

ML ERASMUS, NOV 14, 2022

---

optimization part

$C(\beta)$  for classification

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

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

$$\begin{aligned} \frac{\partial C}{\partial \beta_1} &= - \sum_{i=0}^{n-1} x_i (y_i - p(x_i)) \\ &= g_1 = 0 \end{aligned}$$

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

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

$$\begin{aligned} \frac{\partial^2 C}{\partial \beta \partial \beta^T} &= X^T W X = H \\ &= \text{Hessian matrix} \end{aligned}$$

$$w_{ii} = p(x_i)(1-p(x_i))$$

How to find optimal  $\hat{\beta}$ ?

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

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

$$\|\beta^{(n+1)} - \beta^{(n)}\|_2 \leq 5 \times 10^{-8} \quad (\text{convergence criterion})$$

$$\begin{aligned} C(\beta) &= C(\hat{\beta}) = C(\beta^{(n+1)}) \\ &= C(\beta^{(n)}) + g^{T(n)}(\beta - \beta^{(n)}) \\ &\quad + \frac{1}{2}(\beta - \beta^{(n)})^T H(\beta - \beta^{(n)}) + \dots \end{aligned}$$

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

$$C(\beta) \approx \underbrace{C(\beta^{(n)})}_{C_0} + g^{T(n)} b^{(n)}$$

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

This is of the form

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

$x$  is unknown

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

$$\Rightarrow Ax = -g \quad (g^{(n)} = g(\beta^{(n)}))$$


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

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

$$x^{(n+1)} = x^{(n)} - A(x^{(n)})^{-1} g^{(n)}$$

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

- need  $g^{(n)}$
- need  $H^{-1}(p^{(n)})$


 replace  $H^{-1}(p^{(n)})$  with  
 a parameter called the  
 learning rate  $\gamma^{(n)}$

Two main issues

① evaluation of gradients

- Gradient descent (GD)
- Stochastic gradient descent (SGD)
- GD/SGD with momentum

② - optimization  $\gamma^{(n)}$

- different scheduler  
(no gradient info)
- adaptive learning  
rates

- Ada grad
- RMS prop
- ADAM

$$\gamma^{(n)} \approx H^{-1}(\beta^{(n)})$$

For Logistic regression  
with  $\beta_0$  and  $\beta_1$

$$H(\beta) = \begin{bmatrix} \frac{\partial g_0}{\partial \beta_0} & \frac{\partial g_0}{\partial \beta_1} \\ \frac{\partial g_1}{\partial \beta_0} & \frac{\partial g_1}{\partial \beta_1} \end{bmatrix}$$

Taylor expand

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

$$C(\beta^{(n)} - \gamma g^{(n)}) \approx$$

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

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

$C(p^{(n)})$  is not a function of  $\gamma$

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

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

$$Hg = \lambda g \quad \text{1. } g^T g = 1$$

$$\gamma = \frac{1}{\lambda} \quad \begin{array}{l} \lambda \text{ is an} \\ \text{eigenvalue} \\ \text{of } H \end{array}$$

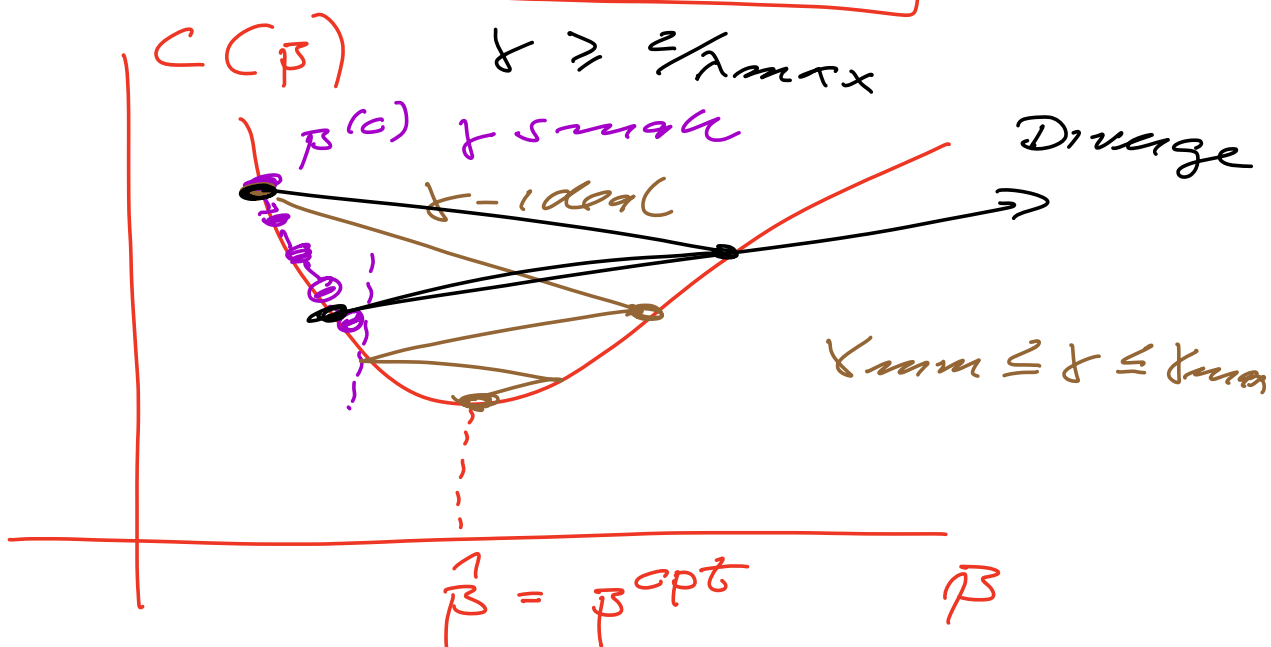
Smallest learning rate

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

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

in general

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



Simple GD

$$\beta^{(n+1)} = \beta^{(n)} - \gamma \underbrace{g^{(n)}(\beta^{(n)})}_{\text{Full Gradient}}$$

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