

Erasmus+ course  
on machine  
learning, November  
12, 2023

Taylor expand cost function  $C(\beta)$  as  
function of  $\hat{\beta} - \beta^{(m)}$

$$\nabla C(\beta) = 0$$

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

$$H(\beta) = \frac{\partial^2 C}{\partial \beta \partial \beta^T}$$

$$g(\beta) = \nabla C(\beta)$$

$$(f(x) = x^T b - \frac{1}{2} x^T A x) \\ \nabla f(x) = 0 \Rightarrow Ax = b$$

Logistic regression

$$g(\beta) = x^T (P - y) \quad 1 + \frac{1}{1 + e^{-x^T \beta}} = x^T w x \\ w_{ii} = p_i(1 - p_i)$$

Expand around

$$\beta^{(n)} - \gamma g^{(n)} \quad | \quad g^{(n)} = g(\beta^{(n)})$$



replaces  $H$

$\gamma$  = learning rate

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

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

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

↗  
2nd derivative

Take derivative w.r.t  $\gamma$ ;  $g^{(n)} \rightarrow g$   
 $H^{(n)} \rightarrow H$

$$\gamma = \frac{g^T g}{g^T H g} \quad (\text{parallel to steepest descent})$$

$$\beta^{(n)} = \frac{g^{T(n)} g^{(n)}}{\underbrace{g^{T(n)} H^{(n)} g^{(n)}}_{\gamma}} \quad g^{(n)}$$

assume  $\gamma$

$$Hg = \lambda g$$

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

The convergence criterion

$$\text{it's } \gamma < \frac{2}{\lambda_{\max}}$$

$\lambda_{\max}$  is largest eigenvalue  
of  $H$

Gradient descent

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

ordinary least squares

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

algorithm of ordinary least squares

- initialize randomly  $\beta^{(0)}$   
=  $[\beta_0^{(0)}, \beta_1^{(0)}, \dots, \beta_{p-1}^{(0)}]$

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

- - - -  $y \in \mathbb{R}^n$

specify iteration step, - after -  
define  $\gamma = 0.01$

for  $i = 1, \text{iter}$

calculate  $g$

Do update

$\beta \leftarrow \beta - \gamma g$

stop if  $|\beta^{(i+1)} - \beta^{(i)}| < \epsilon \sim 10^{-8}$

end for

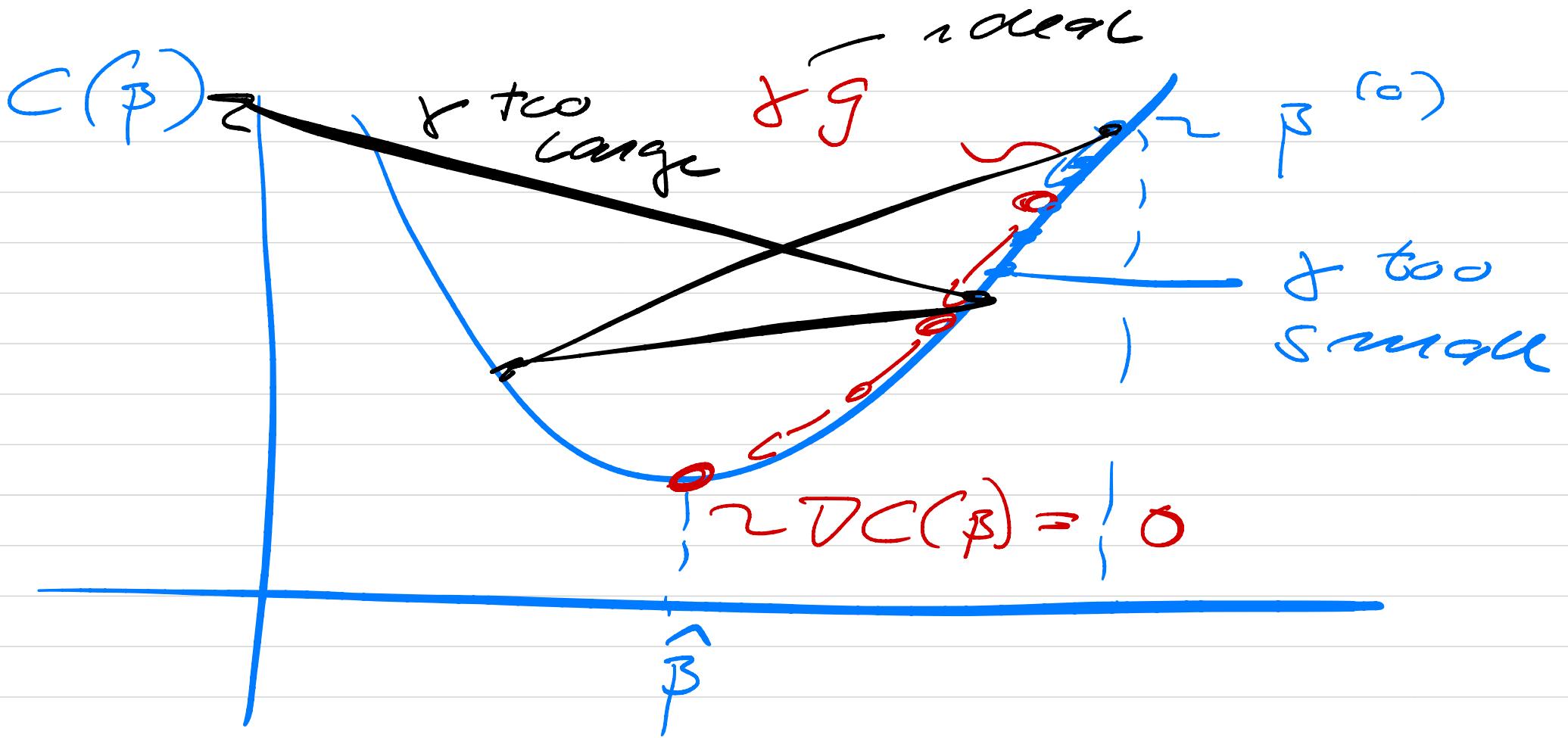
For ordinary least squares

$$H = \mathbf{X}^T \mathbf{X} \quad (\text{no dependence on } \beta)$$

Logistic regression

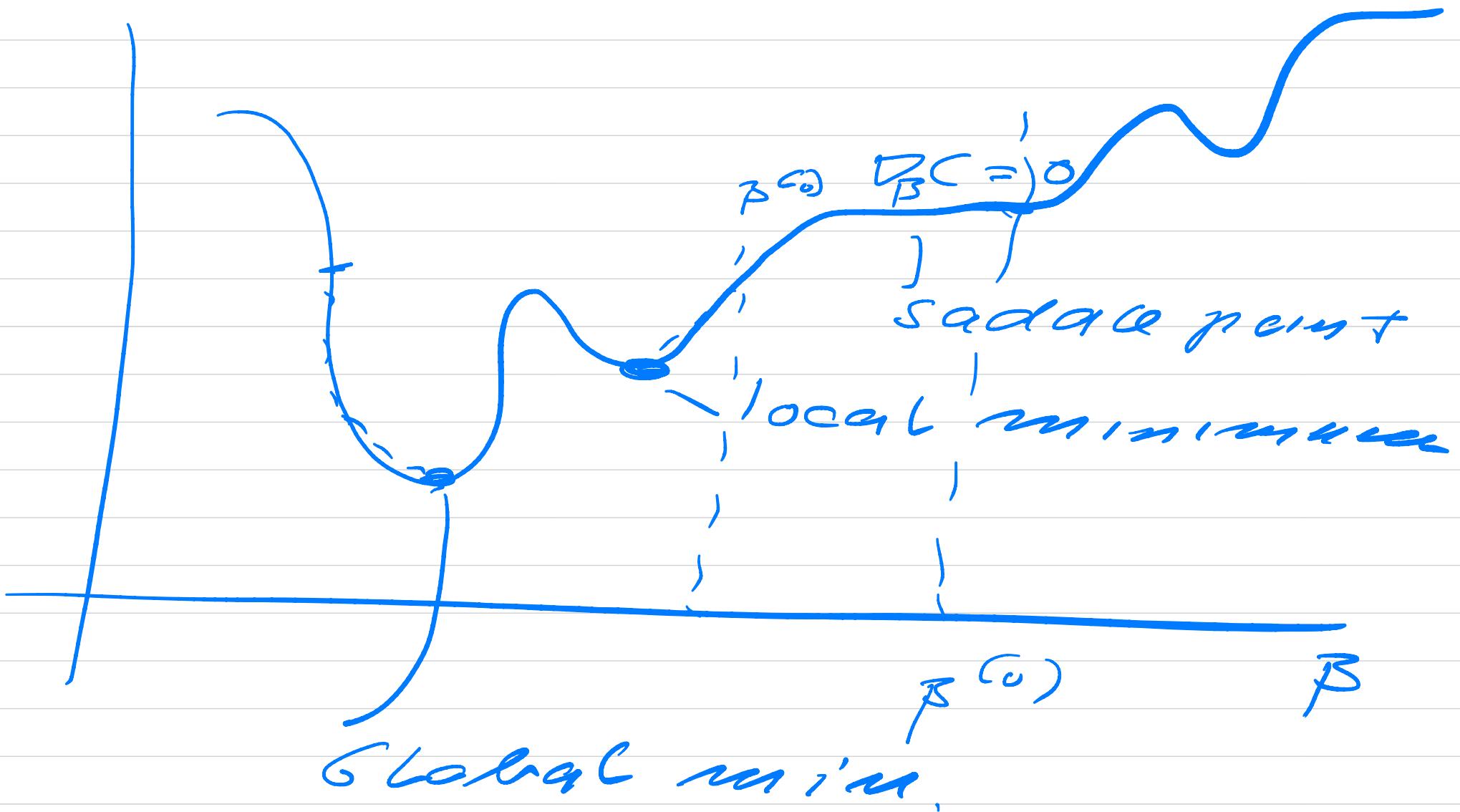
$$H = \mathbf{X}^T \mathbf{W} \mathbf{X}$$

$\nwarrow$  imphat dependence  
on  $\beta$



Goodfellow et al., chapters  
 4 and 8

# General problem



steepest descent

$$f(x) = \frac{1}{2} x^T A x - x^T b$$

↑ known

$$\frac{\partial f(x)}{\partial x} = 0 \quad Ax - b \Rightarrow$$

$$Ax = b$$

solve iteratively

define residual

$$r = b - Ax$$

start with guess  $r_0$  and  $x_0$

exact solution when  $r=0$

$$x_0 = 0, r_0 = +b$$

in general

$$r_{k+1} = b - Ax_{k+1}$$

$$x_{k+1} = x_k + \alpha_k r_k$$

$$r_{k+1} = b - A(x_k + \alpha_k r_k)$$

$$= \frac{(b - Ax_k)}{r_k} - \alpha_k A r_k$$

$$r_{k+1} = r_k - \alpha_k A r_k$$

$$r_{k+1} = 0 \Rightarrow$$

$$r_k = \alpha_k A r_k \quad | \times r_k^T$$

$$\alpha_k = \frac{r_k^T r_k}{r_k^T A r_k} \quad | x_k - \alpha_k g_k$$