

Sep 30, 2020

New topic: least squares, QR factorization

Model:  $b_i \approx a_i^T x + z \quad i = 1, \dots, m$

$\uparrow$  outcome       $\uparrow$  features       $\uparrow$  intercept

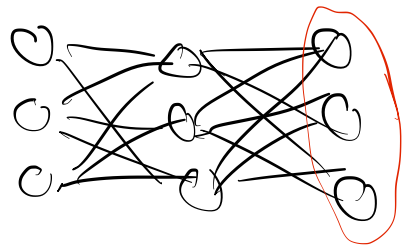
Error:  $\sum_{i=1}^m (a_i^T x + z - b_i)^2 \Rightarrow \left\| \begin{pmatrix} a_1^T \\ \vdots \\ a_m^T \\ 1 \end{pmatrix} \begin{pmatrix} x \\ z \end{pmatrix} - b \right\|_2^2$

Linear least squares:  $\min_x \|Ax - b\|_2^2$

$\begin{matrix} m & n & n & m \\ \boxed{A} & \boxed{x} & - & \boxed{b} \end{matrix}$

Examples:

(1) stats/ML  $a_i =$  data about page  $i$ ,  $b_i =$  # clicks



Final layer  
 $= a_i$

# IJALM

(2) polynomial fitting  $a_i^T = (1 \ z_i \ z_i^2 \ z_i^3)$   $b_i = f(z_i)$

$(z_i, f(z_i))$

$$\min_x \|Ax - b\|_2^2$$

Why squared errors?

$$\|\cdot\|_2 \quad \|\cdot\|_\infty \quad \dots$$

- ① not unreasonable
- ② computationally nice
- ③ statistical interpretations
- ④ building block

Gauss-Markov theorem

$$b_i = a_i^T x^* + \varepsilon_i \quad \left\{ \begin{array}{l} E(\varepsilon_i) = 0 \\ \text{Var}(\varepsilon_i) < \infty \end{array} \right.$$

$$\hat{x} = \arg \min_x \|Ax - b\|_2^2 \quad \left\{ \begin{array}{l} \text{Cov}(\varepsilon_i, \varepsilon_j) = 0 \end{array} \right.$$

- (1)  $E(\hat{x}) = x^*$  (unbiased)
- (2) least variance amongst all linear estimators

(BLUE)

$$b_i = a_i^T x + \varepsilon_i, \quad \varepsilon_i \sim N(0, \sigma^2) \text{ i.i.d.}$$

$$b_i - a_i^T x = \varepsilon_i, \quad f(z) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{1}{2}z^2/\sigma^2\right)$$

Likelihood  $(b_i, A, x)$

$$= \prod_{i=1}^m f(b_i - a_i^T x) = \left( \frac{1}{\sqrt{2\pi}\sigma} \right)^m \exp\left(-\frac{1}{2\sigma^2} \left( \sum_{i=1}^m (b_i - a_i^T x)^2 \right)\right)$$

$$\text{max likelihood} \Rightarrow \min \|Ax - b\|_2^2$$

What is the answer?

$$\hat{x} = \arg \min_x \|Ax - b\|_2$$

$$(Ax - b)^T (Ax - b)$$

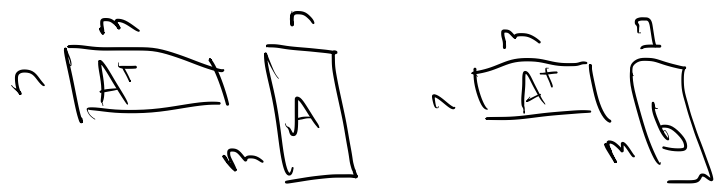
$$= x^T A^T A x - 2x^T A^T b + b^T b$$

Take gradient w.r.t x ...

$$2A^T A \hat{x} - 2A^T b = 0$$

$$A^T A \hat{x} = A^T b$$

Normal equations



Non singular?

$$\hat{x} = (A^T A)^{-1} A^T b$$

Moore-Penrose pseudoinverse,  $A^+$

Numbers  $b_1, \dots, b_m$ . Summary stat  $\bar{x}$

$$\frac{1}{n} \sum_{i=1}^n (b_i - \bar{x})^2 = \frac{1}{n} \left\| \begin{pmatrix} 1 \\ \vdots \\ 1 \end{pmatrix} x - b \right\|_2^2$$

$$A = e$$

$$\hat{x} = (e^T e)^{-1} (e^T b) = \frac{1}{n} \sum b_i = \text{mean}$$

residual  $r = A\hat{x} - b$   $\|r\|_2^2$

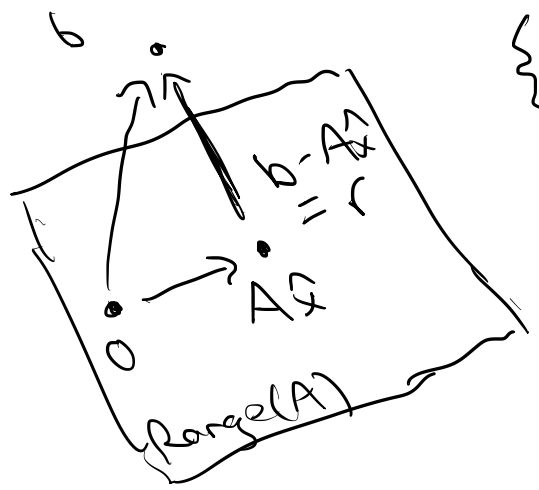
$$r = \frac{1}{n} \sum_{i=1}^n (\hat{x} - b_i)^2 = \text{variance}$$

$$A^T A \hat{x} = A^T b \Leftrightarrow A^T (A \hat{x} - b) = 0$$

$$\min r \Leftrightarrow b - A \hat{x} \in \text{Null}(A^T)$$

$$\Leftrightarrow r \perp \text{Range}(A)$$

$$\{Ax\}$$



How do we compute  $\hat{x}$ ?

Idea 1: just a linear system!

$$(A^T A) \hat{x} = A^T b \quad A^T A \text{ SPD} \Rightarrow A^T A = L L^T \text{ (cholesky)}$$

$$\text{Flops: } O(mn^2) + O(mn) + O(n^3) \Rightarrow O(mn^2)$$

$$\min_x \left\| \begin{matrix} m \\ A \end{matrix} x - \begin{matrix} 1 \\ b \end{matrix} \right\|_2^2 \quad r=0 \quad \|r\|_2^2 \geq 0 \quad \hat{x} = A^{-1} b$$

$O(m^3)$

Conditioning?

$$A = U \Sigma V^T \quad A^T A = V \Sigma^2 V^T \quad \kappa_2(A^T A) = \frac{\sigma_1^2}{\sigma_n^2} = \kappa_2(A)^2$$

Idea 2: SVD

$$A = U \Sigma V^T \quad \hat{x} = A^+ b = (A^T A)^{-1} A^T b$$

$$A^T A = V \Sigma^2 V^T$$

$$\begin{aligned} (A^T A)^{-1} A^T b &= V \Sigma^{-2} \cancel{V^T} V \Sigma U^T b \\ &= V \Sigma^{-1} U^T b \end{aligned}$$

Can be used but there are faster algorithms  
(very useful for analysis!)

Idea 3: QR factorization

"full QR"

$${}^m \boxed{A} = {}^m \boxed{Q} \begin{matrix} {}^n \\ \boxed{R} \\ \hline \boxed{0} \end{matrix}$$

$$Q^T Q = Q Q^T = I$$

R upper tri;

$$\|Ax - b\|_2^2 = \|Q^T (Ax - b)\|_2^2$$

$$\cancel{Q^T Q} \boxed{\begin{matrix} R \\ 0 \end{matrix}}$$

$$= \left\| \begin{pmatrix} R \\ 0 \end{pmatrix} x - Q^T b \right\|_2^2$$

$$Rx = Q^T b \quad (1:n, 1:n)$$

"thin economy / reduced QR"

$$A = \begin{pmatrix} Q_1 & Q_2 \end{pmatrix} \begin{pmatrix} R \\ 0 \end{pmatrix} = Q_1 R$$

$${}^m \boxed{A} = {}^m \boxed{Q} \boxed{R}$$

$$Q^T Q = I \quad (Q Q^T \neq I)$$

$$R \hat{x} = Q^T b$$

① Thin QR  $O(mn^2)$

$$\hat{x} = (A^T A)^{-1} A^T b = (R^T \cancel{Q^T Q} R)^{-1} R^T Q^T b$$

②  $c = Q^T b$

$$= R^{-1} \cancel{R^{-T} R^T} Q^T b$$

③  $R \hat{x} = c$

How to compute QR?

From what we know...

$$A^T A = R^T Q^T Q R = R^T R = L L^T \quad L = R^T$$

① Cholesky  $\Rightarrow L^T = R$

②  $A = QR \Rightarrow Q = AR^{-1}$

"Cholesky QR"

Graph  $\Leftrightarrow A$

$$d = Ae \quad d_i = \sum_j A_{ij}$$

$$\rightarrow v = \operatorname{arg\,min}_i d_i$$

$A(:, v)$  nonzeros  $\Rightarrow$  neighbors  $N$

$$d[N] = 1$$

$$\operatorname{arg\,min}_{i \neq v} d_i$$