

Lecture January 22

Linear Regression

$$y_i = X_i \beta + \varepsilon_i$$

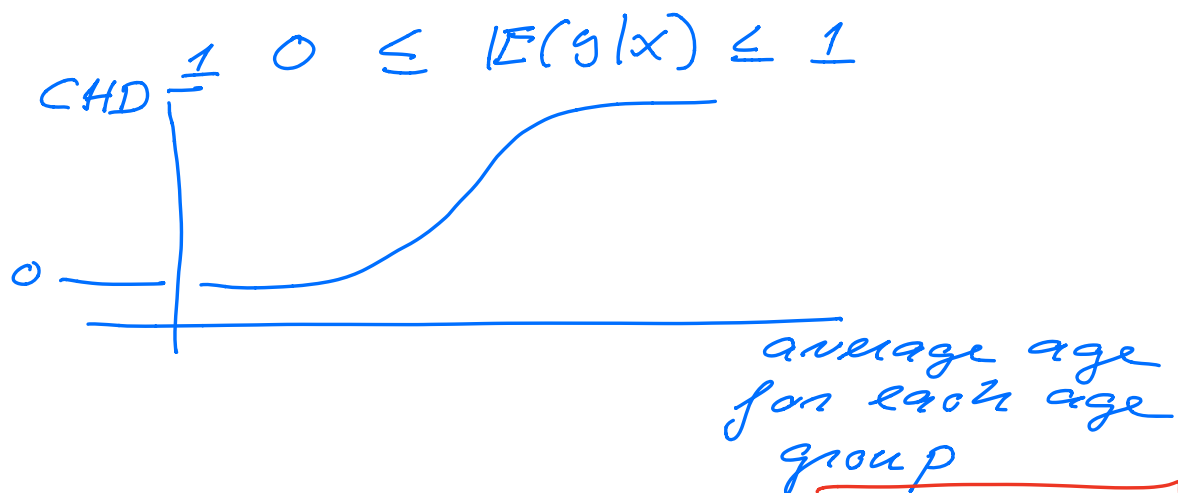
$$E[y_i] = X_i \beta = \beta_0 + \beta_1 x_i$$

(simplest choice)

$$= E(y|x, \beta)$$

This is a continuous function.

Binary case



$$p(x) = p(y|x) = E(y|x) = \frac{e^{\beta_0 + \beta_1 x}}{1 + e^{\beta_0 + \beta_1 x}}$$

Linear Regression

$$y = \boxed{X\beta} + \varepsilon$$

$$\underline{\varepsilon \sim N(0, \sigma^2)}$$

Logistic regression

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

y takes two values 0, 1

$$y = 1 \Rightarrow \varepsilon = \underline{1 - p(x)} \quad \text{with } p(x)$$

$$y = 0 \Rightarrow \varepsilon = \underline{-p(x)} \quad \text{with probability } 1 - p(x)$$

$$\sum_{i=1}^2 p(x_i) = p(x) + (1 - p(x)) = 1$$

$$E[\varepsilon] = \bar{\varepsilon} = (1 - p(x))p(x) - p(x)(1 - p(x)) \equiv 0$$
$$\left(E[x] = \sum p(x_i) x_i \right)$$

$\text{var}[\varepsilon^2] = p(1-p)$,
corresponds to the binomial distribution.

$$y_i = 1 \quad ; \quad p(y_i | x_i, \beta)$$

$$y_i = 0 \quad ; \quad 1 - p(y_i | x_i, \beta)$$

$$p(y_i | x_i, \beta) = p(x_i)$$

probability:

$$p(x_i)^{y_i} (1 - p(x_i))^{1-y_i}$$

$$C(\beta) = \prod_{i=0}^{n-1} p(x_i)^{y_i} (1 - p(x_i))^{1-y_i}$$

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

$$\beta^{\text{opt}} = \underset{\beta \in \mathbb{R}^p}{\text{argmin}} [-C(\beta)]$$

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

$\beta_0 + \beta_1 x_i$

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

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

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

\Rightarrow (in matrix-vector form)

$$\boxed{X^T(y - p)} \quad \begin{array}{l} p, y \in \mathbb{R}^n \\ X \in \mathbb{R}^{n \times p} \end{array}$$

$$p = \{ (p(x_0), p(x_1) \dots p(x_{n-1})) \}$$

$$X^T(y - p) \in \mathbb{R}^p$$

$$\beta \in \mathbb{R}^p \quad (p=2, \beta_0, \beta_1)$$

$$X^T(y - p) = 0$$

is a non-linear equation
in β .

2nd-derivative of $C(\beta)$
wrt to β

$$\frac{\partial^2 C(\beta)}{\partial \beta \partial \beta^T} = \frac{X^T W X}{= \text{Hessian}} = H$$

$$\left(W = p(y_i | x_i; \beta)(1 - p(y_i | x_i; \beta)) \right. \\ \left. \text{only diagonal elements} \right)$$

$$W \in \mathbb{R}^{m \times m}$$

$$H \in \mathbb{R}^{p \times p}.$$

H is a positive-definite matrix.

Newton-Raphson for root searching $(x^T(g-p)=0)$

$$f(x) = 0$$

$$x^{(k+1)} = x^{(k)} - \left[\frac{f(x)}{f'(x)} \right]_{x=x^{(k)}}$$

$$\|x^{(k+1)} - x^{(k)}\|_2 \leq \delta \sim 10^{-10}$$

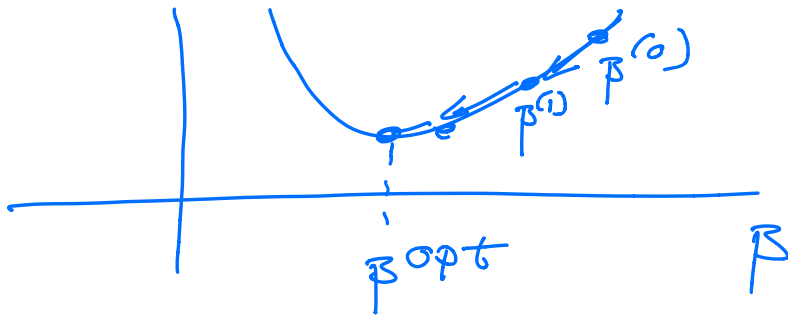
$$\beta^{(k+1)} = \beta^{(k)} - \left(\underset{\substack{\uparrow \\ p \text{ large}}}{H^{-1}} \frac{\partial C}{\partial \beta} \right)_{\beta=\beta^{(k)}}$$

p large

$$\boxed{\beta^{(k+1)} = \beta^{(k)} - \delta_k \left[\frac{\partial C}{\partial \beta} \right]_{\beta=\beta^{(k)}}}$$

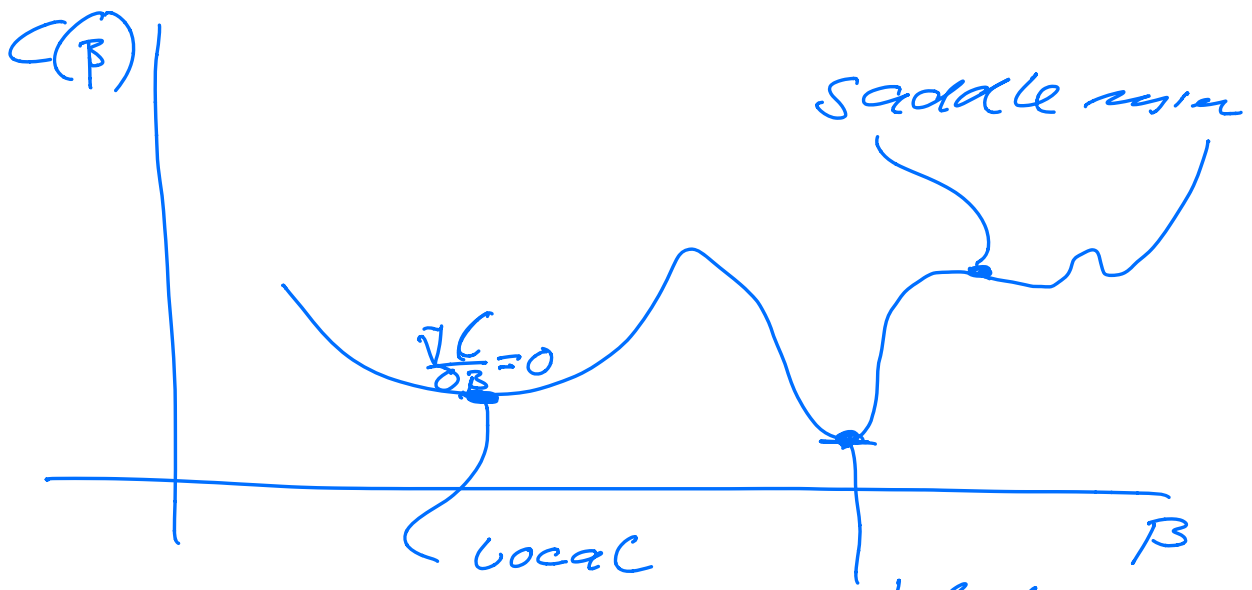
new hyperparameter
= Learning rate,

$$\gamma_k \Rightarrow \gamma = \{10^{-5}, 10^{-4}, 10^{-3}, \dots\}$$



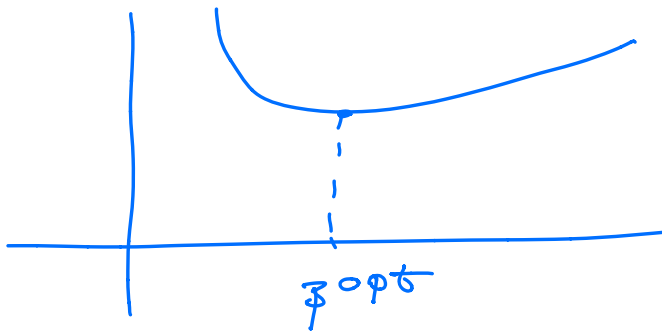
accuracy score

$$\text{score} = \frac{\sum_{i=0}^{n-1} I(y_i = \hat{y}_i)}{\# \text{ events}}$$



min

global
min

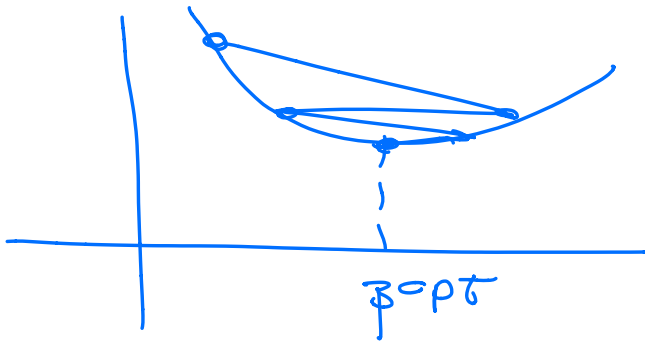


$$\gamma < \frac{2}{\lambda_{\max}}$$

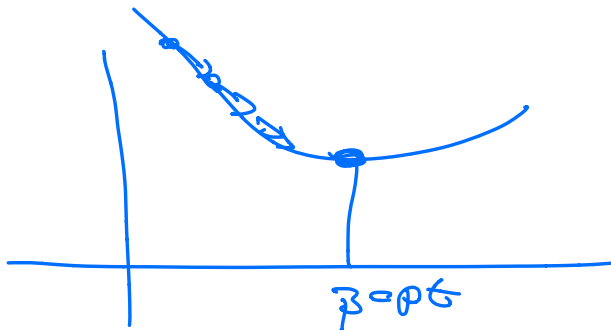
if $\gamma^{opt} < \gamma < 2\gamma^{opt}$ longest
value of

$$\gamma^{opt} = \frac{1}{\lambda_{\max}}$$

Hessian
matrix



$$\gamma < \gamma^{opt}$$



.. . opt

11 $x > x^{-1}$, no convergence