

Ordinary Least Squares

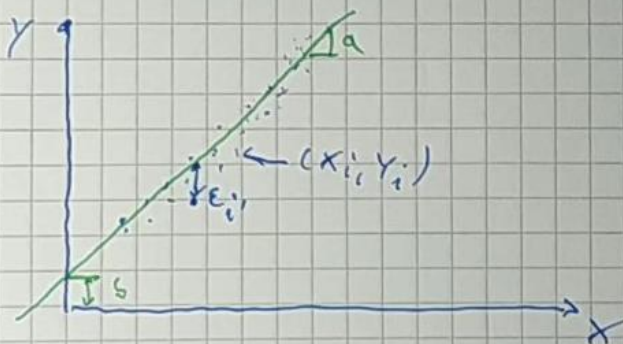
- simplest case: only one feature: $X \in \mathbb{R}$, $Y \in \mathbb{R}$

model: $Y = X \cdot \beta + \epsilon + b$

\uparrow \uparrow \uparrow
 slope noise intercept

assumes: $\epsilon \sim \mathcal{N}(0, \sigma^2)$

\uparrow \uparrow
 mean variance is independent of X_i or i
 same for all instances



$$Y = f(X, \eta)$$

$$\Leftrightarrow p(Y | X)$$

$$Y = X \cdot \beta + b + \epsilon$$

$$r = Y - (X\beta + b) \sim \epsilon \sim \mathcal{N}(0, \sigma^2)$$

deterministic eq of features and random number ϵ

posterior prob of ~~model~~ residual r is Gaussian

- determine optimal parameters β, b by maximum likelihood principle
 $\hat{=}$ choose $\hat{\beta}, \hat{b}$ such that likelihood of TS $p(TS | \hat{\beta}, \hat{b})$ is maximized,

simplify via the i.i.d assumption \Rightarrow probability of TS factorizes over i

$$p(TS | \beta, b) = \prod_{i=1}^N p((X_i, Y_i) | \beta, b) = \prod_{i=1}^N \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{1}{2} \frac{(Y_i - (X_i\beta + b))^2}{\sigma^2}\right)$$

instead of maximizing this, we minimize the negative log-likelihood (NLL)

$$-\log p(TS | \beta, b) = \sum_{i=1}^N \left[\frac{(Y_i - (X_i\beta + b))^2}{2\sigma^2} + \log \sqrt{2\pi\sigma^2} \right]$$

drop, because argument is independent of σ^2

$$\boxed{\hat{\beta}, \hat{b} = \underset{\beta, b}{\operatorname{argmin}} \sum_{i=1}^N (Y_i - (X_i\beta + b))^2}$$

ordinary least squares problem

To solve: set derivative of the loss L zero:

$$\text{Loss}(\{(x_i, y_i)\}_{i=1}^N) = \sum_{i=1}^N (y_i - (x_i \beta + b))^2$$

$$\frac{\partial \text{Loss}}{\partial b} = -2 \sum_{i=1}^N (y_i - (x_i \beta + b)) \stackrel{!}{=} 0$$

$$\sum_{i=1}^N y_i = \sum_{i=1}^N x_i \beta + \underbrace{\sum_{i=1}^N b}_{N \cdot b} \quad | : N$$

$$\frac{1}{N} \sum_{i=1}^N y_i = \bar{y} = \beta \cdot \bar{x} + b \quad \bar{x} = \frac{1}{N} \sum_{i=1}^N x_i$$

\Rightarrow regression line always goes through the center of mass (\bar{x}, \bar{y})

\Rightarrow when we center the data before doing regression, the regression line always goes through the origin $\Rightarrow \boxed{b = 0}$

\Rightarrow in the sequel, we assume centered data

$$\boxed{X \leftarrow X - \bar{X}, \quad Y \leftarrow Y - \bar{Y}}$$

centered loss:
$$\text{Loss} = \sum_{i=1}^N (y_i - x_i \beta)^2$$

$$\frac{\partial \text{Loss}}{\partial \beta} = -2 \sum_i x_i (y_i - x_i \beta) \stackrel{!}{=} 0$$

$$\left(\sum_{i=1}^N x_i^2 \right) \cdot \beta = \sum_{i=1}^N x_i y_i$$

$$\boxed{\hat{\beta} = \frac{\sum_i x_i y_i}{\sum_i x_i^2}}$$

Generalize to multiple features: $x_i \in \mathbb{R}^D$

x_i / β scalar product
column vector \mathbb{R}^D

$$\text{Loss} = \sum_i (y_i - x_i \beta)^2 = (Y - X \beta)^T (Y - X \beta) \quad \text{in matrix notation}$$

$$\frac{\partial \text{Loss}}{\partial \beta} = -2 X^T (Y - X \beta) \stackrel{!}{=} 0$$

$$\boxed{(X^T X) \beta = X^T Y} \quad \text{normal equations of OLS problem}$$

$$X: \begin{bmatrix} \text{---} \\ \text{---} \\ \text{---} \\ \text{---} \end{bmatrix}_N^D \quad Y: \begin{bmatrix} \text{---} \\ \text{---} \\ \text{---} \\ \text{---} \end{bmatrix}_N^1 \quad \beta: \begin{bmatrix} \text{---} \\ \text{---} \\ \text{---} \\ \text{---} \end{bmatrix}^D$$

- assume "true generative process" for data : $Y_i = X_i \beta^* + \epsilon_i$ $\epsilon_i \sim \mathcal{N}(0, \sigma^2)$
 β^* and ϵ^* unknown, must be learned
 variance independent of i

- if data is centered : $\frac{1}{N} \sum_i X_i = 0$ $\frac{1}{N} \sum_i Y_i = 0 \Rightarrow \hat{\beta} = 0$ regardless of β

\Rightarrow assume throughout that data is centered

- maximum likelihood loss \rightarrow negative log \rightarrow squared loss

$$\hat{\beta} = \arg \min_{\beta} \sum_i (Y_i - X_i \beta)^2 = \arg \min_{\beta} (Y - X\beta)^T (Y - X\beta)$$

- setting the derivative $\frac{\partial}{\partial \beta}$ to zero : normal equations
 $(X^T X) \cdot \beta = X^T Y$

- formal solution

$$\hat{\beta} = \underbrace{(X^T X)^{-1} X^T}_{X^+} Y$$

X^+ : pseudo-inverse, Moore-Penrose inverse

$X^+ = \underbrace{(X^T X)^{-1}}_S X^T$ is a generalization of the inverse matrix X^{-1} (which only exists if X is square ($N=D$)) to the case $N > D$

for the pseudo-inverse to exist, the scatter matrix $S = X^T X$ must be invertible

S is square $D \times D$ and is invertible if features are not redundant.

features are redundant if one feature can be expressed as a linear combination of the others :

$$X_j = \sum_{j' \neq j} w_{j'} X_{j'} \quad X_j \text{ is redundant}$$

one column of X

if no feature is redundant, X has full rank $\Rightarrow S^{-1}$ exists

Algorithms to solve the normal equations:

① Cholesky decomposition: factorize $S = R^T \cdot R$

R : upper triangular matrix

R^T : lower — " —

$$\begin{pmatrix} & & \neq 0 \\ & \neq 0 & \\ \oplus & & \end{pmatrix}$$

always exists if X (and therefore S)

have full rank, simple algorithm

$$\Rightarrow \text{normal equations} \quad \underbrace{R^T \cdot R \cdot \beta}_{=z} = X^T Y$$

two phases: — solve $R^T \cdot z = X^T Y$ via forward substitution for z
— solve $R \cdot \beta = z$ via backward substitution for β

$D=3$

$$\begin{aligned} r_{11} \cdot z_1 &= (X^T Y)_1 \\ r_{21} \cdot z_1 + r_{22} \cdot z_2 &= (X^T Y)_2 \\ r_{31} \cdot z_1 + r_{32} \cdot z_2 + r_{33} \cdot z_3 &= (X^T Y)_3 \end{aligned} \quad \left| \quad \begin{aligned} z_1 &= \frac{(X^T Y)_1}{r_{11}} \\ z_2 &= \frac{(X^T Y)_2 - r_{21} \cdot z_1}{r_{22}} \\ z_3 &= \frac{(X^T Y)_3 - r_{31} \cdot z_1 - r_{32} \cdot z_2}{r_{33}} \end{aligned} \right.$$

advantages: — easy

— S can be computed without ever creating full matrix X , is a good thing when N is very big

$$S = 0$$

for $i = 1 \dots N$

for $j = 1 \dots D$

for $j' = 1 \dots D$

$$S_{jj'} = S_{jj'} + X_{ij} \cdot X_{ij'}$$

only need a single row of X at a time

disadvantage: numerically less stable than other alg.

condition number of a matrix: $\kappa = \|X\| \|X^+\| \approx \frac{\lambda_{\max}}{\lambda_{\min}}$
same norm of X

- all eigenvalues of S non-negative
- if a feature is redundant, $\lambda_{\min} = 0$
- if it is close to redundant, ~~and~~

$$\lambda_{\min} \ll \lambda_{\max}$$

$\Rightarrow \kappa = \frac{\lambda_{\max}}{\lambda_{\min}}$ is large

\Rightarrow ~~compute~~ computing S^{-1} results in large numerical errors

$$S = X^T X \quad \Rightarrow \quad \kappa(S) = \kappa(X)^2$$

why is this a problem: rule of thumb for the numerical error of solving a linear system:

$$A \cdot b = c, \quad \text{where } \rightarrow$$

~~error in b~~
 number of valid digits in $b = \frac{\kappa(A) \cdot \mu}{-\log_{10}(\kappa \cdot \mu)}$
 \leftarrow machine precision of the floating point data type
 $\mu_{\text{float32}} = 10^{-7}$ $\mu_{\text{float64}} = 10^{-16}$

\Rightarrow if X has bad condition, e.g. $\kappa = 10^8$

S has very ^{bad} condition $\kappa(S) = \kappa(X)^2 = 10^{16}$

\Rightarrow even if we use float64 ("double") $\Rightarrow \kappa(S) \cdot \mu = 1$

\Rightarrow at best one valid digit \Rightarrow alternative alg. works with X directly \Rightarrow 8 valid digits

Alg (2) singular value decomposition of X

can always express $X = U \cdot \Lambda \cdot V^T$

\uparrow \uparrow \uparrow
 $N \times D$ $D \times D$ $D \times D$
 orthogonal $N \times D$ diagonal $D \times D$ orthogonal $D \times D$

elements of $\Lambda = \begin{pmatrix} \lambda_1 & & 0 \\ & \lambda_2 & \\ 0 & & \ddots \\ & & & \lambda_D \end{pmatrix}$ are called "singular values" of X
 $\hat{=}$ generalization of eigenvalues to rectangular matr.

solve normal eq. using SVD: $X^T = V \cdot \Lambda \cdot U^T$

$$S = X^T X = U \cdot \Lambda \cdot \underbrace{U^T \cdot U}_{=I} \cdot \Lambda \cdot V^T = V \cdot \Lambda^2 \cdot V^T$$

$$S^{-1} = (X^T X)^{-1} = \underbrace{V^{-T}}_{(V^T)^{-1} = (V^{-1})^T = V} \cdot \Lambda^{-2} \cdot \underbrace{V^{-1}}_{V^T} = V \cdot \Lambda^{-2} V^T$$

$$X^+ = (X^T X)^{-1} \cdot X^T = V \cdot \Lambda^{-2} \underbrace{V^T \cdot V}_{=I} \cdot \Lambda \cdot U^T = V \Lambda^{-1} U^T$$

$$\Rightarrow \boxed{\hat{\beta} = X^+ Y = V \Lambda^{-1} U^T \cdot Y}$$

compute SVD of X , invert the singular values $\Lambda^{-1} = \begin{pmatrix} \frac{1}{\lambda_1} & & 0 \\ & \ddots & \\ 0 & & \frac{1}{\lambda_D} \end{pmatrix}$, calculate $\hat{\beta}$

advantages: - work with X directly \Rightarrow numerically very stable

- also works for redundant features: some (or several) $\lambda_i = 0$
 \Rightarrow do not touch them when computing $\hat{\beta}$

- disadvantages of SVD:
 - needs X explicitly, problematic for big N
 - SVD decomposition alg. very complicated to implement
 - \Rightarrow always use library routine, not reprogram it

(3) LSQR algorithm [Paige & Saunders 1982]: when X is very big
 [SVD works on standard computers up to 1000×1000 matrices]

In practice, large feature matrices often have special structure,

- e.g.
- sparse; most elements of X are 0 (but X still full rank)
 - Toeplitz: every row is a shifted version of previous row

\Rightarrow handle X of size $10^6 \times 1000$ and above
 least-squares alg. takes advantage of the special structure
 decomposition

$$X = U B V^T$$

$\underbrace{\quad\quad\quad}_{\text{orthonormal}} \quad \underbrace{\quad\quad\quad}_{\text{upper bi-diagonal matrix}} \quad \underbrace{\quad\quad\quad}_{\text{orthonormal}}$

$$B = \begin{pmatrix} b_{11} & b_{12} & & 0 \\ & \ddots & \ddots & \\ 0 & & b_{p-1,p} & \\ & & & b_{pp} \end{pmatrix}$$

trick 1: compute the decomposition such that only two rows / columns of U and V are needed simultaneously at any time

trick 2: alg. never accesses X directly, but only via matrix-vector product subroutines; if it needs $X \cdot a \Rightarrow$ call `x-times-vec(a)`
 $X^T \cdot b \Rightarrow$ call `xt-times-vec(b)`

user must provide these subroutines to the LSQR program, these 146 routines are tailored to special structure of X