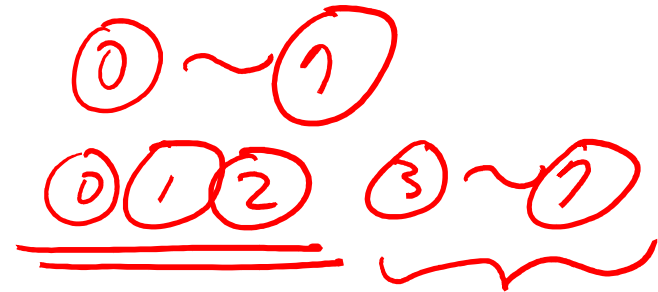


Computations with Large Matrices

- $Ax=b$ in its many variations
 - Ordinary elimination might compute an accurate x , or maybe not
 - Too many equations ($m > n$) and no solution \rightarrow least²
 - Square matrix might be singular A^{-1} ?
 - Solution might be impossible to compute (A is extremely ill-conditioned or simply too large)
 - In deep learning we have too many solutions: we want one that will generalize well to unseen test data
- Separate sources of difficulty
- Identify the problem
- Suggest a course of action



Good Problems


- Every matrix A has a pseudoinverse A^+
 - Inverse for every matrix, but this might not help

$A^+ = A^{-1}$
if invertible
- Elimination will succeed (with row changes, $A \setminus b$)
 - Square and invertible, reasonable size, not large condition number

$A \setminus b$
- $(m > n = r)$ normal equations to find the least squares solution
 - columns of A are independent and not too ill-conditioned
 - b is probably not in the column space of A

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

Difficult Problems

- $(m < n)$ many solutions (underdetermined) ✓
- Columns of A may be in bad condition → 
 - Too large condition number
 - x is not well determined → orthogonalize the columns by a Gram-Schmidt or Householder algorithm
- A may be nearly singular
 - $A^T A$ will have a very larger inverse, Gram-Schmit may fail
 - Different approach to add a penalty term → make $A^T A$ more positive (common in inverse problem)
- A is way too big (no elimination)
 - Random sampling of the columns ✓
 - Results are never certain, but the probability of going wrong is low

Numerical Linear Algebra

By Trefethen and Bau

- Fundamentals
 - Reaching the SVD and Eckart-Young
- QR Factorization and Least Squares
 - All 3 ways: A^+ , $(A^T A)^{-1} A^T$, QR
- Conditioning and Stability
 - Condition numbers, backward stability, perturbations
- Systems of Equations
 - Direct elimination: PA=LU, Cholesky's $S=A^T A$
- Eigenvalues ✓
 - Reduction to tridiagonal-Hessenberg-bidiagonal; QR with shifts
- Iterative Methods
 - Arnoldi, Lanczos, GMRES, conjugate gradients, Krylov

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

$$\mathbf{Ax} = \lambda \mathbf{x}$$

$$\mathbf{Sq} = \lambda \mathbf{q}$$

$$\mathbf{Av} = \sigma \mathbf{u}$$

-
- Krylov Subspaces and Arnoldi Iteration
 - Eigenvalues from Arnoldi
 - Linear Systems by Arnoldi and GMRES ✓
 - Symmetric Matrices: Arnoldi becomes Lanczos
 - Eigenvalues of Tridiagonal T by QR Iteration
 - Computing the SVD
 - Conjugate Gradient for $Sx=b$
 - Preconditioning for $Ax=b$ ✓

Least Squares: Four Ways

- ✓ SVD of A leads to its pseudoinverse $A^+ \rightarrow \hat{\mathbf{x}} = A^+ \mathbf{b}$
- ✓ $A^T A \hat{\mathbf{x}} = A^T \mathbf{b}$ can be solved directly when A has independent columns
- Gram-Schmidt idea produces orthogonal columns in $Q \rightarrow A=QR$
 - $A^T A \hat{\mathbf{x}} = A^T \mathbf{b} \rightarrow (QR)^T (QR) \hat{\mathbf{x}} = (QR)^T \mathbf{b} \rightarrow R^T (Q^T Q) R \hat{\mathbf{x}} = R^T Q^T \mathbf{b}$
 - $\rightarrow R \hat{\mathbf{x}} = Q^T \mathbf{b}$ (safe to solve and fast)
- Minimize $\|\mathbf{b} - A\mathbf{x}\|^2 + \delta^2 \|\mathbf{x}\|^2 \rightarrow$ that penalty changes the normal equations to $(A^T A + \delta^2 I) \mathbf{x}_\delta = A^T \mathbf{b}$
 - Now the matrix is invertible and \mathbf{x}_δ goes to $\hat{\mathbf{x}}$ as $\delta \rightarrow 0$

$A^T A$ and $A^T C A$

weight
covariance



- Samples in applied mathematics
 - Stiffness matrix in mechanical engineering
 - Conductance matrix in circuit theory
 - (weighted) graph Laplacian in graph theory
 - Gram matrix (inner products of columns of A) in mathematics
- Characteristics
 - Symmetry: attractive
 - Size may be a problem
 - Condition number: square of the condition number of A
 - In large problems, expensive and often dangerous to compute
- We try not to compute them \times
 - Orthogonal matrices and triangular matrices are good ones

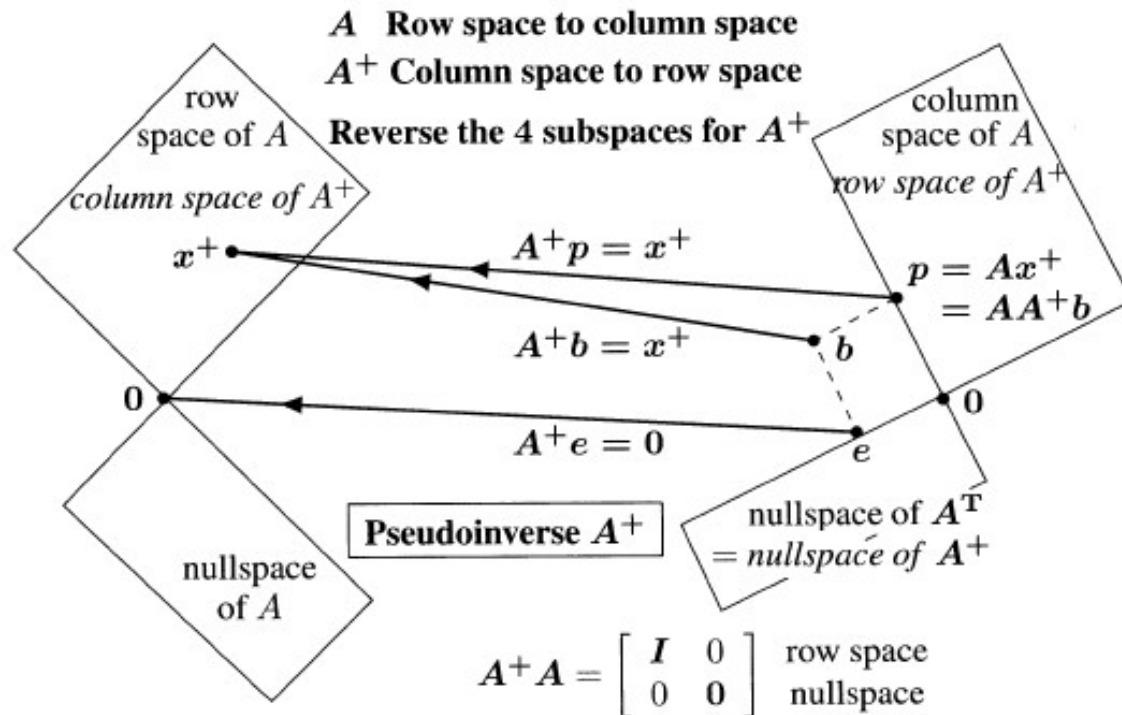
$A^T A$

Q

R

A^+ is the pseudoinverse of A

- Suitable “pseudoinverse” when A has no inverse
 - Rule 1: if A has independent columns, then $A^+ = (A^T A)^{-1} A^T$
 - Rule 2: if A has independent rows, then $A^+ = A^T (A A^T)^{-1}$
 - Rule 3: A diagonal matrix Σ is inverted where possible, otherwise, Σ^+ has zeros



Least Square Solutions to $Ax=b$ is $x^+=A^+b$

- $x^+=A^+b$ is the minimum norm least square solution
 - $x = x^+=A^+b$ makes $\|b-Ax\|^2$ as small as possible
 - If another \underline{x} achieves that minimum then $\|x^+\| < \|\underline{x}\|$

Example: What is the shortest least squares solution to $\begin{bmatrix} 3 & 0 \\ 0 & 0 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} 6 \\ 8 \end{bmatrix}$?

$Ax=b$
 $x=A^{-1}b$
 $x^+=A^+b$

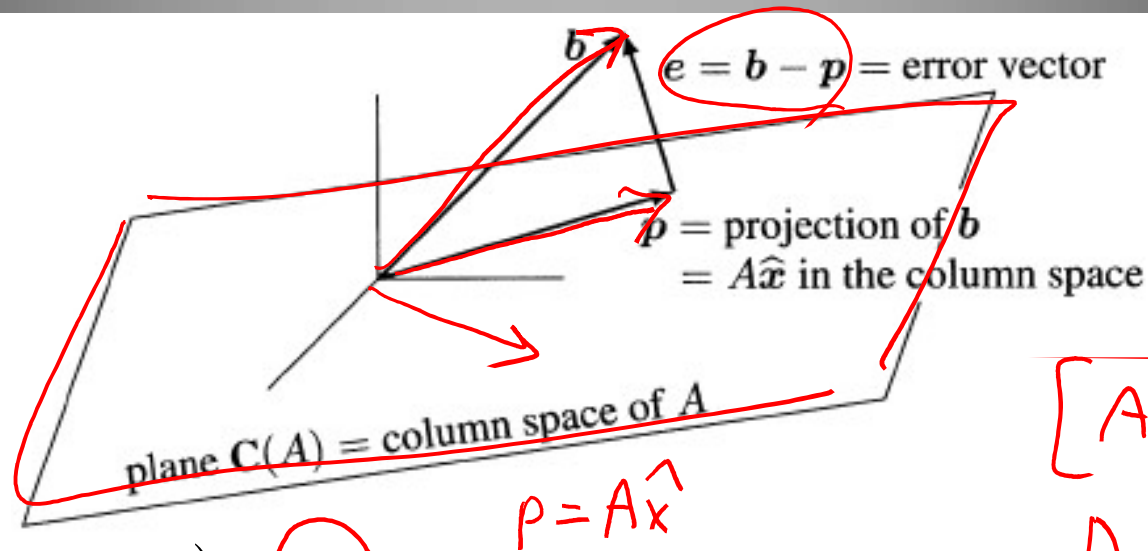
- SVD solve the least squares problem in one step A^+b
 - Computational cost?

squared error $\|b - Ax\|^2 = \|b - U\Sigma\Sigma^T x\|^2 = \|U^T b - \Sigma V^T x\|^2$

$w = \Sigma^+ U^T b$
 $w = V^T x^+ = \Sigma^+ U^T b \rightarrow x^+ = V \Sigma^+ U^T b = A^+ b$

$= \begin{bmatrix} 1/3 & 0 \\ 0 & 0 \end{bmatrix} b$
 $= \begin{bmatrix} 2 \\ 0 \end{bmatrix}$

Normal Equations



$$\mathbf{e} = \mathbf{b} - \mathbf{p} (= \mathbf{b} - \mathbf{A}\hat{\mathbf{x}}) \perp \mathbf{A}\hat{\mathbf{x}}$$

$$\Rightarrow (\mathbf{A}\hat{\mathbf{x}})^T (\mathbf{b} - \mathbf{A}\hat{\mathbf{x}}) = \hat{\mathbf{x}}^T \mathbf{A}^T (\mathbf{b} - \mathbf{A}\hat{\mathbf{x}}) = 0 \rightarrow \mathbf{A}^T (\mathbf{b} - \mathbf{A}\hat{\mathbf{x}}) = 0$$

$$\mathbf{A}^T \mathbf{A}\hat{\mathbf{x}} = \mathbf{A}^T \mathbf{b} \quad (\text{normal equation for } \hat{\mathbf{x}})$$

$$\hat{\mathbf{x}} = (\mathbf{A}^T \mathbf{A})^{-1} \mathbf{A}^T \mathbf{b} \quad (\text{least squares solution to } \mathbf{A}\mathbf{x} = \mathbf{b})$$

$$\mathbf{p} = \mathbf{A}\hat{\mathbf{x}} = \mathbf{A}(\mathbf{A}^T \mathbf{A})^{-1} \mathbf{A}^T \mathbf{b} \quad (\text{projection of } \mathbf{b} \text{ onto the column space of } \mathbf{A})$$

$$\mathbf{P} = \mathbf{A}(\mathbf{A}^T \mathbf{A})^{-1} \mathbf{A}^T \quad (\text{projection matrix that multiplies } \mathbf{b} \text{ to give } \mathbf{p})$$

$$\begin{bmatrix} \mathbf{A} \end{bmatrix} \begin{bmatrix} \mathbf{x} \end{bmatrix} = 0$$

$$x_1 a_1 + x_2 a_2 + \dots$$

Least Squares with a Penalty Term

- To prove that the limit is A^+ for every matrix A

Minimize $\|\mathbf{b} - \mathbf{A}\mathbf{x}\|^2 + \delta^2 \|\mathbf{x}\|^2 \rightarrow$ Solve $(\mathbf{A}^T \mathbf{A} + \delta^2 \mathbf{I})\hat{\mathbf{x}} = \mathbf{A}^T \mathbf{b}$ $\delta \rightarrow 0$

$$\mathbf{A} = \mathbf{U}\mathbf{\Sigma}\mathbf{\Sigma}^T$$

$$\hat{\mathbf{x}} = (\mathbf{A}^T \mathbf{A} + \delta^2 \mathbf{I})^{-1} \mathbf{A}^T \mathbf{b} \rightarrow \hat{\mathbf{x}} = \mathbf{A}^+ \mathbf{b}$$

$$\mathbf{A}^T \mathbf{A} + \delta^2 \mathbf{I} = \mathbf{V}\mathbf{\Sigma}^T (\mathbf{U}^T \mathbf{U}) \mathbf{\Sigma} \mathbf{V}^T + \delta^2 \mathbf{I} = \mathbf{V}(\mathbf{\Sigma}^T \mathbf{\Sigma} + \delta^2 \mathbf{I}) \mathbf{V}^T$$

$$(\mathbf{A}^T \mathbf{A} + \delta^2 \mathbf{I})^{-1} \mathbf{A}^T = \mathbf{V}(\mathbf{\Sigma}^T \mathbf{\Sigma} + \delta^2 \mathbf{I})^{-1} (\mathbf{V}^T \mathbf{V}) \mathbf{\Sigma}^T \mathbf{U}^T = \mathbf{V} \left[(\mathbf{\Sigma}^T \mathbf{\Sigma} + \delta^2 \mathbf{I})^{-1} \mathbf{\Sigma}^T \right] \mathbf{U}^T$$

$$\lim_{\delta \rightarrow 0} (\mathbf{A}^T \mathbf{A} + \delta^2 \mathbf{I})^{-1} \mathbf{A}^T = \lim_{\delta \rightarrow 0} \mathbf{V} \left[(\mathbf{\Sigma}^T \mathbf{\Sigma} + \delta^2 \mathbf{I})^{-1} \mathbf{\Sigma}^T \right] \mathbf{U}^T = \mathbf{V} \mathbf{\Sigma}^+ \mathbf{U}^T = \mathbf{A}^+$$

Randomized Linear Algebra

- Matrix multiplication

- Sampling matrix S

- $C=AS$ and $R=S^T B \rightarrow CR=AS S^T B \sim AB$

- It will not be true that SS^T is close to I

- It will be true that the expected value of SS^T is I ↴

- Random matrix multiplication with the correct mean AB

- ✓ • Norm-squared sampling minimizes the variance

- Applications of random matrix multiplication

- Interpolative approximation $A \sim \underline{CMR}$

- Application of A by a low rank matrix

- Approximation of the SVD of A

$$AB \rightarrow CR$$

↙ ↘
 AS $S^T B$

$$\frac{\|a_j\| \|b_j^T\|^2}{c}$$

P_j