# Computational data mining

Fatemeh Mansoori

# Ax = b

- 1) *A* is square and invertible and its condition number $\sigma_1/\sigma_n$ is not large
  - The elimination will succeed.
  - We have $PA = LU$ or $A = LU$
- 2) m > n = r : There are too many equations Ax = b to expect a solution
  - If the columns of *A* are independent (invertible $A^TA$) and not too ill-conditioned then we solved the normal equation $A^T A \hat{x} = A^T b$
  - Vector b is probably not in column space of A, and Ax = b is impossible.
  - $A\hat{x}$ is the projection of b onto column space of A.
- 3) m < n. Ax = b have many solutions.
  - A has non zero null space.
  - Want to choose best x
  - Possible choices are x+ and x1
  - x+=A$^+$b. A$^+$ gives the minimum l2 norm solution with nullspace component = zero
  - x1 = minimum l1 norm solution. (this solution is often sparse: many zero components)

# $Ax = b$

- 4) column of A may be in bad condition or near singular
  - $A^TA$ will have very large inverse
  - Orthogonize the columns by Gram-Schmidt algorithm
- 5) column of A may be in bad condition or near singular
  - Gram-Schmidt may fail.
  - A different approach is to add a penalty term
- 6) A is too big
  - Best solution is random sampling of column

# Least Squares problem

- Ax = b
  - Have distinct solution : If A is square and rank r = n and b is on the column space of A

  - unsolvable linear equations :
    - A is singular
    - b is not on the column space of A
    - Need to produce the best solution $\widehat{x}$
    - Least square method choose $\widehat{x}$ to make $\|b - A\widehat{x}\|^2$ as small as possible
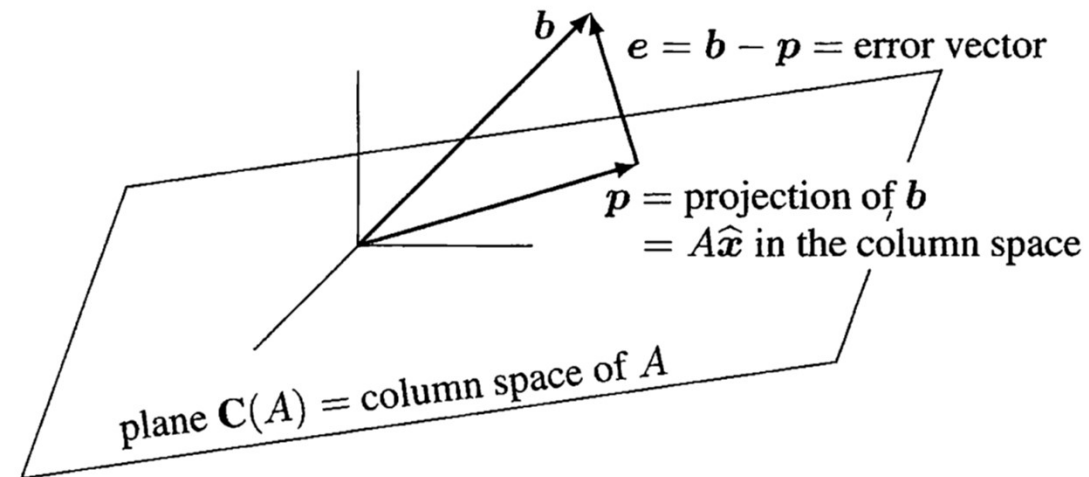    - Minimizing means that its derivation are zero

# Least square problem

- Suppose you have m measurement mi
- The measurements are noisy.
- You want to find c + dx = b where the best line you can fit to the data
- The solution is to try to minimize the ||b-Ax||. Where :

  - $m_i$s are in the matrix A . A=$\begin{bmatrix} 1 & m_1 \\ 1 & m_2 \\ \vdots & \vdots \\ 1 & m_n \end{bmatrix}$

  - x=$\begin{bmatrix} c \\ d \end{bmatrix}$

# Least Squares problem

- b is not in the column space of *A*
  - *Ax = b* has no solution
  - project *b* onto the column space of *A*
  - $A^{\mathrm{T}} A \widehat{x} = A^{\mathrm{T}} b.$
  - To invert $A^T A$ we need to know that *A* has independent columns
  - Show that e is peroendecular
  to all vector Ax in the column
  space



$b$

$e = b - p = $ error vector

$p = $ projection of $b$
$= A\widehat{x}$ in the column space

plane $\mathbf{C}(A) = $ column space of $A$

# Least Squares problem

- *A* now has independent columns : *r* = *n.* That makes $A^T A$ positive definite and invertible

- We could check that our $\widehat{x}$ is the same vector $x^+ = A^+ b$ that came from the pseudoinverse

- There are no other $\widehat{x}$ because the rank is assumed to be *r* = *n*
  - nullspace of *A* only contains the zero vector

# When is $A^TA$ Invertible?

- $A^TA$ is invertible exactly when $A$ has independent columns
- Always $A$ and $A^TA$ have the same nullspace
  - Why?

# Pseudoinverse of *A*

- If *A* is invertible then $A^+$ is $A^{-1}$

- If *A* is *m* by *n* then $A^+$ is *n* by *m*

- When *A* multiplies a vector *x* in its row space, this produces *Ax* in the column space
  - Those two spaces have equal dimension *r* (the rank).
  - Restricted to these spaces *A* is always invertible and $A^+$ inverts *A*.
  - $A^+Ax = x$ exactly when *x* is in the row space
  - $AA^+b = b$ when *b* is in the column space

- The nullspace of $A^+$ is the nullspace of $A^T$

  - It contains the vectors y in $R^m$ with $A^Ty = 0$

  - vectors $y$ are perpendicular to every $Ax$ in the column space

  - For these $y$ we accept $x^+ = A^+y = 0$ as the best solution to the unsolvable equation $Ax = y$

# "pseudoinverse" when *A* has no inverse

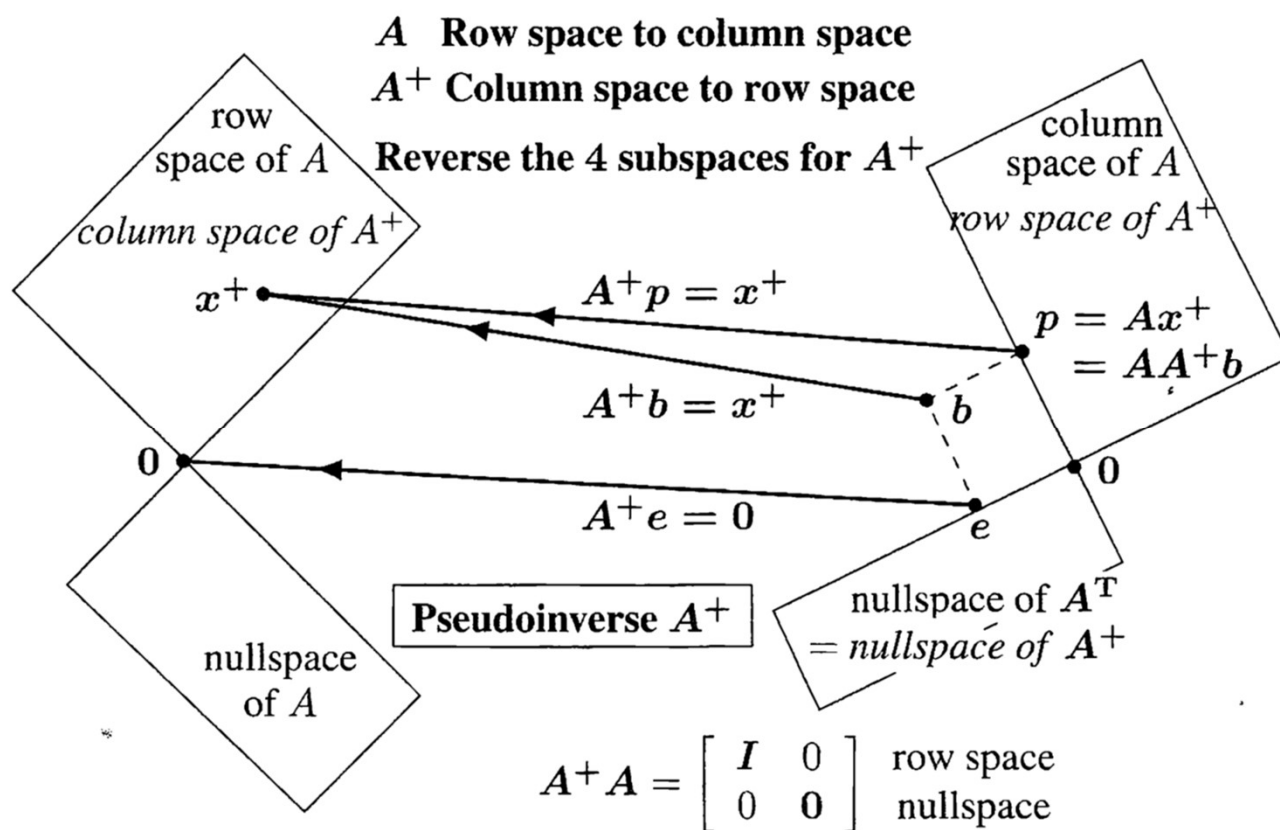**Rule 1**  If $A$ has independent columns, then $A^+ = (A^T A)^{-1} A^T$ and so $A^+ A = I$.

**Rule 2**  If $A$ has independent rows, then $A^+ = A^T (A A^T)^{-1}$ and so $A A^+ = I$.

**Rule 3**  A diagonal matrix $\Sigma$ is inverted where possible—otherwise $\Sigma^+$ has zeros

$$\Sigma = \begin{bmatrix} \sigma_1 & 0 & 0 & 0 \\ 0 & \sigma_2 & 0 & 0 \\ 0 & 0 & 0 & 0 \end{bmatrix} \qquad \Sigma^+ = \begin{bmatrix} 1/\sigma_1 & 0 & 0 \\ 0 & 1/\sigma_2 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \end{bmatrix}$$

On the four subspaces

$$\Sigma^+ \Sigma = I \quad \Sigma \Sigma^+ = I$$

$$\Sigma^+ \Sigma = 0 \quad \Sigma \Sigma^+ = 0$$

**All matrices**

$$\text{The pseudoinverse of } A = U\Sigma V^{\mathrm{T}} \text{ is } A^+ = V\Sigma^+ U^{\mathrm{T}}.$$



$A$  Row space to column space
$A^+$ Column space to row space

Reverse the 4 subspaces for $A^+$

row space of $A$
column space of $A^+$

column space of $A$
row space of $A^+$

$x^+$

$A^+ p = x^+$

$p = Ax^+$
$= AA^+ b$

$A^+ b = x^+$

$b$

$0$

$A^+ e = 0$

$e$

$0$

Pseudoinverse $A^+$

nullspace of $A^{\mathrm{T}}$
$= $ nullspace of $A^+$

nullspace of $A$

$$A^+ A = \begin{bmatrix} I & 0 \\ 0 & 0 \end{bmatrix} \quad \begin{matrix} \text{row space} \\ \text{nullspace} \end{matrix}$$

- pseudoinverse *A+* solves the least squares equation $A^{T}Ax = A^{T}b$ in one step

$$A^{+}b = V\Sigma^{+}U^{\mathrm{T}}b \text{ is best possible}$$

# When r < n

- $x^+ = A^+b$ is the minimum norm least squares solution
- When $A$ has independent columns and rank $r = n$
  - $X^+$ is the only least square solution
- if there are nonzero vectors $x$ in the nullspace of $A$ ($r < n)$ they can be added to $x+$
  - The error $||b - A(x^+ + x)||$ is not affected when $Ax = 0$ but the length $|| x^+ + x||$ *will grow*

- Example: find the shortest square solution to $\begin{bmatrix} 3 & 0 \\ 0 & 0 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} 6 \\ 8 \end{bmatrix}$

# Gram schmidt

- The column of A are still to be independence
- Third approach to find the x is orthogonalzing the columns of A
  - Why : ? The operation count is doubled compared to the $A^{\mathrm{T}}A\widehat{x} = A^{\mathrm{T}}b,$
  - Because : orthogonal vectors provide numerical stability
  - Stability is become import when $A^{\mathrm{T}}A$ is nearly singular.
  - The condition number of $A^{\mathrm{T}}A$ is its norm $|| A^{\mathrm{T}}A||$ times $|| (A^{\mathrm{T}}A)^{-1}||$
  - When $\sigma_1^2/\sigma_n^2$ is large, it is wise to orthogonalize the columns of A in advance
  -

# singular matrix or ill conditioned

- Near singular matrix or ill conditioned
  - Matrix which has its determinate close to zero and whose inverse is unreliable
- The extend of ill-conditioned matrix id defined by its condition number
- Matrix A is ill-conditioned if it is invertible but can become non-invertible(singular) if some of its entries are changed ever so slightly
  - Condition number of A is a measure of how ill-conditioned A is
  - The bigger the condition number is the more ill-conditioned A is
  - Well-conditioned matrices have condition number close to 1

- Solving linear systems whose coefficient matrices are ill-conditioned is tricky
  - A small change in the data (right hand side vector), can lead to radically different answer.

- Example

$$A = \begin{bmatrix} 4.5 & 3.1 \\ 1.6 & 1.1 \end{bmatrix}, \quad b = \begin{bmatrix} 19.249 \\ 6.843 \end{bmatrix}, \quad b_1 = \begin{bmatrix} 19.25 \\ 6.84 \end{bmatrix}.$$

```python
A = np.array([[4.5,3.1],
              [1.6,1.1]])

b1 = np.array([19.249, 6.843]).reshape(2,1)
b2 = np.array([19.25, 6.84]).reshape(2,1)

x1 = np.matmul(np.linalg.inv(A), b1)
print(x1)

x2 = np.matmul(np.linalg.inv(A), b2)
print(x2)


u,s,v = np.linalg.svd(A)
condition = s[0]/s[1]
print(condition)
```

```
[[3.94]
 [0.49]]
[[2.9]
 [2. ]]
3362.999702646099
```

# Gram-Schmidt

Independent columns $a_1, \ldots, a_n$ lead to orthonormal $q_1, \ldots, q_n$

$q_1 = a_1/\|a_1\|$.

**Gram-Schmidt step**  **Orthogonalize**  $A_2 = a_2 - (a_2^\mathsf{T} q_1)\, q_1$
**Normalize**  $q_2 = A_2/\|A_2\|$

$A_3 = a_3 - (a_3^\mathsf{T} q_1)\, q_1 - (a_3^\mathsf{T} q_2)\, q_2$  **Normalize**  $q_3 = \dfrac{A_3}{\|A_3\|}$

$a_1 = \|a_1\|\, q_1$

$a_2 = (a_2^\mathsf{T} q_1)\, q_1 + \|A_2\|\, q_2$

$a_3 = (a_3^\mathsf{T} q_1)\, q_1 + (a_3^\mathsf{T} q_2)\, q_2 + \|A_3\|\, q_3$

$$\begin{bmatrix} a_1 & a_2 & a_3 \end{bmatrix} = \begin{bmatrix} q_1 & q_2 & q_3 \end{bmatrix} \begin{bmatrix} r_{11} & r_{12} & r_{13} \\ 0 & r_{22} & r_{23} \\ 0 & 0 & r_{33} \end{bmatrix}$$

**Gram-Schmidt produces orthonormal $q$'s from independent $a$'s. Then $A = QR$.**

$\widehat{x} = (A^{\mathrm{T}} A)^{-1} A^{\mathrm{T}} b$ is $(R^{\mathrm{T}} R)^{-1} R^{\mathrm{T}} Q^{\mathrm{T}} b$. This is exactly $\widehat{x} = R^{-1} Q^{\mathrm{T}} b$.

```python
Q,R = np.linalg.qr(A)
print(Q)
print(R)
qb1 = np.matmul(Q,b1)
qb2 = np.matmul(Q,b2)
print("---------")
print(qb1)
print(qb2)
print("---------")
print(np.matmul(np.linalg.inv(R),np.array([[-20.42],[-0.001]])))
print(np.matmul(np.linalg.inv(R),np.array([[-20.42],[0]])))
```

```
[[3.94]
 [0.49]]
[[2.9]
 [2. ]]
3362.999702646099
-0.00999999999999957
[[-0.94221469 -0.33500967]
 [-0.33500967  0.94221469]]
[[-4.77598157e+00 -3.28937617e+00]
 [ 0.00000000e+00 -2.09381042e-03]]
---------
[[-2.04291617e+01]
 [-1.02596711e-03]]
[[-2.04290989e+01]
 [-4.18762085e-03]]
---------
[[3.94662327]
 [0.47759816]]
[[4.27556088]
 [0.        ]]
```

What is the advantage of Gram Schmidt in solving Least square problem when A is ill conditioned ?

What is the condition number of Q

# Gram-Schmidt with Column Pivoting

- straightforward description of Gram-Schmidt worked with the columns of A in their original order a1, a2, a3, …

- This could be dangerous!
  - Then roundoff error could wipe us out

- We need column exchanges to pick the largest remaining column. Change the order of columns as we go.

**Old**    Accept column $a_j$ as next. Subtract its components in the directions $q_1$ to $q_{j-1}$

**New**    When $q_{j-1}$ is found, subtract the $q_{j-1}$ component from **all remaining columns**

# Algorithm for Gram Schmidt with Column Pivoting

$$i \quad = \text{argmax} \, \|A_{j-1}(:, \ell)\| \text{ finds the largest column not yet chosen for the basis}$$

$$\boldsymbol{q}_j = A_{j-1}(:, i)/\|A_{j-1}(:, i)\| \text{ normalizes that column to give the new unit vector } \boldsymbol{q}_j$$

$$Q_j = \begin{bmatrix} Q_{j-1} & \boldsymbol{q}_j \end{bmatrix} \text{ updates } Q_{j-1} \text{ with the new orthogonal unit vector } \boldsymbol{q}_j$$

$$\boldsymbol{r}_j = \boldsymbol{q}_j^{\mathrm{T}} A_{j-1} \text{ finds the row of inner products of } \boldsymbol{q}_j \text{ with remaining columns of } \boldsymbol{A}$$

$$R_j = \begin{bmatrix} R_{j-1} \\ \boldsymbol{r}_j \end{bmatrix} \text{ updates } R_{j-1} \text{ with the new row of inner products}$$

$$A_j = A_{j-1} - \boldsymbol{q}_j \boldsymbol{r}_j \text{ subtracts the new rank-one piece from each column to give } A_j$$

# Least Squares with a Penalty Term

- If A has dependent columns and Ax = 0 has nonzero solutions, then $A^TA$ cannot be invertible
  - This is where we need $A^+$
- A gentle approach will "regularize" least squares

**Penalty term**   Minimize $||Ax-b||^2+\delta^2\,||x||^2$   Solve $(A^TA + \delta^2 I)\,\widehat{x} = A^Tb$

- This fourth approach to least squares is called ridge regression
- The x in the above equation approaches shortest solution   $x^+ = A^+b$

- The difficulty of computing A+ is to know if a singular value is zero or very small

The diagonal entry in $\Sigma^+$ is zero or extremely large

$$\text{From 0 to } 2^{10} \quad \begin{bmatrix} 2 & 0 \\ 0 & 0 \end{bmatrix}^+ = \begin{bmatrix} 1/2 & 0 \\ 0 & 0 \end{bmatrix} \quad \text{but} \quad \begin{bmatrix} 2 & 0 \\ 0 & 2^{-10} \end{bmatrix}^+ = \begin{bmatrix} 1/2 & 0 \\ 0 & 2^{10} \end{bmatrix}$$

# Pseudoinverse $A^+$ is the Limit of $(A^{\mathrm{T}}A + \delta^2 I)^{-1}A^{\mathrm{T}}$

• Suppose A is 1 by 1 matrix *(just a single number $\sigma$)*

For $\delta > 0$ $\quad (A^{\mathrm{T}}A + \delta^2 I)^{-1}A^{\mathrm{T}} = \left[ \dfrac{\sigma}{\sigma^2 + \delta^2} \right]$ is 1 by 1

limit is zero if $\sigma = 0$. The limit is $\dfrac{1}{\sigma}$ if $\sigma \neq 0$.

# Pseudoinverse $A^+$ is the Limit of $(A^{\mathrm{T}}A + \delta^2 I)^{-1}A^{\mathrm{T}}$

- Now suppose a diagonal matrix $\Sigma$.
- seeing the 1 by 1 case at every position along the main diagonal

$\Sigma$ has positive entries $\sigma_1$ to $\sigma_r$ and otherwise all zeros

$(\Sigma^{\mathrm{T}}\Sigma + \delta^2 I)^{-1}\Sigma^{\mathrm{T}}$ has positive diagonal entries $\dfrac{\sigma_i}{\sigma_i^2 + \delta^2}$ and otherwise all zeros.

Positive numbers approach $\frac{1}{\sigma_i}$. Zeros stay zero. When $\delta \to 0$ the limit is again $\Sigma^+$.

prove that the limit is $A^+$ for *every matrix* $A$

$$\underset{\delta \to 0}{\text{limit}} \quad V\left[(\Sigma^{\text{T}}\Sigma + \delta^2 I)^{-1}\Sigma^{\text{T}}\right] U^{\text{T}} = V\Sigma^+ U^{\text{T}} = A^+.$$

# Question

- Is it true that $(AB)^+ = B^+ A^+$.

- $A = \begin{bmatrix} 1 & 0 \end{bmatrix}$ and $B = \begin{bmatrix} 1 \\ 1 \end{bmatrix}$

$$(A^{\mathbf{T}})^+ = (A^+)^{\mathbf{T}}$$

$$(A^{\mathbf{T}} A)^+ = A^+ (A^{\mathbf{T}})^+.$$

If $C$ has full column rank and $R$ has full row rank then $(CR)^+ = R^+C^+$ is true.

- By using above equation the pseudoinverse of A could be computed with out using SVD

- Which matrices have A+ = A ?