

Lecture FYS-
STK3155/4155,
September 16, 2024

FYS-STK 3155/4155, September 16

Kaggle.com
Linear regression

$$y_i = f(x_i) + \varepsilon_i \sim N(\mu, \sigma^2)$$

$x_i \in (-\infty, \infty)$ $[x_a, x_b]$

$$y_i \in (-\infty, \infty) \quad [y_a, y_b]$$

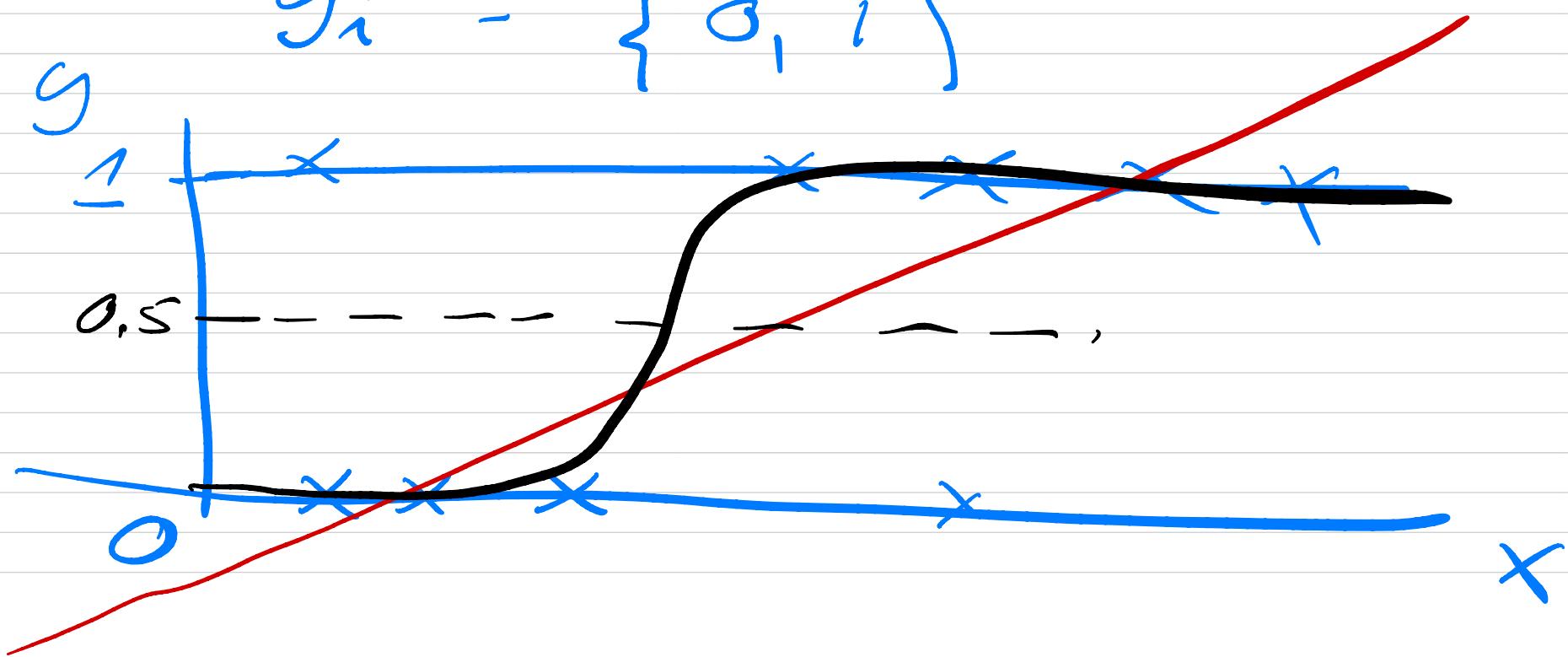
$$y_i \cong \underbrace{\sum_j x_{ij} \beta_j}_{\tilde{y}_i} + \varepsilon_i$$

First order polynomial

$$y_i \approx \beta_0 + \beta_1 x_i + \varepsilon_i$$

What about y_i given by discrete values?

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



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

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

$$\int_{x \in D} p(x) dx = 1 \quad \left\{ \sum_{i \in D} p(x_i) = 1 \right.$$

$$x \in [0, 1] \quad p(1) = 1$$

$$p(x_i) \leq p(x_j) \quad x_i \leq x_j$$

Example : sigmoid function

$$P(x) = \frac{e^{\beta x}}{1 + e^{\beta x}}$$

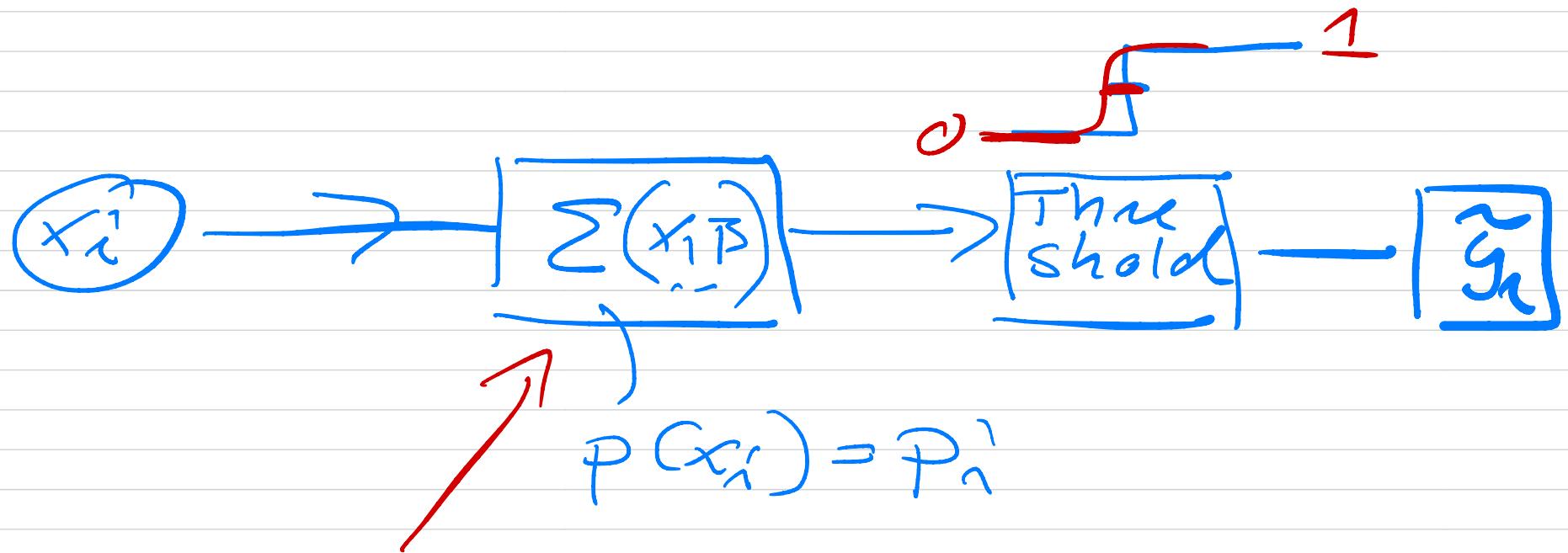
$$x \in [0, a] \quad \beta = 1$$

$$\int_0^a dx \frac{e^x}{1 + e^x} = \ln(1 + e^a)$$

$$y = 1 ; P(x) = \frac{e^{\beta x}}{1 + e^{\beta x}}$$

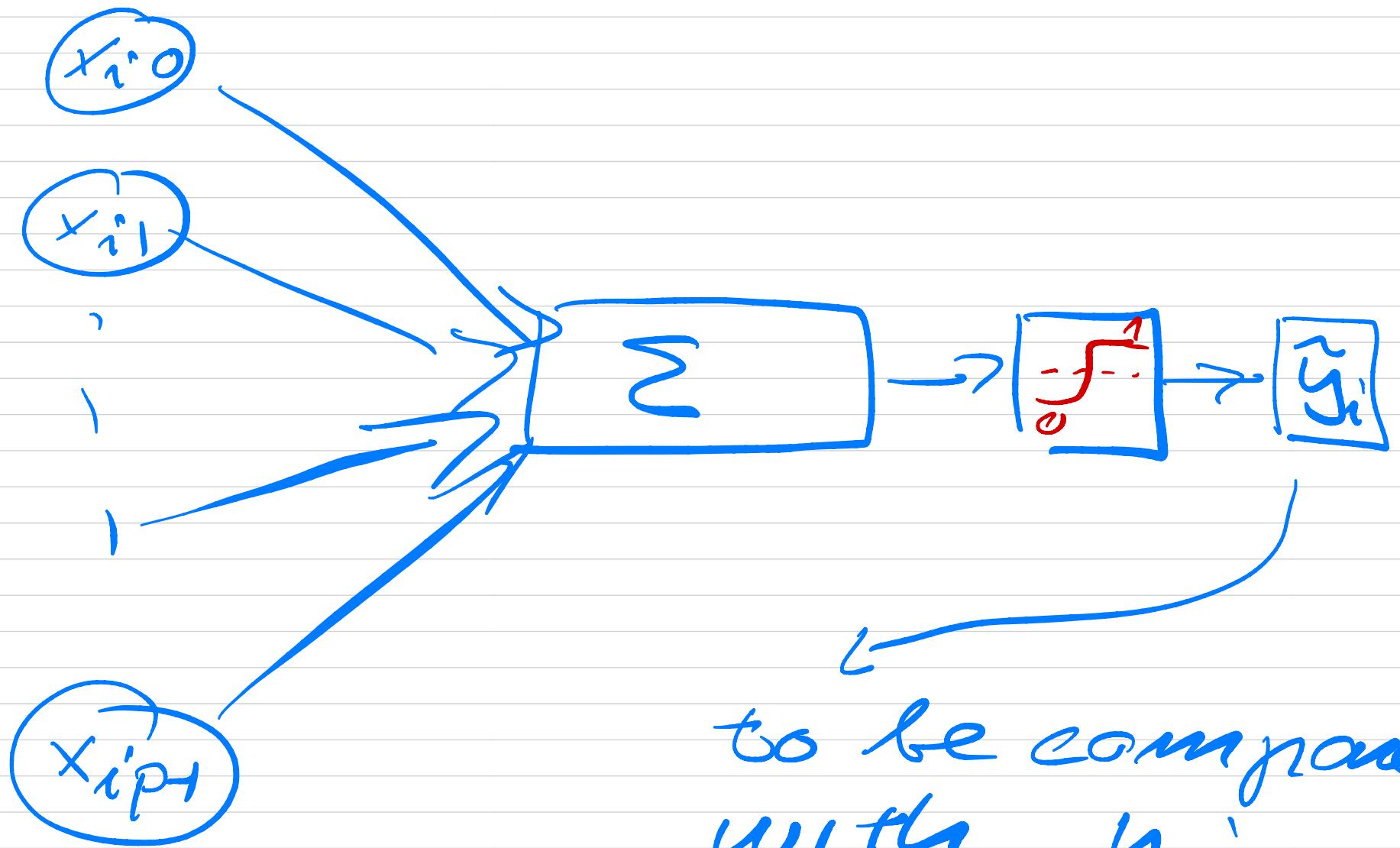
$$y = 0 ; 1 - P(x) = \frac{1}{1 + e^{\beta x}}$$

Graphical representation



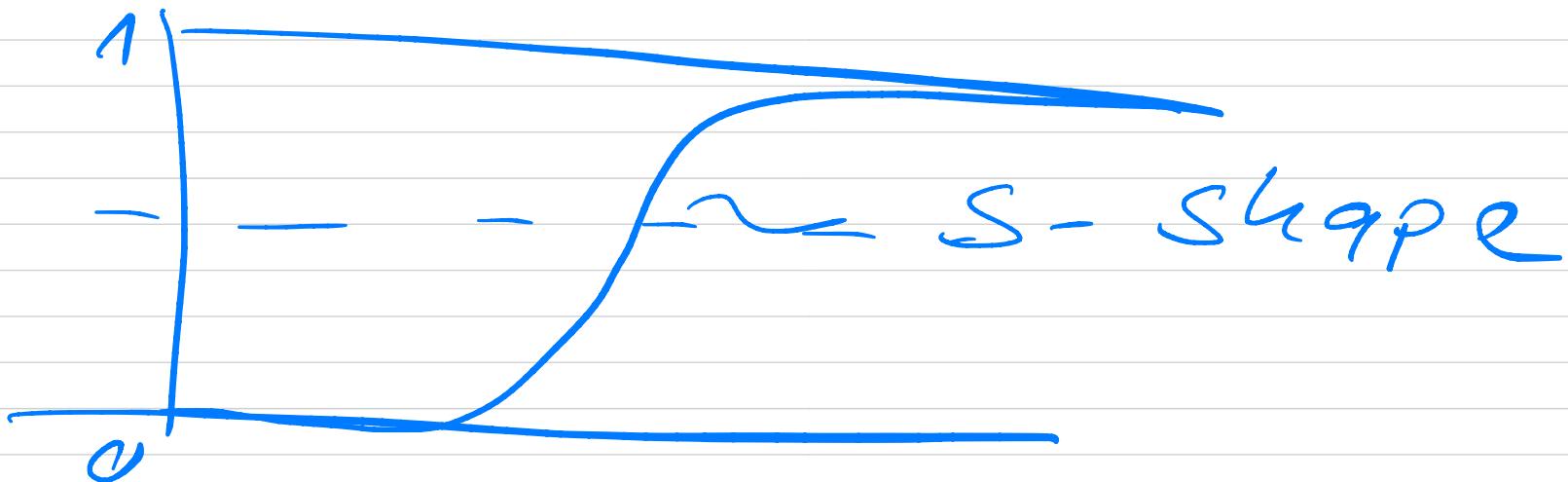
Neural networks
activation
function

$$X = \begin{bmatrix} x_{00} & x_{01} & \cdots & x_{0, p-1} \\ x_{10} \\ x_{20} \\ \vdots \\ x_{i0} & x_{i1} & x_{i2} & \cdots & x_{i, p-1} \\ \vdots \\ x_{n-10} & - & - & - & x_{n-1, p-1} \end{bmatrix}$$



to be compared
with y_i

$$P(x_i) = \frac{e^{\beta_0 + \beta_1 x_i}}{1 + e^{\beta_0 + \beta_1 x_i}}$$



$$y_i' = 1 ; P(x_i) = \frac{e^{\beta_0 + \beta_1 x_i}}{1 + e^{\beta_0 + \beta_1 x_i}}$$

$$y_i' = \pi_i + \varepsilon_i' = P_1'$$

what probability does ε_1'
follow?

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

$$y_i' = 0 = p_i' + \varepsilon_i'$$

$$\varepsilon_i' = 1 - p_i' \quad \text{when } y_i' = 1$$

$$\varepsilon_i' = -p_i' \quad \text{when } y_i' = 0$$

$$E[\varepsilon_i] = p_i \underbrace{(1-p_i)}_{\varepsilon_i} + (-p_i)(1-p_i)$$

$$(E[X] = \sum_{i=0}^{n-1} x_i p(x_i))$$

$$= 0$$

$$\text{Var}[\varepsilon_i] = (1-p_i)^2 p_i + (-p_i)^2 (1-p_i)$$

$$= p_i(1-p_i)$$

ε_i follows a binomial distribution.

How do we find β ?

y_i are i.i.d.

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

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

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

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

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

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

\Rightarrow

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

$$\log \left[\frac{e^{\beta_0 + \beta_1 x_1'}}{1 + e^{\beta_0 + \beta_1 x_1'}} \right]$$

$$= (\beta_0 + \beta_1 x_1') - \log(1 + e^{\beta_0 + \beta_1 x_1'})$$

\Rightarrow

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

$$\frac{\partial C}{\partial \beta_0} = 0 \quad \text{and} \quad \frac{\partial C}{\partial \beta_1} = 0$$

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

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

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

in general for the
gradient $g = \frac{\partial C}{\partial \beta} =$

$$\frac{\partial c}{\partial p} = -X^T(y - p)$$

↙ ↘

$y, p \in \mathbb{R}^n$

$X \in \mathbb{R}^{m \times n} \Rightarrow$

$$g = \frac{\partial c}{\partial p} \in \mathbb{R}^P$$

\hat{p} is included in p ?

in linear regression

$$\frac{\partial^2 C}{\partial \beta \partial \beta^T} = \frac{2}{n} \underbrace{X^T X}_{\text{Hessian matrix}}$$

Hessian matrix

independent of β

$$\frac{\partial^2 C}{\partial \beta \partial \beta^T} = X^T W X \text{ depends on } \beta,$$

Log Reg

$$W_{ii} = p_i(1-p_i)$$
$$W_{ij} (i \neq j) = 0$$

How do we solve the eqs for
 $\hat{\beta}$?

Newton-Raphson's method.

β - scalar , $\frac{\partial c}{\partial \beta} \Rightarrow \frac{dc}{d\beta}$

Taylor-expansion around
the optimal $\hat{\beta}$

$$\hat{\beta} - \beta^{(m)}$$

$$c(\hat{\beta}) = c(\beta^{(n)}) +$$

$$\underbrace{g(\beta^{(n)})}_{\text{(1)}} (\hat{\beta} - \beta^{(n)})$$

$$\left. \left(\frac{dc}{d\beta} \right) \right|_{\beta = \beta^{(n)}}$$

$$+ \frac{1}{2} \underbrace{\frac{\partial^2 c}{\partial \beta^2}}_{\text{(2)}} \left(\hat{\beta} - \beta^{(n)} \right)^2 + \dots$$

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

$$C(\tilde{\beta}) \approx C(\beta^{(n)}) + \underbrace{g^{(n)} b}_{g(\beta^{(n)})}$$

$$+ \frac{1}{2} \gamma^2 \underbrace{H^{(n)}}_{H(\beta^{(n)})}$$

$$\frac{dC}{db} = g^{(n)} + H^{(n)} b = 0$$

$$b = \tilde{\beta} - \beta^{(n)} = -(H^{(n)})^{-1} g^{(n)}$$

$$\tilde{\beta} = \beta^{(n+1)}$$

$$\boxed{\beta^{(m+1)} = \beta^{(m)} - (H^{(m)})^{-1} g^{(m)}}$$

$$m = 0, 1, \dots$$

- (i) start with a guess $\beta^{(0)}$
- (ii) iterate over m till

$$|\beta^{(m+1)} - \beta^{(m)}| \leq \varepsilon \approx 10^{-8}$$

in general, $H(\beta^{(n)})$ is
a matrix, $g(\beta^{(n)})$ is a
vector

$$g \in \mathbb{R}^P$$

$$H \in \mathbb{R}^{P \times P}$$

$$\beta \in \mathbb{R}^P$$

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

Replace $H(\beta^{(n)})$ with
a parameter $\gamma^{(n)}$

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

learning rate

γ

Can we find an optimal

γ ?

Expand $c(p)$ around

$$p^{(n)} - \gamma g^{(n)}$$

$$c(p^{(n)} - \gamma g^{(n)}) = c(p^{(n)})$$

$$- \gamma (g^{(n)})^T (g^{(n)})$$

$$+ \frac{1}{2} \gamma^2 [g^{(n)}]^T H^{(n)} g^{(n)} + \dots$$

$$g^{(n)} \rightarrow g$$

$$H^{(n)} \rightarrow H$$

Derivative wrt χ

$$\frac{\partial C}{\partial \chi} = 0 \Rightarrow$$

$$\chi = \frac{g^T g}{g^T H g} \quad (\chi^{(n)})$$

Assume

$$Hg = \lambda g$$

$$\chi = \frac{g^T g}{g^T \lambda g} = \frac{1}{\lambda} \Rightarrow$$

$$x_{\min} = \frac{1}{\lambda_{\max}}$$

$$x_{\max} = \frac{1}{\lambda_{\min}}$$

If you can diagonalize

$$H^{(m)}$$
$$\delta < \frac{2}{\lambda_{\max}}$$