

Lecture FYS-
STK3155/4155,
September 21, 2023

$$E[y_i] = \sum_j x_{ij} \beta_j = x_i * \beta$$

$$\text{var}[y_i] = \sigma^2$$

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

$$y_i' = f(x_i) + \varepsilon_i'$$

$$y_i' \in (-\infty, +\infty)$$

$$x_i \in (-\infty, +\infty)$$

assumption: y_i are iid.

$$y_i' \sim N(x_i * \beta, \sigma^2)$$

$$P(y_i | x_i | \beta)$$

$$D = \{ (x_0, y_0), (x_1, y_1), \dots, (x_{n-1}, y_{n-1}) \}$$

$$= \frac{1}{\sqrt{2\pi\sigma^2}} \exp \left\{ -\frac{(y_i - x_i \cdot \beta)^2}{2\sigma^2} \right\}$$

$$P(D|\beta) = \prod_{i=0}^{n-1} P(y_i | x_i | \beta)$$

$$\hat{\beta} = \arg \max_{\beta \in \mathbb{R}^p} P(D|\beta)$$

$$\hat{\beta} = \arg \max_{\beta \in \mathbb{R}^P} \log P(D|\beta)$$

MLE
maximum likelihood estimator

$$\hat{\beta} = \arg \min_{\beta \in \mathbb{R}^P} [-\log P(D|\beta)]$$

$$\log \prod_{i=1}^n q_i = q_1 q_2 q_3 \dots q_n$$

$$\sum_{i=1}^n \log q_i$$

$$-\log(P(D|\beta)) =$$

$$-\log \left[\prod_{i=0}^{n-1} \frac{1}{\sqrt{2\pi\sigma^2}} \exp \left[-\frac{(y_i - x_i * \beta)^2}{2\sigma^2} \right] \right]$$

$$= \overbrace{\frac{n}{2} \log(2\pi\sigma^2)}^{\alpha}$$

$$+ \sum_{i=0}^{n-1} \frac{(y_i - x_i * \beta)^2}{2\sigma^2} = C(\beta)$$

$$\int \sum_{i=0}^{n-1} \log(2\pi\sigma^2)^{1/2}$$

$$C(\beta) = \alpha + \frac{1}{2\sigma^2} \| (y - X\beta) \|_2^2$$

$$\frac{\partial C(\beta)}{\partial \beta} = 0 \Rightarrow$$

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

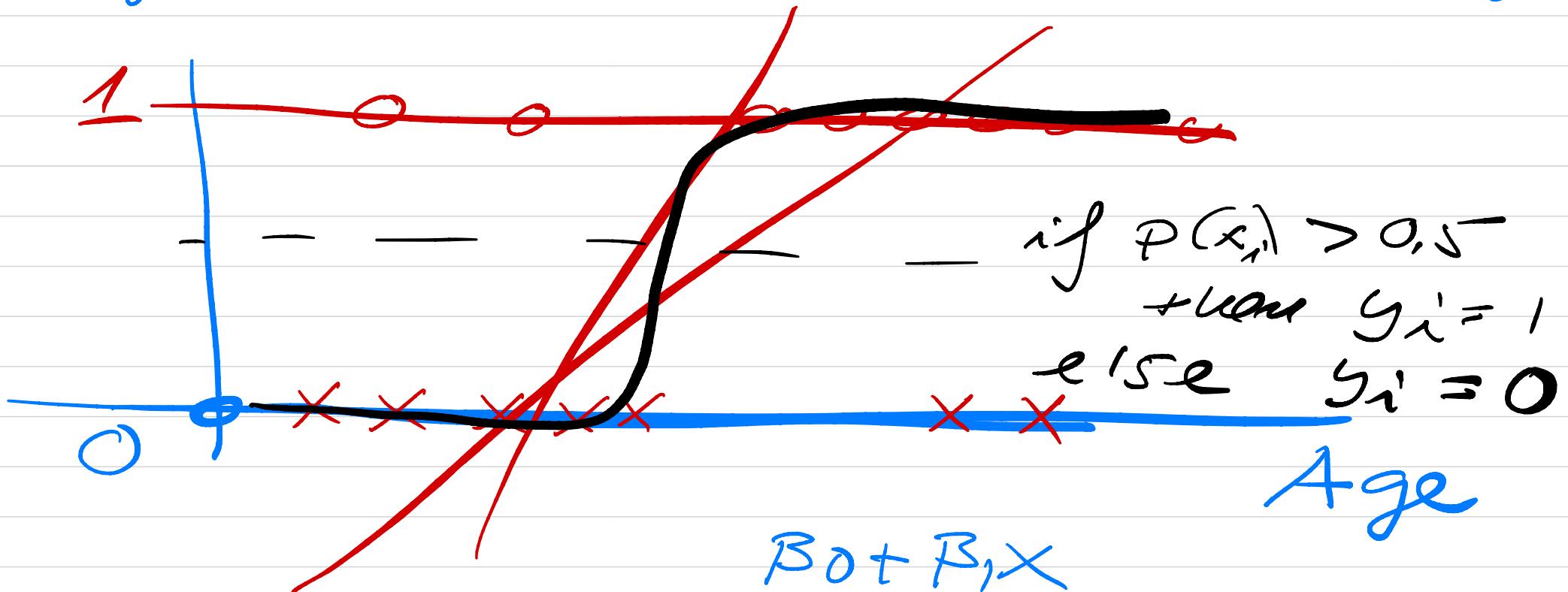
Logistic Regression :
Classification problem

Discrete output; binary

$$y_i = \{0, 1\}$$

$$y_i = p(x_i) + \epsilon_i$$

$f(x) \Rightarrow p(x)$, probability



$$p(x) = \frac{e}{1 + e^{\beta_0 + \beta_1 x}} \quad (\text{Sigmoid})$$

assumption

$$y(x) = p(x) + \varepsilon$$

$$p(x) \cong \frac{e^{\beta_0 + \beta_1 x}}{1 + e^{\beta_0 + \beta_1 x}}$$

$$\Rightarrow P(y_i | x_i | \beta) = p_i'$$

$$D = \{(x_0, y_0), (x_1, y_1), \dots, (x_{n-1}, y_{n-1})\}$$

$$P(D|\beta) = \prod_{i=0}^{n-1} p_i'$$

$$C(\beta) = -\log P(D|\beta)$$

$$\frac{\partial C}{\partial \beta} = 0$$

$y_i = 1$, probability p_i

$$y_i = 1 = p_i + \varepsilon_i \Rightarrow$$

$$\varepsilon_i = 1 - p_i$$

$$y_i = 0 = p_i + \varepsilon_i \Rightarrow \varepsilon_i = -p_i$$

what's the probability of ε ?

$$E[\varepsilon] = \sum_{i=0}^{n-1} p_i \varepsilon_i$$

$$p_i' = P$$

$$= (1-P)P + (-P)(1-P) = 0$$

$y_i = 0$ $P_{y_i=0} = 1 - p_i'$ p_i' is the probability for $y_i = 1$

$$\text{var}[\varepsilon] = (1-P)^2 P$$

↙
Binomial $D + (-P)^2(1-P)$

$$= P(1-P)$$

$$P(D|\beta) = \prod_{i=0}^{n-1} P_i^{y_i} (1-P_i)^{1-y_i}$$

$$\beta_0 + \beta_1 x_i$$

$$P_i = \frac{e}{1 + e^{\beta_0 + \beta_1 x_i}}$$

$$f(x_i) \simeq \tilde{y}_i = x_i * \beta$$

$$\hat{\beta} = \underset{\beta \in \mathbb{R}^P}{\arg \min} \sum \log P(D|\beta)$$

$$C(\beta) = - \sum_{i=0}^{n-1} [y_i \log p_i + (1-y_i) \log (1-p_i)]$$

$$\begin{aligned} y_i \log p_i &= y_i \log \left(\frac{e^{\beta_0 + \beta_1 x_i}}{1 + e^{\beta_0 + \beta_1 x_i}} \right) \\ &= y_i (\beta_0 + \beta_1 x_i) - y_i \log (1 + e^{\beta_0 + \beta_1 x_i}) \end{aligned}$$

$$C(\beta) = - \sum_{i=0}^{n-1} \left[y_i (\beta_0 + \beta_1 x_i) - \log (1 + e^{\beta_0 + \beta_1 x_i}) \right]$$

$$\frac{\partial C(\beta)}{\partial \beta_0} = - \sum_{i=0}^{n-1} (y_i - \beta_i) = g_0$$

$$\frac{\partial C(\beta)}{\partial \beta_1} = 0 = - \sum_{i=0}^{n-1} x_i (y_i - \beta_i) = -g_1$$

for the more general case

$$\frac{\partial C(\beta)}{\partial \beta} = 0 \Rightarrow x^T(P-y) = g$$

$x \in \mathbb{R}^{p \times n}$ $y, P \in \mathbb{R}^n$ $g \in \mathbb{R}^p$

$$\frac{\partial^2 C}{\partial \beta \partial \beta^T} = H = X^T W X$$

$$W = \begin{cases} w_{ii} = p_i(1-p_i) \\ 0 \quad \text{if } i \neq j \end{cases}$$

$$\beta \text{ is in } p_i = \frac{e^{\beta_0 + \beta_1 x_i}}{1 + e^{\beta_0 + \beta_1 x_i}}$$

$$\frac{\partial C}{\partial \beta} = 0 \Rightarrow X^T(\beta - y) = 0$$

non-linear equation in β

$$g = \bar{x}^T(\bar{P} - \bar{y}) = g(\bar{\beta}) + \text{triangle } P(\bar{\beta})$$

$$H = \bar{x}^T \underbrace{w}_{w(\beta)} \bar{x} = H(\beta)$$

Taylor expansion of $C(\hat{\beta})$
around $\hat{\beta}^n$ $\hat{\beta} - \beta^n$
 n refers to iteration - $n -$

$$\frac{\partial C}{\partial \beta} = \nabla_{\beta} C = 0$$

$$CC(\hat{\beta}) = CC(\beta^{(n)}) + (g^{(n)})^T (\hat{\beta} - \beta^{(n)})$$

↑
known

$$+ \frac{1}{2} (\hat{\beta} - \beta^{(n)})^T H(\beta^{(n)}) (\hat{\beta} - \beta^{(n)})$$

+ ...

$\left. \frac{\partial^2 C}{\partial \beta \partial \beta^T} \right|_{\beta=\beta^{(n)}} = H(\beta^{(n)})$

$$b = \hat{\beta} - \beta^{(n)} \Rightarrow$$

$$C(\hat{\beta}) = C(\beta^{(n)}) + (\mathbf{g}^{(n)})^T \mathbf{b}$$

$$+ \frac{1}{2} \mathbf{f}^T H^{(n)} \mathbf{b} + \dots$$

Keep to second order in b
only

$$\frac{\partial C}{\partial b} = \mathbf{g}^{(n)} - H^{(n)} \cdot \mathbf{b} = 0$$

$$b = \hat{\beta} - \beta^{(n)} \Rightarrow$$

$$= \beta^{(n+1)} \equiv \hat{\beta} - (H^{(n)})^{-1} \mathbf{g}^{(n)}$$

For logistic regression

$$H^{(n)} = H(\beta^{(n)}) = \bar{X}^T W(\beta^{(n)}) X$$

$$g^{(n)} = \bar{X}^T (P(\beta^{(n)}) - y)$$

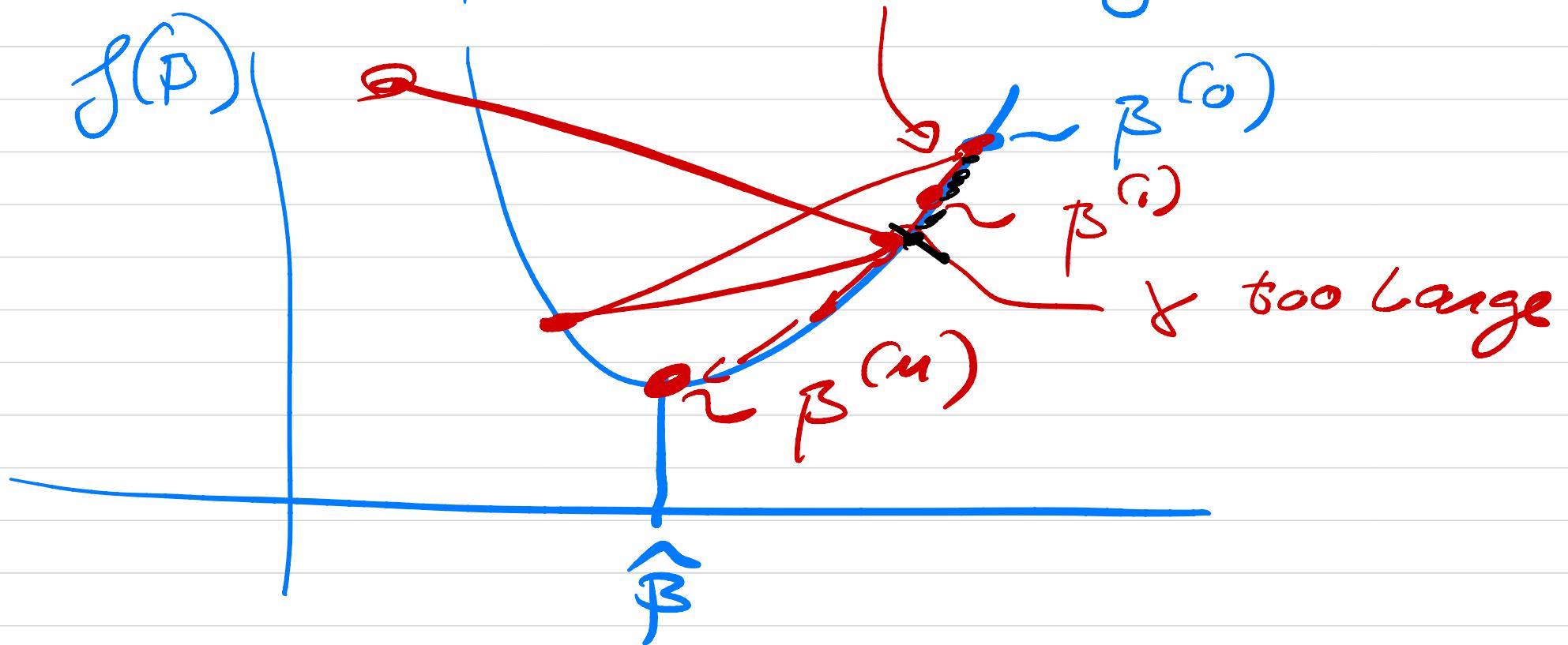
$$\beta^{(n+1)} = \beta^{(n)} - (\bar{X}^T W(\beta^{(n)}) \bar{X})^{-1} \times \\ \bar{X}^T (P(\beta^{(n)}) - y)$$

Newton-Raphson's method

start with a guess $\beta^{(0)}$ and
iterate n-times till $g=0$

$H^{-1} \Rightarrow \gamma^{(n)} = \text{learning rate}$

$$\beta^{(n+1)} = \beta^{(n)} - \gamma^{(n)} \cdot g^{(n)}$$



optimal rate $\gamma < \frac{2}{\lambda_{\max}}$

λ_{\max} is the largest eigenvalue
of H