

→ Multiple input columns, 1 output column

→ $\hat{y}_n = \beta_0 + \beta_1 x_{n1} + \beta_2 x_{n2} + \dots + \beta_m x_{nm}$

↑ Hyperplane in n-dim coordinate system.

| | x_1 | $x_2 \dots x_m$ | y_n | \hat{y} |
|----------|----------|-----------------------|----------|-------------|
| 1 | x_{11} | $x_{12} \dots x_{1m}$ | y_1 | \hat{y}_1 |
| 2 | x_{21} | $x_{22} \dots x_{2m}$ | y_2 | \hat{y}_2 |
| \vdots | \vdots | $\vdots \dots \vdots$ | \vdots | \vdots |
| n | x_{n1} | $x_{n2} \dots x_{nm}$ | y_n | \hat{y}_n |

→ $\hat{y} = \begin{bmatrix} \hat{y}_1 \\ \hat{y}_2 \\ \vdots \\ \hat{y}_n \end{bmatrix} = \begin{bmatrix} \beta_0 + \beta_1 x_{11} + \beta_2 x_{12} + \dots + \beta_m x_{1m} \\ \beta_0 + \beta_1 x_{21} + \beta_2 x_{22} + \dots + \beta_m x_{2m} \\ \vdots \\ \beta_0 + \beta_1 x_{n1} + \beta_2 x_{n2} + \dots + \beta_m x_{nm} \end{bmatrix}_{n \times 1}$

→ $y = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix}_{n \times 1}$ $\hat{y} = \begin{bmatrix} \hat{y}_1 \\ \hat{y}_2 \\ \vdots \\ \hat{y}_n \end{bmatrix}_{n \times 1}$

→ $e^T e = [y_1 - \hat{y}_1 \ y_2 - \hat{y}_2 \dots y_n - \hat{y}_n] \begin{bmatrix} y_1 - \hat{y}_1 \\ y_2 - \hat{y}_2 \\ \vdots \\ y_n - \hat{y}_n \end{bmatrix}$

$= (y_1 - \hat{y}_1)^2 + (y_2 - \hat{y}_2)^2 + \dots + (y_n - \hat{y}_n)^2$

② $E = e^T e = \sum_{i=1}^n (y_i - \hat{y}_i)^2$ → minimize

Loss function

① $\hat{y} = X\beta$

$X = \begin{bmatrix} 1 & x_{11} & x_{12} & \dots & x_{1m} \\ 1 & x_{21} & x_{22} & \dots & x_{2m} \\ \vdots & \vdots & \vdots & \dots & \vdots \\ 1 & x_{n1} & x_{n2} & \dots & x_{nm} \end{bmatrix}_{n \times (m+1)}$ $\beta = \begin{bmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_m \end{bmatrix}_{(m+1) \times 1}$

→ $E = e^T e = (y - \hat{y})^T (y - \hat{y})$

$= (y^T - \hat{y}^T) (y - \hat{y})$

$= y^T y - y^T \hat{y} - \hat{y}^T y + \hat{y}^T \hat{y}$

③ $E = y^T y - 2\hat{y}^T y + \hat{y}^T \hat{y}$

Let $y = A, \hat{y} = B$

$A^T B = B^T A$

$(A^T B)^T = B^T (A^T)^T = B^T A$

$A^T B = (A^T B)^T$ \downarrow $A^T B = C$ \downarrow $C = C^T$ Prove symmetric

$\hat{y}^T \hat{y} = y^T X \beta$

$= (1 \times n) \times (n \times m) \times (m \times 1)$

$= 1 \times m \times n \times 1$

$= 1 \times 1 \rightarrow \text{scalar} \therefore A^T B \text{ is symmetric}$

→ $E = y^T y - 2\hat{y}^T y + \hat{y}^T \hat{y} \rightarrow \text{sub } \hat{y} = X\beta$

$= y^T y - 2y^T X\beta + (X\beta)^T (X\beta)$

④ $E = y^T y - 2y^T X\beta + \beta^T X^T X \beta$

* find such value of β matrix for which E is minimum

$\frac{dE}{d\beta} = 0$

→ $E = y^T y - 2y^T X\beta + \beta^T X^T X \beta$

Let $A = X^T X$, given A is symmetric

$(X^T X)^T = X^T (X^T)^T = X^T X$

$\frac{dE}{d\beta} = 0 - 2y^T X + 2\beta^T A = 0$

$= 0 - 2y^T X + 2\beta^T X^T X = 0$

$2\beta^T X^T X = 2y^T X$

$\beta^T X^T X = y^T X$

$\beta^T X^T X (X^T X)^{-1} = y^T X (X^T X)^{-1}$

$\beta^T I = y^T X (X^T X)^{-1}$

$\beta^T = y^T X (X^T X)^{-1}$

$\beta = \left[\frac{y^T X (X^T X)^{-1}}{A \quad B} \right]^T$

$\beta = (X^T X)^{-1} X^T y$

⑤ $\beta = (X^T X)^{-1} X^T y$ OLS Method

$(X^T X)^{-1}$ is symmetric

$(X^T X)^T = (X^T X)^{-1}$

Let $X^T X = A$

$AA^T = I$

$(A^T)^T A^T = I$

$(A^T)^T A = I$

$(A^T)^T A A^T = I A^T$

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

$\therefore (X^T X)^{-1} = (X^T X)^{-1}$

from sklearn.datasets import ...

What is the problem with OLS?

OLS → $(X^T X)^{-1}$

O(n⁴) slow / computationally intensive for high dimensionality data.

one of the reasons to opt for gradient descent