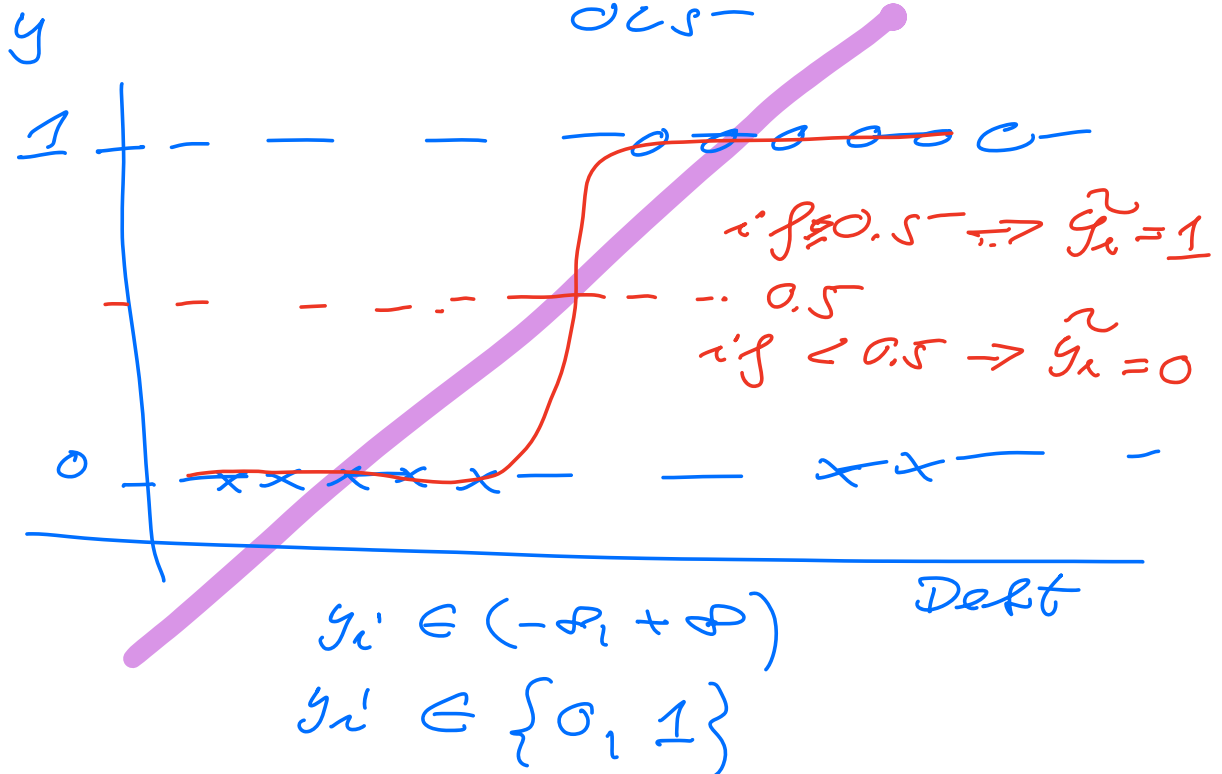
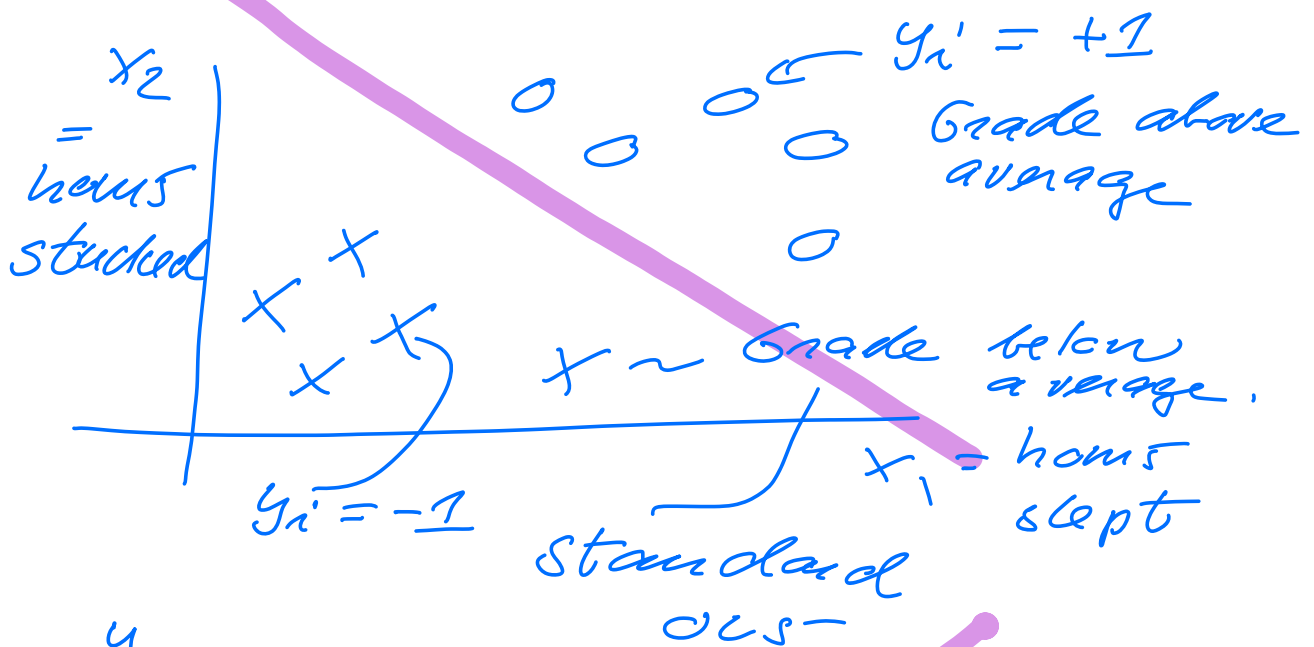


Comp Sci, December 12, 2022

$$y_i = \underset{\substack{\uparrow \\ \text{deterministic}}}{f(x_i)} + \varepsilon_i' \quad \varepsilon_i' \sim N(0, \sigma^2)$$



Replace $f(x_i) \rightarrow p(x_i)$

$$0 \leq p(x_i) \leq 1$$

$$\text{simple } p(x) = \frac{e^x}{1 + e^x}$$

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

Two parameters, β_0, β_1 ,
in principle $x \in (-\infty, +\infty)$

$p(x) \leq 0.5$, then $y = 0$

$p(x) > 0.5$, then $y = 1$

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

$$p(y_i | x_i, \beta) = p(y_i = +1 | x_i, \beta) = p_i$$

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

$$\sum_{i=0}^1 p(y_i | x_i, \beta) = 1$$

logit transformation

$$g(x_i) = \log\left(\frac{p_i}{1-p_i}\right)$$

$$= \beta_0 + \beta_1 x_i'$$

Linear Regression

$$y \approx X\beta + \varepsilon = \hat{y} + \varepsilon$$

$$y \sim N(X\beta, \sigma^2)$$

Replace with

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

which distribution does y follow? which distribution does ε follow?

$$y = 1, \text{ then } p(x) = p$$

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

$y = 0$ has probability $1 - p$

if $y = 1$ has

$\varepsilon = 1 - p$
 with probability p

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

if $y = 0 \Rightarrow \varepsilon = -p$

with probability $1-p$
assume that ε has mean
value $= 0$

$$E[\varepsilon] = (1-p)p - p(1-p) = 0$$

$$E[x] = \sum_{i \in D} p(x_i) x_i$$

$$\begin{aligned} \text{var}[\varepsilon^2] &= (1-p)^2 p \\ &\quad + (-p)^2 (1-p) \\ &= p(1-p) \end{aligned}$$

$\Rightarrow \varepsilon$ follows a Binomial
distribution,
assumption

y_i are i.i.d.,

independent and
identically distributed

$y_i = 1$ has probability $p(x_i)$

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

$y_i = 0$ has probability

$$1 - p(x_i) = 1 - p_i$$

we assume binomial

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

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

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

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

optimal β

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

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

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

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

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

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

$$C(\beta) = - \sum_i \left[y_i' (\beta_0 + \beta_1 x_i') - \log(1 + e^{\beta_0 + \beta_1 x_i'}) \right]$$

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

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

in general we have

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

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

$$X \in \mathbb{R}^{n \times p}$$

Exercise: show that

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

W is a diagonal
only matrix

$$W_{ii} = p_i(1 - p_i)$$

Linear regression:

$$\begin{aligned} \text{OLS} \quad \frac{\partial^2 C}{\partial \beta \partial \beta^T} &= X^T X \cdot \frac{2}{n} \\ \hat{\beta} &= (X^T X)^{-1} X^T y \end{aligned}$$

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

our model

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

Due to non-linearity of β , we don't get a simple analytical expression for β .

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

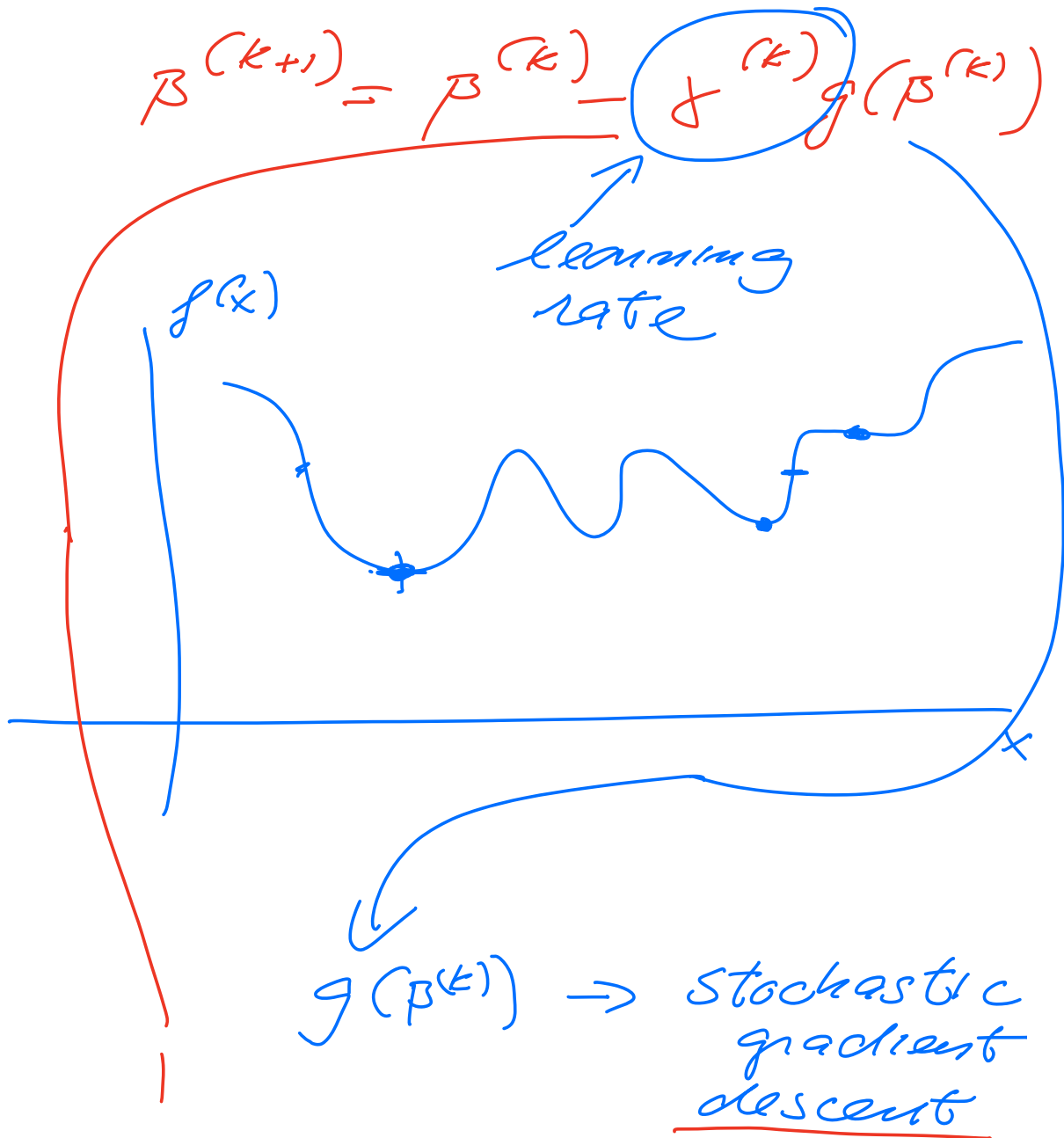
BASIC - method: Newton-Raphson (iterative solution)

$$\beta^{(k+1)} = \beta^{(k)} - \underbrace{H^{-1}(\beta^{(k)})}_{\text{Hessian}} g(\beta^{(k)})$$

$$H = \begin{bmatrix} \frac{\partial^2 \mathcal{L}}{\partial \beta_0 \partial \beta_0} & \frac{\partial^2 \mathcal{L}}{\partial \beta_0 \partial \beta_1} & \dots & \frac{\partial^2 \mathcal{L}}{\partial \beta_0 \partial \beta_{p-1}} \\ \vdots & \vdots & & \vdots \\ \frac{\partial^2 \mathcal{L}}{\partial \beta_{p-1} \partial \beta_0} & \dots & \dots & \dots \end{bmatrix}$$

in "all" algorithms

$$\beta^{(k+1)} = \beta^{(k)} - \underbrace{\alpha^{(k)}}_{\text{learning rate}} g(\beta^{(k)})$$



(efficient
evaluations
of gradients)

↳ momentum
(memory)

GD

momentum
SGD

↳ Learning rate
magic?

- fixed $\eta^{(k)} = \eta$

$\eta \in [10^{-5}, 10^{-9} \dots 10^{-1}]$

- schedulers for $\eta^{(k)}$

$\eta^{(k)} \sim \eta_0 \exp(-k\tau)$

- adaptive learning
rates

- ADAM

- Root Mean square
propagation

(RMSprop)

— ADAGRAD

Taylor expansion of

$C(\beta)$ around $\hat{\beta} - \beta^{(k)}$
 $\hat{\beta} \rightarrow \beta$

$$C(\beta) = C(\beta^{(k)})$$

$$+ \underbrace{g^T(\beta^{(k)})}_{g^T(k)} (\beta - \beta^{(k)})$$

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

$$+ O((\beta - \beta^{(k)})^3)$$

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

approximate to second
derivative

$$C(\beta) \approx C(\beta^{(k)}) + g_{(k)}^T b + \frac{1}{2} b^T H(\beta^{(k)}) b$$

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

$$\beta = \beta^{(k)} - H^{-1}(\beta^{(k)}) g_{(k)}$$

$$f(x) = c + \frac{1}{2} x^T A x + g^T x$$

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

$$x = -A^{-1}g$$