

## Lecture October 1

Repeat

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

$$\hat{\beta} - \beta^{(n)} = -g^{(n)} \cdot \eta$$

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

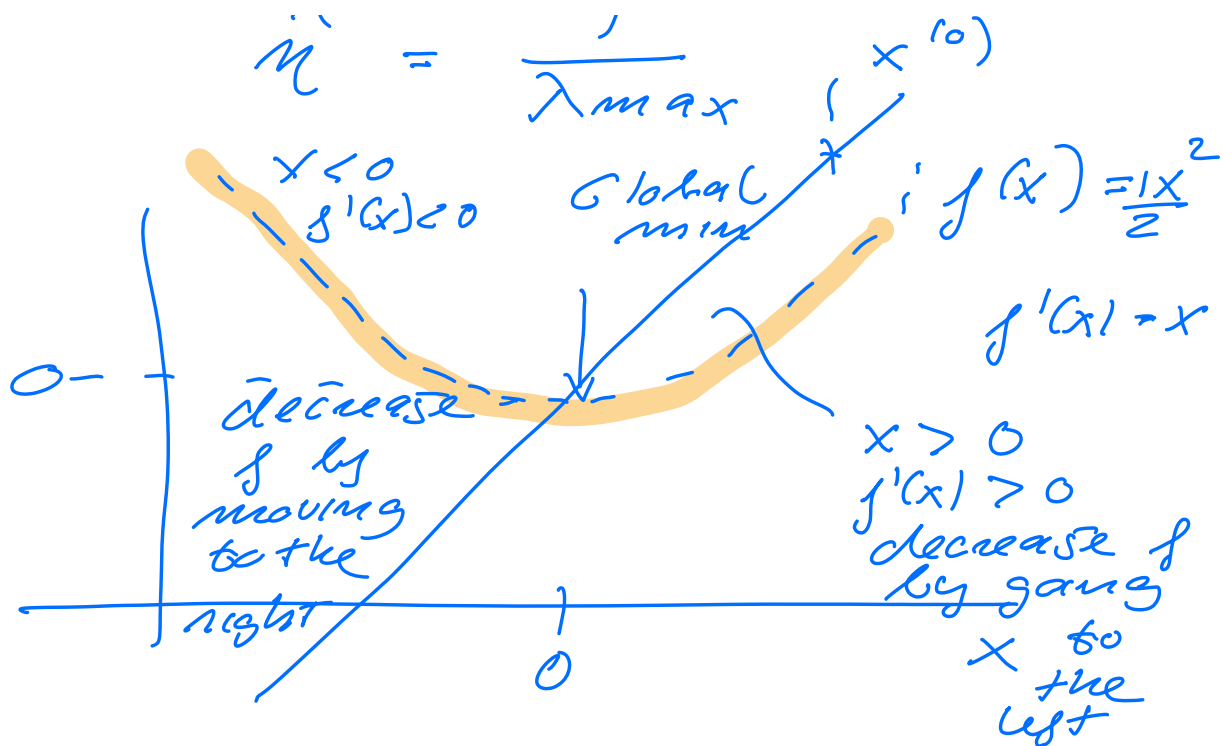
$$\frac{\partial C}{\partial \eta} = 0 \quad \Rightarrow$$

$$\hat{\eta} = \frac{g^T g}{g^T H g}$$

$$Hg = \lambda g$$

^

1



$$f'(x) = 0 \quad \left( = \frac{\partial C(\beta)}{\partial \beta} \right)$$

$$f''(x) > 0 \quad \text{local min}$$

$$f''(x) < 0 \quad \text{local max}$$

$$f''(x) = 0 \quad \text{inconclusive}$$



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

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

$$\left\{ f(x) = \underbrace{\tilde{C}}_{\text{constant}} + x^T \tilde{b} + x^T A x \frac{1}{2} \right\}$$

$$\frac{\partial x^T A x}{\partial x} = (A + A^T) x$$

$$\frac{\partial x^T \tilde{b}}{\partial x} = \tilde{b}$$

$$\frac{\partial C}{\partial b^T} = 0 = 1/b + g \Rightarrow$$

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

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

$$H^{-1} = \eta$$

$$\frac{\partial f(x)}{\partial x^T} = 0 = Ax + \tilde{b}$$

$$Ax = -\tilde{b} = b$$

$$x = A^{-1} \cdot b$$

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

$$\frac{\partial f(x)}{\partial x^T} = 0 \Rightarrow \boxed{Ax = b}$$

Steepest descent

Defines a residual

$$r = b - Ax$$

Start with a guess

$$x^{(0)}$$

$$r^{(0)} = b - Ax^{(0)}$$

... + 1

Define a recipe for the next value of  $x$

$$x^{(1)} = x^{(0)} + \underbrace{\alpha^{(0)}}_{\text{to be determined}} r^{(0)}$$

$$r^{(k+1)} = b - Ax^{(k+1)}$$

$$r^{(k+1)} = b - A(x^{(k)} + \alpha^{(k)} r^{(k)})$$

$$|r^{(k+1)}| \leq \varepsilon \sim 10^{-10}$$

$$b - A(x^{(k)} + \alpha^{(k)} r^{(k)}) = 0$$

$$\underbrace{b - Ax^{(k)}}_{r^{(k)}} - A\alpha^{(k)} r^{(k)} = 0$$

$$(r^{(k)})^T | \quad r^{(k)} = A\alpha^{(k)} r^{(k)}$$

$$\Rightarrow \alpha^{(k)} = \frac{(r^{(k)})^T r^{(k)}}{(r^{(k)})^T A r^{(k)}}$$

$$x^{(k+1)} = x^{(k)} + \alpha^{(k)} r^{(k)}$$

# Stochastic gradient descent (SGD)

$$\beta^{(n+1)} = \beta^{(n)} - \eta^{(n)} \nabla_{\beta} C(\beta^{(n)})$$

$\nabla_{\beta} C(\beta^{(n)})$  in linear regression involves a loop over  $i = 0, 1, \dots, n-1$

$$\nabla_{\beta} C \sim X^T (X\beta - y) = \sum_{i=0}^{n-1} (y_i - \hat{y}_i) X_{i,:}$$

$X \in \mathbb{R}^{n \times p}$   $X_{i,:} \in \mathbb{R}^p$

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

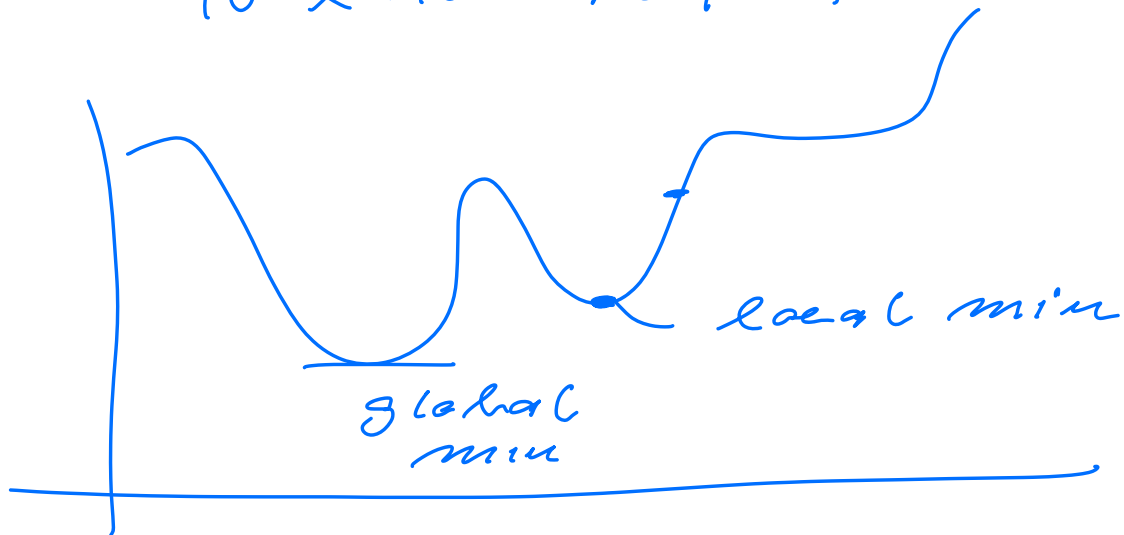
matrix-matrix multiplication

is  $A \in \mathbb{R}^{m \times m} \sim \text{FLOPs} \sim O(m^3)$

$$O(p^2 n)$$

$$n = 10^5 \quad p = 10^3$$

$$10^6 \times 10^5 \text{ Flops,}$$



SGD: select, at random,  
 $n_B < n$  of the training  
 ( $\ll$ ) points.

place these points in  
 a mini-batch  $B_1$

select randomly another  
 $n_B$  points and define  
 mini-batch  $B_2$  ( $n - n_B$ )

Continue till we run  
 out of training points

Gives  $n/n_B$  mini-batches

$B_1, B_2, \dots, B_{(n/n_B)}$

Start with  $B_1$  and compute gradients. And then use this as input to  $B_2$

Continue till you run out of mini-batches = 1 epoch