

Linear regression

Training data

	1, 2, ..., i, j, ..., p	
1		
2		
⋮		
i	$x_{i1} \dots x_{ij} \dots x_{ip}$	y_i
⋮		
n		

Supervised Learning

input $(x_{i1} \dots x_{ij} \dots x_{ip})$

output y_i supervision

want to learn mapping to apply to other testing situations

input \rightarrow output

Model:

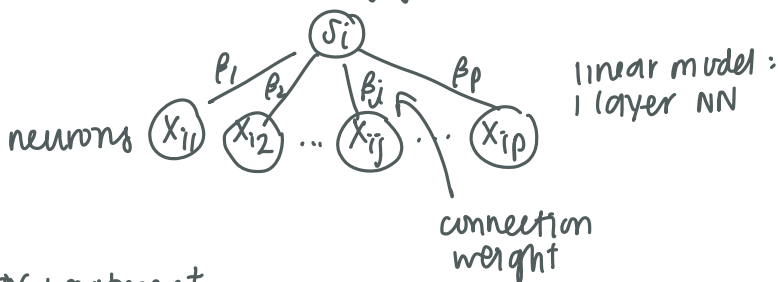
$$y_i = x_{i1}\beta_1 + \dots + x_{ij}\beta_j + \dots + x_{ip}\beta_p + e_i$$

$$s_i = x_{i1}\beta_1 + \dots + x_{ij}\beta_j + \dots + x_{ip}\beta_p$$

intercept/bias
 $s_i = \beta_0 + x_{i1}\beta_1 + \dots + x_{ij}\beta_j + \dots + x_{ip}\beta_p$

$$= \boxed{x_{i0}}\beta_0 + x_{i1}\beta_1 + \dots$$

||
1



Row vector treatment

$$x_i = \begin{pmatrix} x_{i1} \\ \vdots \\ x_{ij} \\ \vdots \\ x_{ip} \end{pmatrix} \quad \beta = \begin{pmatrix} \beta_1 \\ \vdots \\ \beta_j \\ \vdots \\ \beta_p \end{pmatrix}$$

	1, ..., p	
1		
i	$-x_i^T-$	y_i row i
⋮		
n		

$$s_i = \langle x_i, \beta \rangle = x_i^T \beta$$

$$(= \langle x_i, w \rangle + b)$$

\downarrow weight \downarrow β_0

$$e_i = y_i - s_i$$

$$\text{Loss}(\beta) = \frac{1}{2} \sum_{i=1}^n e_i^2 \quad \text{least squares loss function}$$

$$\frac{d\text{Loss}}{d\beta_k} = \sum_{i=1}^n e_i \frac{de_i}{ds_i} \frac{ds_i}{d\beta_k} \quad \text{simple error back propagation chain rule}$$

$$= -\sum_{i=1}^n e_i \cdot x_{ik} \quad \forall k=1, \dots, p$$

$$\text{Loss}'(\beta) = \begin{pmatrix} \frac{d\text{Loss}}{d\beta_1} \\ \vdots \\ \frac{d\text{Loss}}{d\beta_k} \\ \vdots \\ \frac{d\text{Loss}}{d\beta_p} \end{pmatrix} = -\sum_{i=1}^n \begin{pmatrix} x_{i1} e_i \\ \vdots \\ x_{ik} e_i \\ \vdots \\ x_{ip} e_i \end{pmatrix} = -\sum_{i=1}^n \begin{pmatrix} x_{i1} \\ \vdots \\ x_{ik} \\ \vdots \\ x_{ip} \end{pmatrix} e_i = -\sum_{i=1}^n \begin{matrix} p \times 1 \\ 1 \times 1 \end{matrix}$$

Estimating eqn

$$\sum_{i=1}^n x_i (y_i - x_i^T \beta) = 0$$

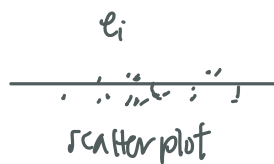
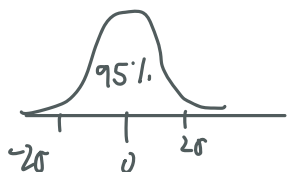
$$\sum_{i=1}^n x_i y_i - \sum_{i=1}^n x_i x_i^T \beta = 0$$

$$\begin{matrix} \boxed{\sum_{i=1}^n x_i x_i^T} \beta_{p \times 1} = \sum_{i=1}^n x_i y_i \\ \downarrow \quad \downarrow \quad \downarrow \quad \downarrow \\ p \times p \quad 1 \times p \quad p \times 1 \quad 1 \times 1 \\ \boxed{\quad} = \boxed{\quad} \end{matrix}$$

$$\hat{\beta}_{p \times 1} = \left(\sum_{i=1}^n x_i x_i^T \right)^{-1} \left(\sum_{i=1}^n x_i y_i \right)$$

Maximum likelihood

$$e_i \sim N(0, \sigma^2)$$



$$\text{likelihood}(\beta) = \prod_{i=1}^n p(y_i | s_i) \quad x_i^T \beta$$

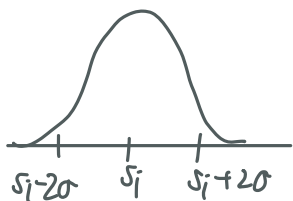
↓ max

Most plausible β

$$[y_i | s_i] \sim N(s_i, \sigma^2)$$

$$y_i = s_i + e_i$$

$$p(y_i | s_i) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(y_i - s_i)^2}{2\sigma^2}} \quad \text{prob density}$$



$$\text{log-likelihood}(\beta) = \sum_{i=1}^n \log p(y_i | s_i)$$

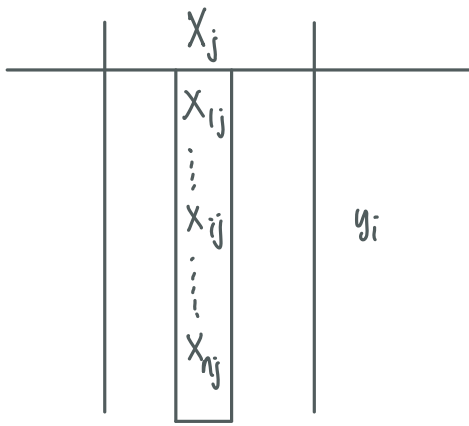
$$= \sum_{i=1}^n \left[-\frac{(y_i - s_i)^2}{2\sigma^2} - \frac{1}{2} \log(2\pi\sigma^2) \right]$$

$$= -\frac{1}{2\sigma^2} \sum_{i=1}^n (y_i - s_i)^2 - \frac{n}{2} \log(2\pi\sigma^2)$$

↓ max

$$\min \sum_{i=1}^n (y_i - s_i)^2 \quad \text{least squares objective function}$$

column vector view



$$X_{n \times p} = (X_1 \dots X_j \dots X_p)$$

$$\begin{pmatrix} y_1 \\ \vdots \\ y_i \\ \vdots \\ y_n \end{pmatrix} = \underbrace{\begin{pmatrix} X_{i1}\beta_1 + \dots + X_{ij}\beta_j + \dots + X_{ip}\beta_p \end{pmatrix}}_{s_i} + \begin{pmatrix} e_1 \\ \vdots \\ e_i \\ \vdots \\ e_n \end{pmatrix}$$

$$Y_{n \times 1} = X_1\beta_1 + \dots + X_j\beta_j + \dots + X_p\beta_p + e$$

pretend X_j as scalars

$$= (X_1 \dots X_j \dots X_p) \begin{pmatrix} \beta_1 \\ \vdots \\ \beta_j \\ \vdots \\ \beta_p \end{pmatrix} + e_{n \times 1}$$

reduce X_j as vector

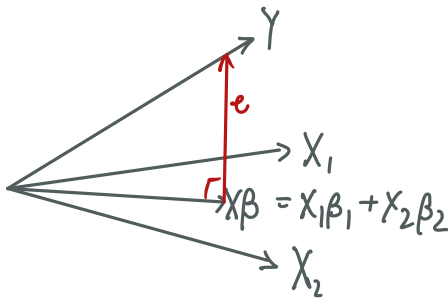
$$= X \beta + e$$

$n \times p \quad p \times 1 \quad n \times 1$

$$\text{Loss}(\beta) = \frac{1}{2} \sum_{i=1}^n e_i^2 = \frac{1}{2} |e|^2$$

$$\frac{d\text{Loss}}{d\beta_k} = - \sum_{i=1}^n e_i \cdot X_{ik}$$

$$= - \langle e, X_k \rangle \quad k=1, \dots, p$$



$$X_k^T e = 0 \quad \forall k=1, \dots, p$$

$$\begin{pmatrix} \vdots \\ X_k^T \\ \vdots \end{pmatrix} e = 0$$

$$\begin{pmatrix} \vdots \\ X_k^T \\ \vdots \end{pmatrix} e$$

X^T
 $p \times n$

$$X^T(Y - X\beta) = 0$$

$$X^T Y - X^T X \beta = 0$$

$$X^T X \beta = X^T Y$$

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

eq. for column vector view

rewrite/visit row vector view

$$\begin{pmatrix} x_i \\ \vdots \\ \dots i \dots n \end{pmatrix} \begin{pmatrix} x_i^T \\ \vdots \\ \dots n \end{pmatrix} = \sum_{i=1}^n x_i x_i^T$$

$$X^T X = \sum_{i=1}^n x_i x_i^T$$

$$X^T Y = \sum_{i=1}^n x_i y_i$$

$$\hat{\beta} = \left(\sum_{i=1}^n x_i x_i^T \right)^{-1} \left(\sum_{i=1}^n x_i y_i \right) \quad \text{eq for row vector view}$$

Logistic regression

	1	...	j	...	p	
1						$y_i \in \{0, 1\}$ $\{-, +\}$
\vdots						
i	x_{i1}	...	x_{ij}	...	x_{ip}	
\vdots						
n						

$$s_i = x_{i1}\beta_1 + \dots + x_{ij}\beta_j + \dots + x_{ip}\beta_p$$

$$= x_i^T \beta$$

extra

$$y_i = \text{Sign}(s_i)$$

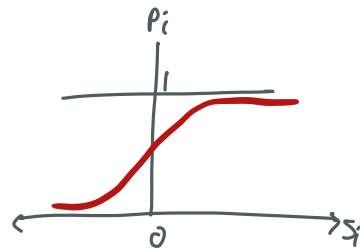
perception

soften

$$p_i = P_Y(y_i=1 | s_i)$$

$$= \text{Sigmoid}(s_i)$$

$$= \frac{e^{s_i}}{1 + e^{s_i}}$$



x least squares

$$\text{Loss}(\beta) = \sum_{i=1}^n (y_i - \text{sigmoid}(s_i))^2$$

✓ maximum likelihood

$$\text{log-likelihood}(\beta) = \sum_{i=1}^n \log P(y_i | s_i)$$