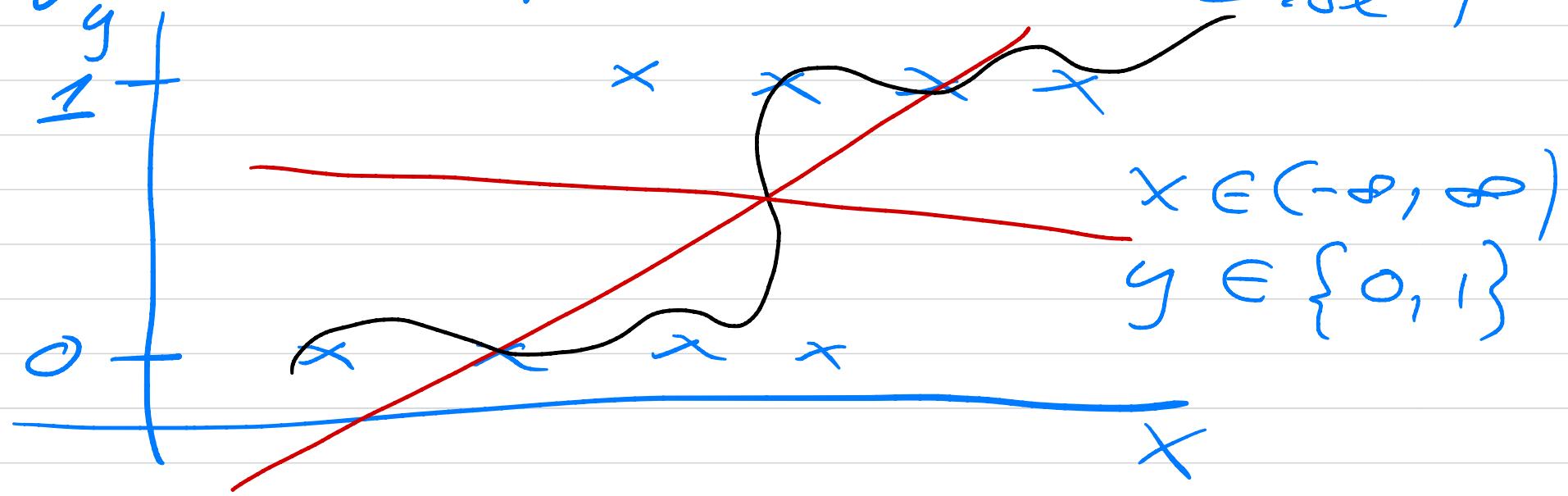


**Erasmus+ lecture on
machine learning,
November 6, 2023**

Logistic Regression (Binary case)



Linear regression

$$y_i \stackrel{\sim}{=} \sum_j x_{ij} \beta_j + \varepsilon_i \quad \text{constraint}$$

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

$$\sum_j x_{ij} \beta_j \simeq f(x_i) \quad \wedge \quad y_i = f(x_i) + \varepsilon_i$$

in linear regression

$$y_i \in (-\infty, \infty) \wedge x_i \in (-\infty, \infty)$$

Logistic regression & classification problems, we want

$f(x_i)$ to represent discrete outputs, e.g. binary

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

$$f(x_i) \Rightarrow p(x_i) = p_i$$

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

Sigmoid function



if $P(x) = S(x) = \frac{1}{1+e^{-x}}$

if $P(x) > 0.5$, then $\hat{y} = 1$

else $P(x) \leq 0.5$, then $\hat{y} = 0$

$$\int_{x \in D} P(x) dx = 1$$

Our assumption

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

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

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

$$\left\{ p(x) = \frac{e^{\beta_0 + \beta_1 x + \beta_2 x + \dots}}{1 + e^{\beta_0 + \beta_1 x + \dots}} \right\}$$

$$p(x) \rightarrow p(y_i' = 1, x_i' / \beta)$$

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

$$P(y_i = 0, x_i | \beta) = 1 - P_i$$

$$\sum_{i=0,1} P(y_i, x_i) = 1$$

Assume that y_i are independent and identically distributed (i.i.d.)

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

what probability distribution does ϵ follow?

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

$$\varepsilon_i = 1 - p_i$$

$$y_i = 0 \Rightarrow \varepsilon_i = -p_i \quad (p_i \rightarrow p)$$

$$\mathbb{E} [\varepsilon] = \sum_i \varepsilon_i (p_i)$$

$$= (1-p)p + (-p)(1-p)$$

$$\text{var} [\varepsilon] = 0 = (1-p)^2 p + (-p)^2 (1-p) \\ = p(1-p) \Rightarrow$$

Ex Binomial distribution with
mean 0 and variance $p(1-p)$

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

Maximum Likelihood estimation (MLE)

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

Cost function $C(\beta) = P(D|\beta)$

$$\vec{\nabla}_{\beta} C(\beta) = 0$$

Equivalent problem

$$\hat{\beta} = \arg \max_{\beta} \log P(D|\beta)$$

Or

$$\hat{\beta} = \arg \min_{\beta} -\log P(D|\beta)$$

cross entropy

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

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

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

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

$$\frac{\partial \log (1 + e^{\beta_0 + \beta_1 x_i})}{\partial \beta_0} = \frac{e^{\beta_0 + \beta_1 x_i}}{1 + e^{\beta_0 + \beta_1 x_i}} = p_i$$

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

$$\frac{\partial C}{\partial \beta} = -X^T(G - P) = X^T(P - g) = g$$

$$X^T \in \mathbb{R}^{D \times n}$$

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

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

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

$$W = \begin{cases} p_i(1-p_i) & i=j \\ w_{ij} = 0 & i \neq j \end{cases}$$

Linear regression (OLS)

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

$$g = x^T (P - y) = 0$$

↑
no β -dependence

$$P(\beta) = \frac{\exp(\beta_0 + \beta_1 x)}{1 + \exp(\beta_0 + \beta_1 x)}$$

non-linear dependence
on β .

No analytical solution for β

Numerical solutions using iterative
schemes like Newton-Raphson

Taylor expand $C(\beta)$ around
the optimal value $\hat{\beta}, \hat{\beta} - \beta^{(n)}$

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

+ ...

$$\frac{\partial C}{\partial \beta} \Big|_{\beta = \beta^{(n)}} = -\hat{g}^{(n)}$$

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

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

$$C(\beta) \simeq C(\beta^{(n)}) + (g^{(n)})^T b$$

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

$$(f(x) = \text{const} \Leftrightarrow x^T x + \frac{1}{2} x^T A x)$$

$$\frac{\partial f}{\partial x} = 0 \Rightarrow Ax = \beta$$

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

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

$$\hat{\beta} = \beta^{(n+1)} = \beta^{(n)} - [H(\beta^{(n)})]^{-1}$$

Newton-Raphson
iterative scheme $\times g(\beta^{(n)})$

$$\beta^{(m+1)} = \beta^{(m)} - [H(\beta^{(m)})]^{-1} (x^T P \beta^{(m)}) - g$$

Bottleneck of all optimization methods resides in $H(\beta^{(m)})$ and its inverse,

In Neural networks it is common to have at least 10^5 or more parameters \Rightarrow

$$H \in \mathbb{R}^{10^5 \times 10^5}$$

$H \rightarrow$ replace with constant

η = learning rate

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

Constant
Adagrad
RMSprop

ADAM — 10^5 citations

⋮
⋮
⋮
⋮