# Lecture 3

- Regression:

| i | $x_i^T$ | $y_i^T$ |
|---|---------|---------|
|   |         |         |

$$P(y \mid x) = \begin{cases} \dfrac{1}{(\sqrt{2\pi \sigma^2})^d} \; \overbrace{\exp\left(-\dfrac{1}{2\sigma^2} \|y-s\|^2\right)}^{\sim \mathcal{N}(s, \sigma^2 I_d)} & \text{Continuous} \;\to\; \text{regression} \\[4mm] \underbrace{\dfrac{e^{\langle y, s\rangle}}{z}}_{\sim \text{Multinomial}(p=softmax(s))} & \text{discrete / categorical} \;\to\; \text{classification} \end{cases}$$

- general model:

$s = f_\theta(x)$



```
        [    Y    ] d       output
            ↑
        [    S    ] d
            ↑ θ = (w,b)
        [    h    ]          hidden
            ↑ non-linear
        [    x    ]          input
        1 ... j .... p
```

$$s = wh + b$$

$$[s]_d = [w][h] + [b]$$

- log - likelihood $(\theta)$ $= \log\limits_{\theta} p(y|x)$

- Regression :

$$\ell(\theta) = -\frac{1}{2\sigma^2} |y - s|^2$$

$$L(\theta) = -\frac{1}{2\sigma^2} \sum_{i=1}^{n} |y_i - s_i|^2$$

- Classification

$$\ell(\theta) = \langle y, s \rangle - \log Z(s)$$

$$L(\theta) = \sum_{i=1}^{n} \left[ \langle y_i, s_i \rangle - \log z(s_i) \right]$$

$$y = \begin{pmatrix} 0 \\ \vdots \\ 0 \\ 1 \\ 0 \\ \vdots \\ 0 \end{pmatrix} \begin{matrix} \\ \\ \\ k \\ \\ \\ c \end{matrix} \qquad s = \begin{pmatrix} s_1 \\ \vdots \\ s_k \\ \vdots \\ s_c \end{pmatrix} \qquad p = \begin{pmatrix} p_1 \\ \vdots \\ p_k \\ \vdots \\ p_c \end{pmatrix} = \text{softmax}(s)$$

So $\langle y, s \rangle = s_k$

$$p_k = \frac{e^{s_k}}{\sum\limits_{c=1}^{c} e^{s_c}} = P(y=k|x)$$

<span style="color:blue">Least Squares</span>

$\text{Loss}(\theta) = - \log - \text{likelihood}(\theta) \Rightarrow \begin{cases} \text{Regression} : \frac{1}{2} |y - s|^2 \\ \text{Classification} : - \log p(y|x) \end{cases}$

$$= - \left( \langle y, s \rangle - \log z(s) \right)$$
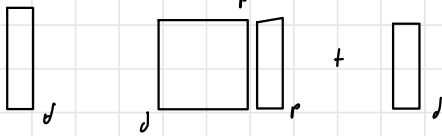
<span style="color:blue">Cross - entropy</span>

• Specialize to $d = 1$.

- Regression:

| | Father height | Mother height | Sons height |
|---|---|---|---|
| $i$ | $x_{i_1}$ | $x_{i_2}$ | $y_i$ |

- Recall linear model:

$$s \quad = \quad w \, x \quad + \quad b$$



$d = 1$

$$s = \langle w, x \rangle + b$$

$$s = x^T \beta + \beta_0 \quad \leftarrow \text{ statistical notation}$$

$$= (x_1 \cdots x_j \cdots x_p) \begin{pmatrix} \beta_1 \\ \beta_2 \\ \vdots \\ \beta_j \\ \vdots \\ \beta_p \end{pmatrix} + \beta_0$$
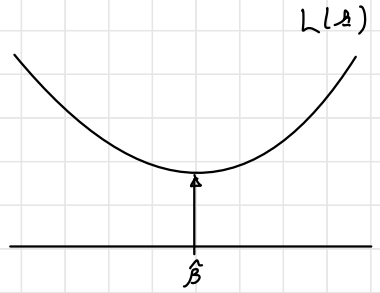
intercept

coefficients

$$= \beta_0 + \sum_{i=1}^{p} x_i \beta_i \qquad * \text{ Can make } \beta_0 \text{ ambiguous:}$$

$$s = x^T \beta$$

- Solving Least Squares Problem:

$$\text{Loss}(\beta) = \frac{1}{2} \sum_{i=1}^{n} (Y_i - S_i)^2$$

$$= \frac{1}{2} \sum_{i=1}^{n} e_i^2$$

$$= \frac{1}{2} \sum_{i=1}^{n} (Y_i - x_i^T \beta)^2$$


$L(\beta)$ — parabola with minimum at $\hat{\beta}$

$$\frac{\partial L}{\partial \beta} = \sum_{i=1}^{n} \frac{\partial L}{\partial e_i} \frac{\partial e_i}{\partial S_i} \frac{\partial S_i}{\partial \beta_k}, \qquad \text{note} \quad S_i = \sum_{j=1}^{p} x_{ij} \beta_j$$

$$= -\sum_{i=1}^{n} e_i x_{ik}$$

- Package:

$$\frac{\partial L}{\partial \beta} = \begin{pmatrix} \frac{\partial L}{\partial \beta_k} \end{pmatrix}_k^{1 \cdots p} = -\sum_{i=1}^{n} \begin{pmatrix} e_i x_{:ik} \end{pmatrix}_k^{1 \cdots p} = -\sum_{i=1}^{n} \underset{p\times 1}{x_i} \, \underset{1\times 1}{e_i}$$

$$= \underset{p\times 1}{\square} \, \underset{1\times 1}{\square}$$

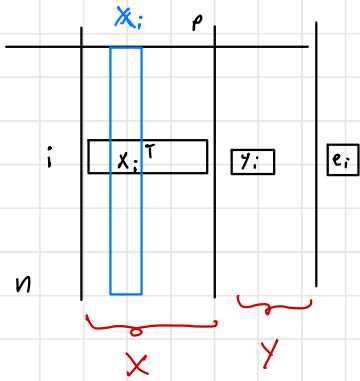$$\frac{\partial L}{\partial \beta} = 0 \implies \sum_{i=1}^{n} x_i e_i = 0$$

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

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

$$\sum x_i x_i^T \beta = \sum x_i Y_i$$

$$\underset{p\times p}{\boxed{\underset{p\times 1}{\square}\,\underset{1\times p}{\square}}}\,\underset{p\times 1}{\square} = \underset{p\times 1}{\square} \implies \hat{\beta} = \left( \sum_{i=1}^{n} x_i x_i^T \right)^{-1} \left( \sum_{i=1}^{n} x_i Y_i \right)$$

$$\sum_{i=1}^{n} x_i x_i^T = (x_1 \cdots x_i \cdots x_n) \begin{pmatrix} x_1^T \\ \vdots \\ x_i^T \\ \vdots \\ x_n^T \end{pmatrix}$$

$X^T$  $X$

meta-rule: Pretend each symbol is a number (scalar)

$$\sum x_i y_i = (x_1, \ldots, x_i, \ldots, x_n) \begin{pmatrix} y_1 \\ \vdots \\ y_i \\ \vdots \\ y_n \end{pmatrix}$$

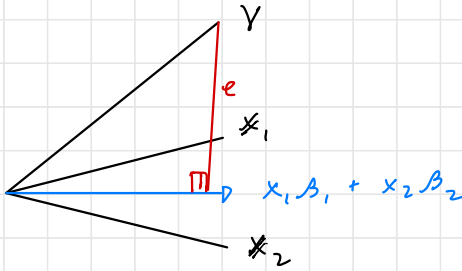$X^T$  $Y$

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

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

$$= \frac{1}{2} |e|^2$$

$$e = \begin{pmatrix} e_1 \\ \vdots \\ e_i \\ \vdots \\ e_n \end{pmatrix} = \begin{pmatrix} y_i - (x_{i1}\beta_1 + \cdots + x_{ip}\beta_p) \end{pmatrix}$$

$$= \frac{1}{2} \left| y - x_1\beta_1 + \cdots + x_j\beta_j + \cdots + x_p\beta_p \right|^2$$

$$= \frac{1}{2} \left| y - (x_1 \cdots x_j \cdots x_p) \begin{pmatrix} \beta_1 \\ \vdots \\ \beta_j \\ \vdots \\ \beta_p \end{pmatrix} \right|^2$$

$$= \frac{1}{2} | y - X\beta |^2$$

- Geometric View:



$$\frac{\partial L}{\partial \beta_k} = - \sum_{i=1}^{n} e_i x_{ik} = - \langle e, x_k \rangle = 0 \quad , \text{ so } \quad e \perp x_k$$

$$x_k^T e = 0$$

So $\begin{array}{c} 1 \\ \vdots \\ k \\ \vdots \\ p \end{array} \left( x_k^T e \right) = 0$

$$\left( x_k^T \right) e = 0$$

$$X^T e = 0$$

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

$$X^T Y - X^T X \beta = 0 \quad , \text{ Thus } \quad \hat{\beta} = (X^T X)^{-1} X^T Y$$

- Suppose $p(y_i | x_i) \sim \mathcal{N}(s_i, \sigma_i^2)$, different variance per observation. (precision)

$$p_\theta(y_i | x_i) = \frac{1}{\sqrt{2\pi\sigma_i^2}} \exp\left( -\frac{1}{2\sigma_i^2} (y_i - s_i)^2 \right)$$

$$\text{Loss}(\beta) = \frac{1}{2} \sum_{i=1}^{n} \frac{1}{\sigma_i^2} (y_i - x_i^T \beta)^2 = \frac{1}{2} \sum_{i=1}^{n} w_i (y_i - x_i^T \beta)^2$$

$$w_i = \frac{1}{\sigma_i^2} \quad \text{weight} \implies \text{weighted least squares}$$

$$\frac{\partial L}{\partial \beta} = -\sum_{i=1}^{n} w_i x_i \underbrace{(y_i - x_i^T \beta)}_{e_i} = 0$$

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

or

$$\tilde{x}_i = \sqrt{w_i} \, x_i$$
$$\tilde{y}_i = \sqrt{w_i} \, y_i$$

$$\Downarrow \quad \text{translate to Least squares}$$

$$\frac{1}{2} \sum (\tilde{y}_i - \tilde{x}_i^T \beta)^2$$

○ Logistic Regression:

| $x_i^T$ | $Y: \in \{0, 1\}$ |
|---|---|
| height, weight | gender |

$$P(Y|x) = \frac{e^{YS}}{Z} \begin{cases} p = \frac{e^s}{Z} = \frac{e^s}{1+e^s} & Y=1 \\ 1-p = \frac{1}{Z} = \frac{1}{1+e^s} & Y=0 \end{cases}$$

Thus $Z = 1 + e^s$

So $p = P(Y=1 | s) = \frac{e^s}{1+e^s} = \text{sigmoid}(s) = \sigma(s)$
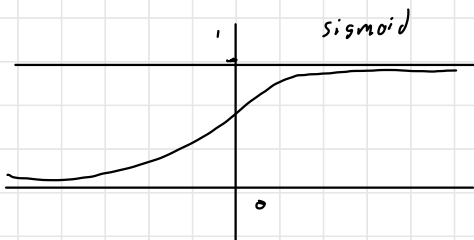
$$1 - p = \frac{1}{1+e^s}$$

$$\frac{p}{1-p} = e^s$$

$$\implies s = \log\left(\frac{1}{1-p}\right) = \text{logit}(p)$$

○ objective function:

$$\log\text{-likelihood}(\beta) = YS - \log(1 + e^s)$$

$$\log\text{-likli}(\beta) = \sum_{i=1}^{n} Y_i s_i - \log(1 + e^{s_i})$$

$$s_i = x_i^T \beta$$


sigmoid

Note derivatives of the tails of sigmoid approach 0.
Derivative is largest near 0

$$\sigma'(s) = \frac{d}{ds}\left(\frac{e^s}{1+e^s}\right)$$

$$= \frac{d}{ds}\left(1 - \frac{1}{1+e^s}\right)$$

$$= \frac{e^s}{[1+e^s]^2} = \frac{e^s}{1+e^s} \cdot \frac{1}{1+e^s}$$

$$\sigma'(s) = p(1-p)$$

- Iterated Reweighted Least Squares.

$$\beta_0 \longrightarrow \beta_1 \longrightarrow \cdots \longrightarrow \beta_t \longrightarrow \beta_{t+1} \longrightarrow \cdots$$

$$\ell(s_i) = y_i s_i - \log(1 + e^{s_i})$$

current $\beta_t \longrightarrow \hat{s}_i = x_i^T \beta_t \longrightarrow \hat{p}_i = \sigma(\hat{s}_i)$

$$\hat{w}_i = \hat{p}_i(1-\hat{p}_i)$$

$$\beta = \beta_t + \Delta\beta$$

$$s_i = x_i^T \beta = x_i^T(\beta_t + \Delta\beta)$$

$$= \underbrace{x_i^T \beta_t}_{\hat{s}_i} + \underbrace{x_i^T \Delta\beta}_{\Delta s_i}$$

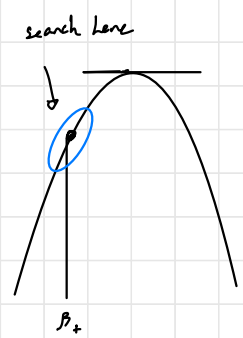$$\ell(s_i) = \ell(\hat{s}_i + \Delta s_i)$$

$$\approx \ell(\hat{s}_i) + \ell'(\hat{s}_i)\Delta s_i + \frac{1}{2}\ell''(\hat{s}_i)\Delta s_i^2$$

$$= \text{const} + \underline{\hat{e}_i}\,\Delta s_i - \frac{1}{2}\underline{\hat{w}_i}\,\Delta s_i^2$$

$$\underline{\ell'(s_i)} = y_i - \frac{e^{s_i}}{1+e^{s_i}} = y_i - p_i = e_i$$

$$\underline{\ell''(s_i)} = -p_i(1-p_i) = -w_i$$

(annotations in diagram): search line, $\beta_t$, $\sigma(s_i)$

$$= -\frac{1}{2}\hat{w}_i \left( \Delta s_i^2 - 2\frac{\hat{e}_i}{\hat{w}_i}\Delta s_i \right) + \text{const}$$

$$= -\frac{1}{2}\hat{w}_i \left( \underbrace{\Delta s_i}_{x_i^T \Delta \beta} - \boxed{\underbrace{\frac{\hat{e}_i}{\hat{v}_i}}_{\hat{y}_i}} \right)^2 + \text{const}$$

$$\text{Loss}(\Delta \beta) = \frac{1}{2}\sum_{i=1}^{n} w_i \left( \hat{y}_i - x_i^T \Delta \beta \right)^2$$

$$\text{Therefore} \quad \Delta \beta_t = \left( \sum_{i=1}^{n} \hat{w}_i \, x_i \, x_i^T \right)^{-1} \left( \sum_{i=1}^{n} \hat{w}_i \, \hat{y}_i \right)$$

Thus $\quad \beta_{t+1} = \beta_t + \Delta \beta_t$

$$\text{until} \quad |\Delta \beta_t| < \varepsilon$$