  onto the columnspace of  $A$

$\Rightarrow$  The linear system  $A\vec{x} = U U^T \vec{b}$  is consistent.

The columns of  $U$  are orthonormal, so  $U$  is an orthogonal matrix  $\Rightarrow U^T U = I$ .

Use the Gram-Schmidt process to construct  $U$ .

(Mention numerical stability.)

## G-S Process:

Goal: Given some set of vectors  $\vec{a}_1, \dots, \vec{a}_n$ , form another sequence of vectors  $\hat{q}_1, \hat{q}_2, \dots, \hat{q}_n$  that are orthonormal.

$$\text{Let: } \hat{q}_1 = \vec{a}_1 / \|\vec{a}_1\| \Rightarrow \|\hat{q}_1\| = 1 \quad \text{All norms are } l_2 \text{ norms.}$$

$$\vec{q}'_2 = \vec{a}_2 - \underbrace{(\vec{a}_2, \hat{q}_1)}_{\text{inner product}} \hat{q}_1$$

is not normalized  $(\vec{x}, \vec{y}) = \vec{x}^T \vec{y}$ .

$\wedge$  means unit vector.

$$\hat{q}_2 = \vec{q}'_2 / \|\vec{q}'_2\|$$

Proceed ...

$$\vec{q}'_3 = \vec{a}_3 - (\vec{a}_3, \hat{q}_1) \hat{q}_1 - (\vec{a}_3, \hat{q}_2) \hat{q}_2$$

$$\hat{q}_3 = \vec{q}'_3 / \|\vec{q}'_3\|$$

So when we are done, by construction,

$$Q = (\hat{q}_1 \hat{q}_2 \dots \hat{q}_n) \quad \text{then} \quad Q^T Q = I$$

And the projection of  $\vec{x}$  onto the column space of  $A$ ,

$$\text{i.e., } \text{span}\{\vec{a}_1, \dots, \vec{a}_n\} \text{ is just } \text{proj}_A \vec{x} = Q Q^T \vec{x}$$

Furthermore, this automatically yields a factorization of the matrix  $A$ :

$$\text{col } Q = \text{col } A$$

$$\text{So any } \hat{q}_j = c_1 \vec{a}_1 + \dots + c_n \vec{a}_n$$

$$\text{and likewise } \vec{a}_j = d_1 \hat{q}_1 + \dots + d_n \hat{q}_n$$

This gives us the QR factorization:

$$A = QR$$

$$(\vec{a}_1 \dots \vec{a}_n) = (\hat{q}_1 \hat{q}_2 \dots \hat{q}_n) \begin{pmatrix} r_{11} & r_{12} & r_{13} & \dots & r_{1n} \\ 0 & r_{22} & & & \\ 0 & & \ddots & & \\ \vdots & & & & r_{nn} \end{pmatrix}$$

The question is: what are the  $r_{ij}$ 's?

$r_{ij}$ 's can be obtained from G-S process.

Ex: From G-S:

$$\hat{q}_1 = \frac{\vec{a}_1}{\|\vec{a}_1\|} \Rightarrow \vec{a}_1 = \|\vec{a}_1\| \hat{q}_1 \Rightarrow r_{11} = \|\vec{a}_1\|$$

$$\vec{q}'_2 = \vec{a}_2 - (\vec{a}_2, \hat{q}_1) \hat{q}_1$$

$$\vec{q}'_2 = \|\vec{a}_2 - (\vec{a}_2, \hat{q}_1) \hat{q}_1\| \hat{q}_2$$

$$\Rightarrow \vec{a}_2 = \underbrace{(\vec{a}_2, \hat{q}_1)}_{r_{12}} \hat{q}_1 + \vec{q}'_2$$

$$= r_{12} \hat{q}_1 + r_{22} \hat{q}_2$$

With sufficient bookkeeping, the G-S process yields the factorization  $A = QR$ .

So now: to solve least squares:

$$\textcircled{1} A\vec{x} = Q \underbrace{Q^T \vec{b}}_{\vec{c}}$$

Alternatively:

$$\textcircled{2} \text{ Given that } A = QR$$

$$\text{Solve } \underbrace{R\vec{x}}_{\text{upper triangular system}} = \underbrace{Q^T \vec{b}}_{\vec{c}}$$

} Preferred method  
Show later on  
HW or exam...

What is the operation count?



# Applications: Linear Regression

Given data:  $(x_i^1, x_i^2, y_i)$  for  $i = 1 \dots n$

Generative model:  $y = a + bx_1 + cx_2 + \epsilon$   
 $\epsilon \sim N(0, \sigma^2)$  error or measurement noise.

You observe  $y_i = a + bx_i^1 + cx_i^2 + \epsilon_i$   
 $\uparrow$  dependent variable  
 $\uparrow \quad \uparrow$  independent variables

Assumption is that noise only appears in  $y_i$ .

Given an estimate of  $a, b, c$ :  $\hat{a}, \hat{b}, \hat{c}$ , then the residuals are given as

$$r_i = y_i - \underbrace{(\hat{a} + \hat{b}x_i^1 + \hat{c}x_i^2)}_{\substack{\text{predicted value, given} \\ \text{estimates } \hat{a}, \hat{b}, \hat{c}.}} \\ \uparrow \\ \text{observed value}$$

Question: How do we estimate  $a, b, c$ ?

One option: minimize the squared residuals:

$$\min_{a, b, c} \|\vec{r}\|_2^2 = \min_{a, b, c} \sum_{i=1}^n (y_i - a - bx_i^1 - cx_i^2)^2$$

This is a least squares problem.

Form the least square system:

$$\min_{a, b, c} \left\| \underbrace{\begin{pmatrix} 1 & x_1^1 & x_1^2 \\ 1 & x_2^1 & x_2^2 \\ \vdots & \vdots & \vdots \\ 1 & x_n^1 & x_n^2 \end{pmatrix}}_X \underbrace{\begin{pmatrix} a \\ b \\ c \end{pmatrix}}_{\vec{a}} - \underbrace{\begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{pmatrix}}_{\vec{y}} \right\| = \min_{\vec{a}} \|\underbrace{X}_{\text{design matrix}} \vec{a} - \vec{y}\|$$

$\square$   $\square$