

Linear Algebra for AI & ML

(September. 17)



LS problem:

$$\min_{x \in \mathbb{R}^n} \|Ax - b\|^2 \quad (2)$$

Assumption: Columns of A are linearly independent.

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

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

\hat{x} which minimizer should satisfy

$$A^T A \hat{x} = A^T b \quad : \text{normal eqns}$$

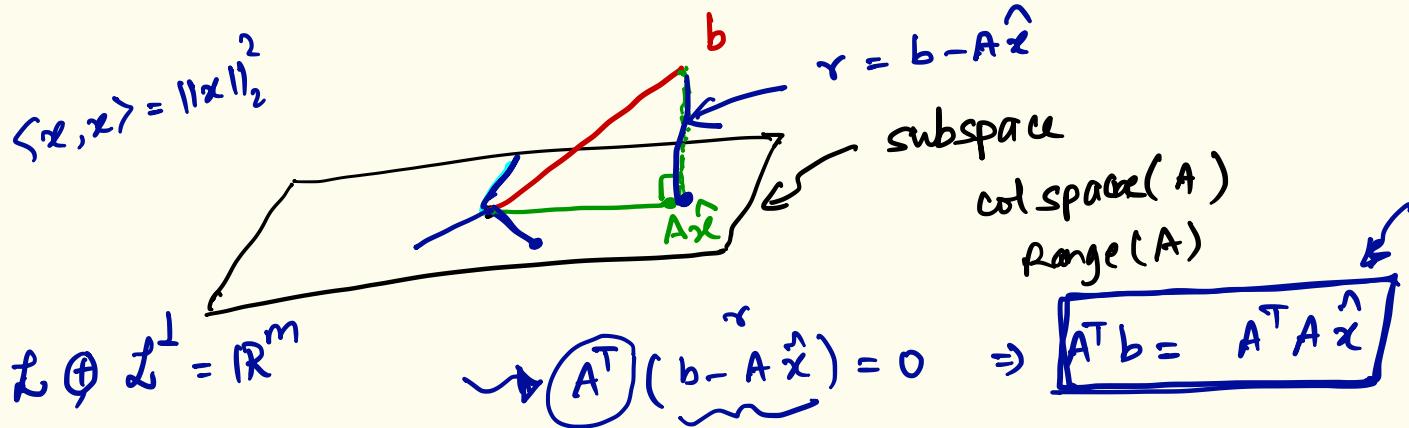
Column interpretation of the LS problem:

$$Ax = x_1[A_1] + x_2[A_2] + \dots + x_n[A_n]$$

Orthogonality principle in the least squares

$$\min_{x \in \mathbb{R}^n} \|Ax - b\|_2^2 \quad (2) \quad = \min_{x \in \mathbb{R}^n} \|x_1 A_1 + \dots + x_n A_n - b\|_2^2$$

\$A_i\$: \$i^{th}\$ column of \$A\$.



$b \notin \text{colspan}(A)$

$\text{Range}(A)$

Find \hat{x} s.t. $A\hat{x}$ is closest to b

$$\|A\hat{x} - b\|_2$$

$$\min_{x \in \mathbb{R}^n} \|Ax - b\|_\infty = \min_{x \in \mathbb{R}^n} \left\{ \max_i |r_i| \right\} \quad \text{LASSO}$$

$r_i = (Ax - b)_i$

$$\min_{x \in \mathbb{R}^n} \|Ax - b\|_1 = \min_{x \in \mathbb{R}^n} \left\{ |r_1| + \dots + |r_m| \right\} \quad \|Ax - b\|_\infty = \max_i |(Ax - b)_i|$$

$$\min_{x \in \mathbb{R}^n} \|Ax - b\|_2^2 = \min_{x \in \mathbb{R}^n} \left\{ (r_1)^2 + \dots + (r_m)^2 \right\} \quad \text{LS}$$

How to compute \hat{x} ??

$$\underbrace{(A^T A)}_{\text{normal eqn}} \hat{x} = A^T b$$

$$\hat{x} = \underbrace{(A^T A)^{-1}}_{\text{be computed like this!}} A^T b$$

\hat{x} will NOT be computed like this!!

$$[A]_{m \times n} = QR$$

$$[]_{m \times n} []_{n \times n}$$

$$(A^T A)^{-1} A^T = R^{-1} Q^T$$

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

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

back substitution.

Given a matrix A , we should avoid
computation of $A^T A$ as well.

(problem case: matrix A has almost
linearly dependent columns)

Least squares date fitting / classification

binary multiclass

$$f: \mathbb{R}^n \rightarrow \mathbb{R}$$

$y = f(x)$

$x \in \mathbb{R}^n ; y \in \mathbb{R}$

↑ independent variables

↑ response variable

$$\text{Date: } \left\{ \begin{matrix} x^{(1)}, & x^{(2)}, & \dots, & x^{(N)} \\ y^{(1)}, & y^{(2)}, & \dots, & y^{(N)} \end{matrix} \right\} \in \mathbb{R}^n$$

$x^{(i)}, y^{(i)}$: Data pair

x_j : j^{th} value
of (i) th
observation

Model: Form a model to describe the relationship between x & y

$$y \approx \hat{f}(x)$$

where $\hat{f}: \mathbb{R}^n \rightarrow \mathbb{R}$ is the model.

Linear in parameters model:

$$\hat{f}(x) = \alpha_1 f_1(x) + \alpha_2 f_2(x) + \dots + \alpha_p f_p(x)$$

where $f_i: \mathbb{R}^n \rightarrow \mathbb{R}$ are basis functions
(feature mappings)

For a fixed $x \in \mathbb{R}^n$; \hat{f} is a linear $\hat{f} =$
in $(\alpha_1, \dots, \alpha_p)^T$; $\hat{f}(x) = (x_1, \dots, x_p)^T \begin{pmatrix} f_1(x) \\ \vdots \\ f_p(x) \end{pmatrix}$

Prediction error:

Our objective is to choose the model f so that it is consistent with the data i.e. $y^{(i)} \approx \hat{f}(x^{(i)})$ for $i=1, 2, \dots, N$

$$\text{Denote } \hat{y}^{(i)} = \hat{f}(x^{(i)})$$

$$\text{Residual: } r^{(i)} = y^{(i)} - \hat{y}^{(i)}$$

In vector notation: $y^d = (y^{(1)}, \dots, y^{(N)})^T$
 $\hat{y}^d = (\hat{y}^{(1)}, \dots, \hat{y}^{(N)})^T$

$$r^d = (y^{(1)} - \hat{y}^{(1)}, \dots, y^{(N)} - \hat{y}^{(N)})^T$$

Least square model fitting:

$$\text{minimize } \|r^d\|_2^2$$

$\alpha_1, \alpha_2, \dots, \alpha_p$

Express $\hat{y}^{(i)} = A_{i1}\theta_1 + \dots + A_{ip}\theta_p \quad i=1, 2, \dots, N$

where $A_{ij} = f_j(x^{(i)}) \quad \text{for } i=1, 2, \dots, N$

$$j=1, 2, \dots, p$$

$$A = [A_{ij}] \in \mathbb{R}^{N \times p}$$

Matrix - vector notation

$$\hat{y}^d = A \begin{bmatrix} \alpha_1 \\ \vdots \\ \alpha_p \end{bmatrix}$$

where $\alpha = \begin{bmatrix} \alpha_1 \\ \vdots \\ \alpha_p \end{bmatrix}$

$$\|x^d\|_2^2 = \|y^d - \hat{y}^d\|_2^2 = \|y^d - Ax^d\|_2^2 \\ = \|Ax^d - y^d\|_2^2$$

$$\min_{\alpha_1, \dots, \alpha_p} \|x^d\|_2^2 = \min_{\alpha_1, \dots, \alpha_p} \|Ax^d - y^d\|_2^2$$

$$\alpha = (A^T A)^{-1} A^T y^d = A^+ y^d$$

If A has linearly dependent columns, then one of the columns (let us assume b^{th} column) can be expressed as a linear combination of the remaining columns.

$$\frac{\|y^d - A\alpha\|^2}{N} : \text{ Mean square error.}$$

For $\hat{\alpha}$ mean square error is minimum.

$$\sqrt{\frac{\|y^d - A\alpha\|^2}{N}} : \text{ Root mean square (RMS).}$$

In general: $\varphi(r^d)$ may be used to determine how "closely" \hat{f} matches

the data $x^{(i)}, y^{(i)}$ $i=1, 2, \dots, N$.

For w: $\varphi(r^d) = \text{RMS}(r^d) = \sqrt{\frac{\|y^d - A\alpha\|^2}{N}}$ $\xrightarrow{\text{LS estimate}} \alpha$.

Least squares fit on a constant: (Interpretation of mean & variance)

$$p=1 \quad \text{and} \quad f_i(x) = 1 \quad \forall x \in \mathbb{R}^n$$

In this case, the model \hat{f} is constant.

$$\hat{f} = \alpha_1 f_1 = \alpha_1$$

$$A \in \mathbb{R}^{N \times 1} = \begin{pmatrix} 1 \\ 1 \\ \vdots \\ 1 \end{pmatrix}_{N \times 1} = 1_N$$

$$\hat{\alpha}_1 = (A^T A)^{-1} A^T y^d$$

$$= \frac{1}{N} 1_N^T y^d$$

$$\hat{\alpha}_1 = \text{avg}(y^d)$$

$$y^d : \begin{bmatrix} y^d \\ \vdots \\ y^{(N)} \end{bmatrix}$$

$$; \text{ RMS } (r^d) = \text{std}(y^d)$$

Univariate \mathbb{R}^n .

$n = 1$

$x \in \mathbb{R}^n$; $y \in \mathbb{R}$

$x \in \mathbb{R}$; $y \in \mathbb{R}$

Straight line fit:

$$f_1(x) = 1, \quad f_2(x) = x$$

$$\hat{f}(x) = \alpha_1 f_1(x) + \alpha_2 f_2(x)$$

$$\hat{f}(x) = \alpha_1 + \alpha_2 x$$

$$A = \begin{bmatrix} 1 & x^{(1)} \\ 1 & x^{(2)} \\ \vdots & \vdots \\ 1 & x^{(n)} \end{bmatrix} \in \mathbb{R}^{N \times 2}$$

$$y^d = \begin{bmatrix} y^{(1)} \\ \vdots \\ y^{(n)} \end{bmatrix}$$

Polynomial fit:

$$n=1 ; f_1(x) = 1, \quad f_2(x) = x; \quad f_3(x) = x^2, \\ \dots, f_p(x) = x^{p-1}; \quad f_{p+1}(x) = x^p$$

Dati: $(x^{(1)}, y^{(1)}), \dots, (x^{(n)}, y^{(n)})$

$$A = \begin{bmatrix} 1 & x^{(1)} & (x^{(1)})^2 & \dots & (x^{(1)})^p \\ 1 & x^{(2)} & (x^{(2)})^2 & \dots & (x^{(2)})^p \\ \vdots & \ddots & \vdots & \ddots & \vdots \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & x^{(n)} & (x^{(n)})^2 & \dots & (x^{(n)})^p \end{bmatrix} \in \mathbb{R}^{N \times (p+1)}$$

Vandermonde matrix

$\left(\text{ill-conditioned matrices} \right) = \frac{\max_{\text{mag}}(A)}{\min_{\text{mag}}(A)}$

log transform of dependent variable.

y is positive and varies over a large range.

$$w = \log y$$

use least squares to develop a model
for w.

$$\hat{w} = \hat{g}(x)$$

$$\hat{y} = e^{\hat{w}} = e^{\hat{g}(x)}$$

Generalization ability:

$x^{(1)} \quad y^{(1)}$, ..., $x^{(N)} \quad y^{(N)}$

