

# Linear regression

# Lecture 2

$$\underline{X} = \left( \begin{array}{c} \\ \\ \\ \end{array} \right) \quad \begin{array}{c} \uparrow \\ \text{n examples} \\ \downarrow \end{array} \quad \begin{array}{c} \uparrow \\ d \text{ features} \\ \downarrow \end{array}$$

$$x_i = \left( \begin{array}{c} \\ \\ \\ \end{array} \right) \quad \begin{array}{c} \uparrow \\ d \text{ features} \\ \downarrow \end{array} \quad \begin{array}{c} \\ \\ \\ \end{array} \quad \begin{array}{c} \uparrow \\ d \times 1 \end{array}$$

$$y_i = x_i^T \beta + \varepsilon_i$$

$\uparrow$  response       $\uparrow$  noise

$$y_i = \sum_{j=1}^d x_{i,j} \beta_j + \varepsilon_i$$

$y_i \in \mathbb{R}$

$\varepsilon_i$  Gaussian noise

$$x_i = \begin{pmatrix} & \\ & \\ & \end{pmatrix} \quad d \quad d \times 1$$

$$x_i^T = \begin{pmatrix} & \\ & \\ & \end{pmatrix} \quad l \times d$$

$$\beta = \begin{pmatrix} & \\ & \\ & \end{pmatrix} \quad d \times 1$$

$$x_i^T \beta \quad l \times d \quad d \times 1 \quad l \times 1$$

$$\underset{\beta \in \mathbb{R}^d}{\text{Min}} \quad \left\{ \frac{1}{n} \sum_{i=1}^n (y_i - x_i^\top \beta)^2 + \lambda \|\beta\|_2^2 \right\}$$

$n$  examples

$d$  features

$$\underline{\lambda = 0} \quad \underset{\beta \in \mathbb{R}^d}{\text{Min}} \quad \frac{1}{n} \sum_{i=1}^n (y_i - x_i^\top \beta)^2$$

least-squares

$$n \gg d$$

$$\beta \in \mathbb{R}^d \quad \beta = \begin{pmatrix} \beta_1 \\ \vdots \\ \beta_d \end{pmatrix}$$

n

$$\min_{\beta \in \mathbb{R}^d} \frac{1}{n} \sum_{i=1}^n (y_i - x_i^\top \beta)^2 + \lambda \sum_{j=1}^d \beta_j^2$$

Constrained Least-squares

$$\left. \begin{array}{l} \text{Min} \\ \beta \in \mathbb{R}^d \end{array} \right\} \frac{1}{n} \sum_{i=1}^n (y_i - x_i^\top \beta)^2$$

$$\text{subject to } 0 \leq \beta_1 \leq \beta_{\max}$$

$$\vdots$$
$$0 \leq \beta_d \leq \beta_{\max}$$

$$\left. \begin{array}{l} \text{Min} \\ \beta \in \mathbb{R}^d \end{array} \right\} f(\beta)$$

(Mathematical) Optimization

$\beta \in \mathbb{R}^d$

$\beta \in \mathcal{B}$

$\mathcal{P}^d$

(grid-search  
brute-force)

Min  $f(\beta)$   
 $\beta \in \mathbb{R}^d$

$f$  differentiable

→ Gradient Descent

Algo .  $\beta = 0$

. Iterate for  $t = 1, \dots, T$

$$\beta^{(t+1)} = \beta^{(t)} - \gamma_t \nabla f(\beta^{(t)})$$

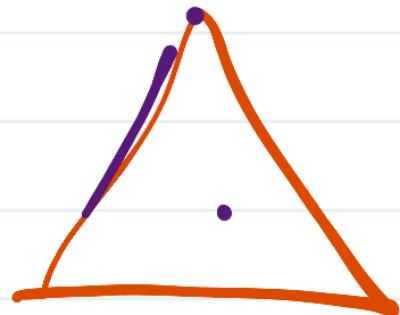
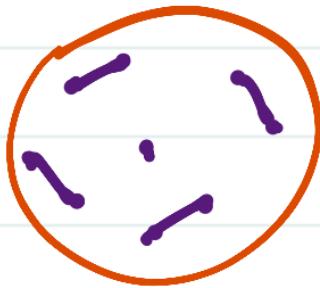
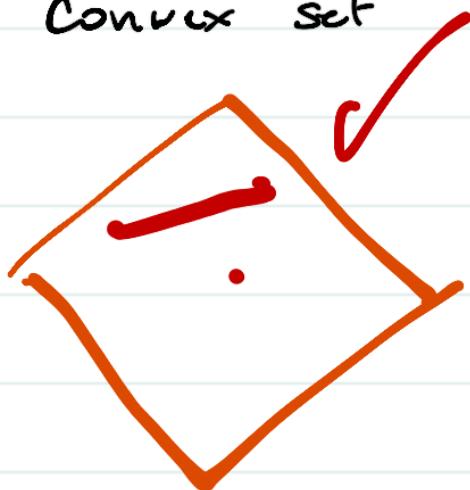
Return  $\beta^{(T)}$

$T$  ?  
 $\beta^{(T)}$  ?  
 $\gamma_t$  ?

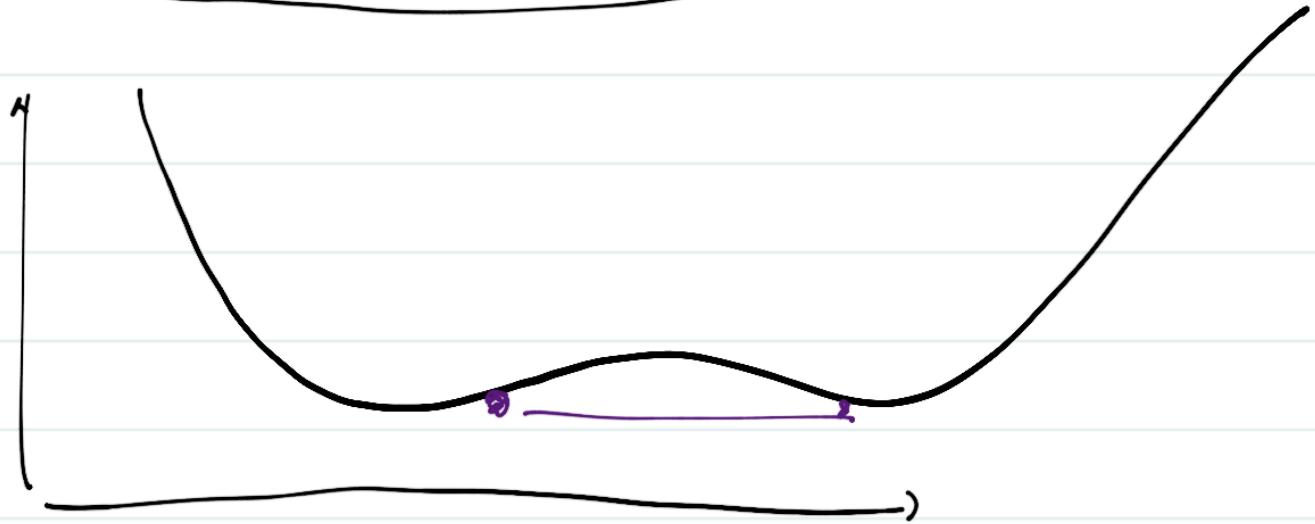
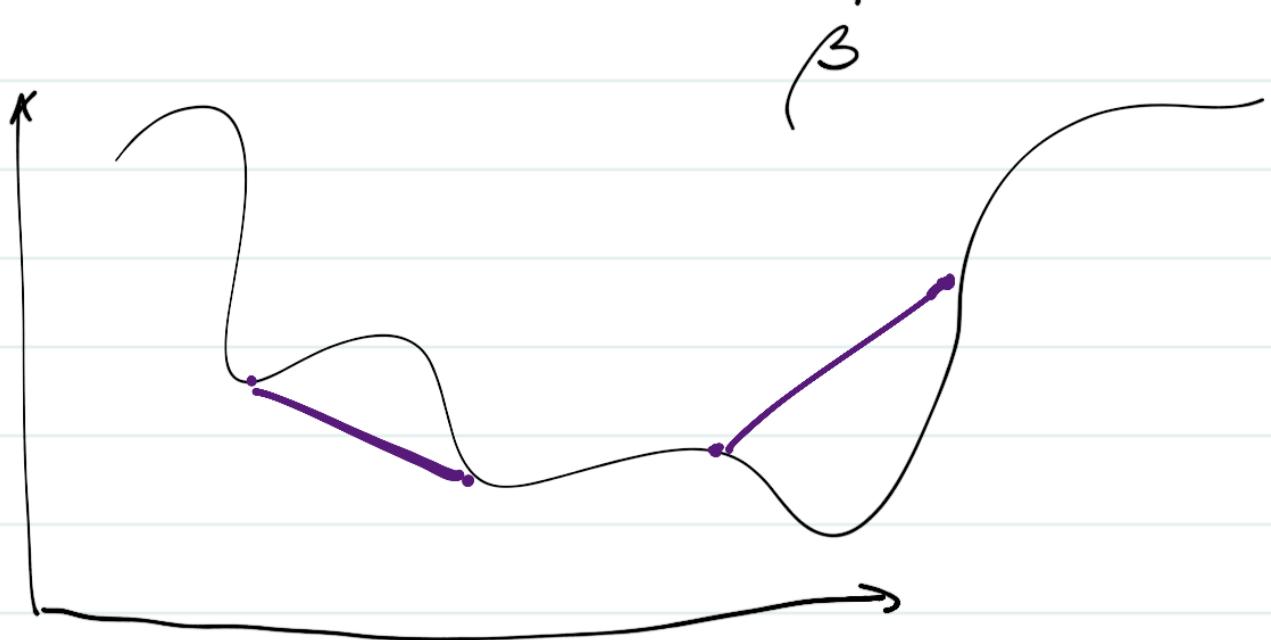
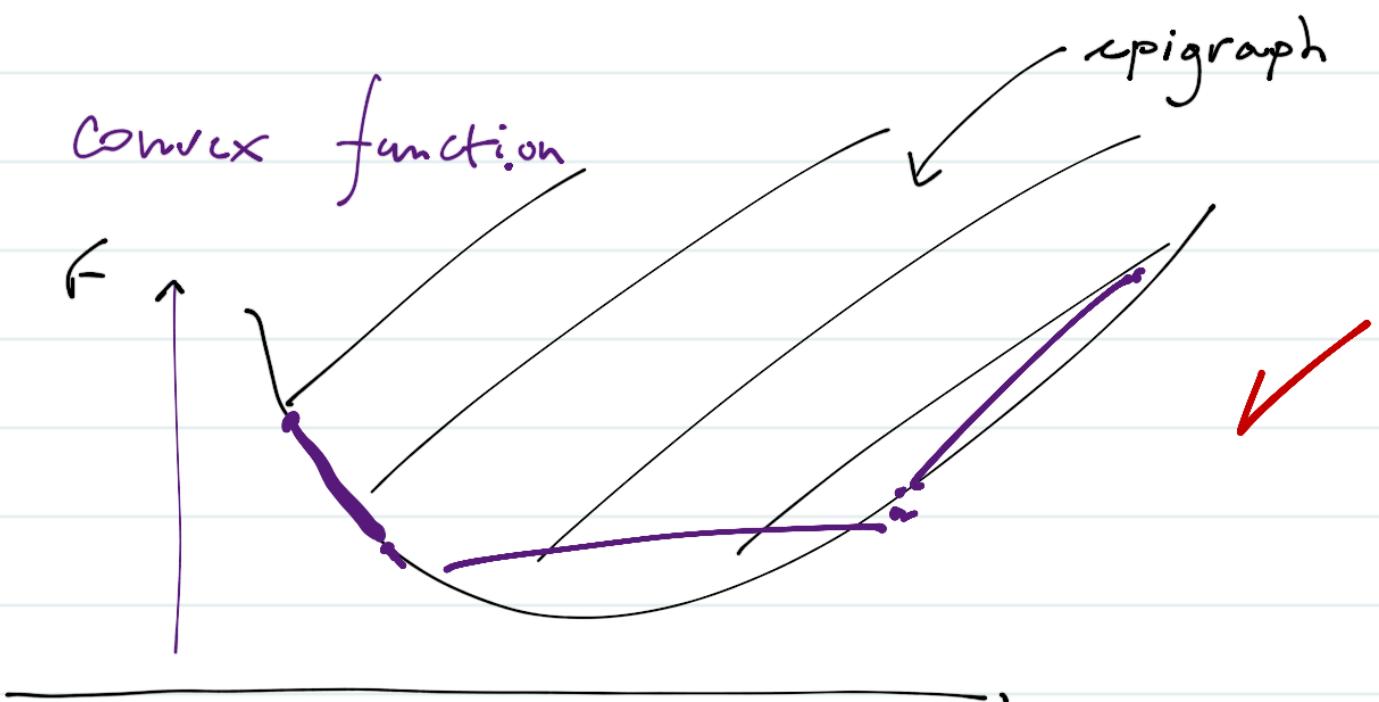
maximum number of iterations  
 what can you say about it?  
 how to set the step-size

$F$  convex

convex set



Convex function

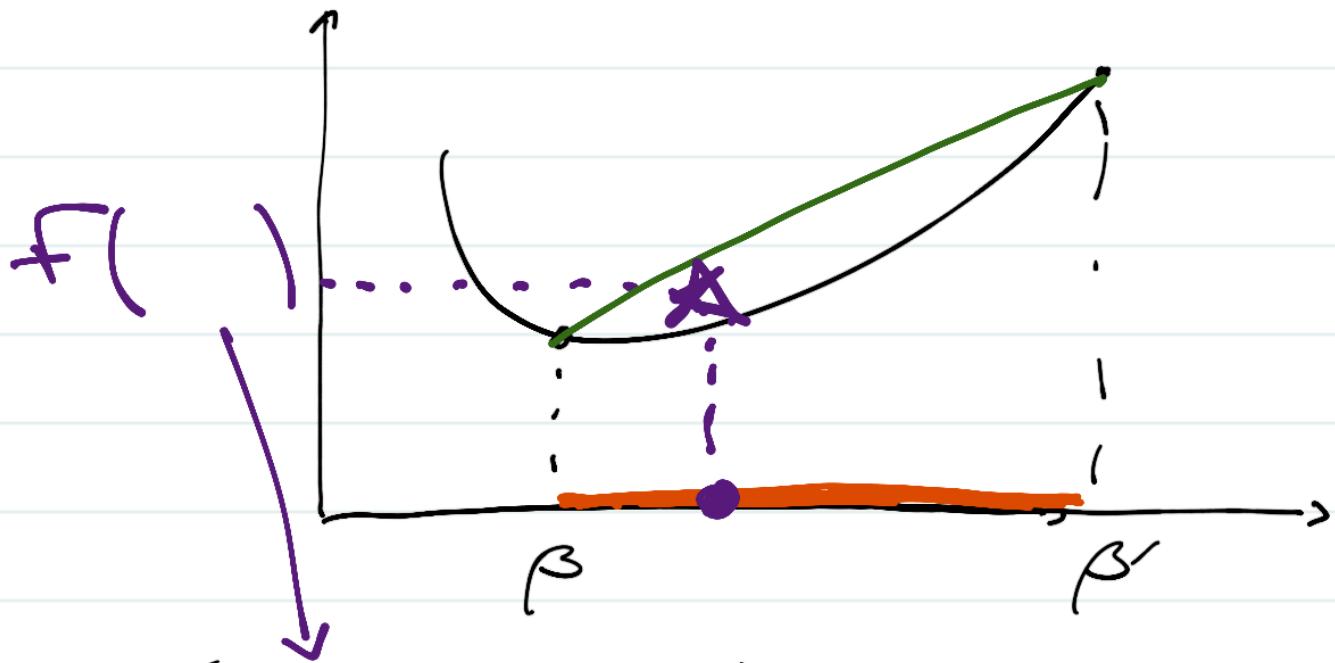


$$f : \mathbb{R}^d \rightarrow \mathbb{R}$$

$$\beta \mapsto f(\beta)$$

For any  $\beta, \beta' \in \mathbb{R}^d$

For any  $\alpha \in [0, 1]$



$$f\left(\boxed{\alpha\beta + (1-\alpha)\beta'}\right)$$

$$\leq \alpha f(\beta) + (1-\alpha) f(\beta')$$

$$\beta = \begin{pmatrix} \cdot \\ \cdot \\ \vdots \\ \cdot \end{pmatrix} \quad \text{or} \quad \alpha\beta = \begin{pmatrix} \alpha x_{1\dots} \\ \cdot \\ \vdots \\ \cdot x_{n\dots} \end{pmatrix} \quad (1)$$

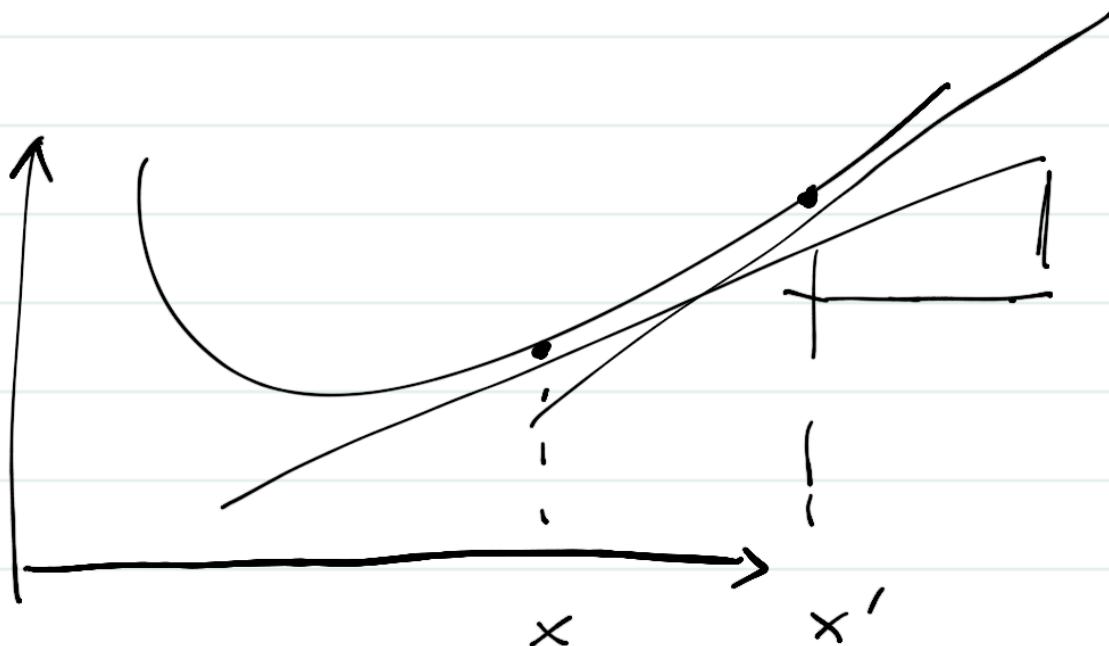
$$\beta' = \begin{pmatrix} \cdot \\ \cdot \\ \vdots \\ \cdot \end{pmatrix} \quad \text{or} \quad (1-\alpha)\beta' = \begin{pmatrix} (1-\alpha)x_{1\dots} \\ \cdot \\ \vdots \\ \cdot (1-\alpha)x_{n\dots} \end{pmatrix} \quad (2)$$

$$\alpha\beta + (1-\alpha)\beta' = (1) + (2)$$

# Smooth functions

A function is  $L$ -smooth with constant  $L$  if  $F$  is continuous, differentiable and for all  $x, x'$

$$\|\nabla F(x) - \nabla F(x')\| \leq L \|x - x'\|$$



# Convergence of gradient descent

Assume that  $F$  is :

- convex

- $L$ -smooth

eta

Then if  $\eta_t = \frac{1}{L}$ , we have

$$F(\beta^{(t)}) - F^* \leq \frac{2L\|\beta^{(0)} - \beta^*\|^2}{t}$$

$$\cdot F^* = \underset{\beta}{\text{Min}} \quad F(\beta)$$

$$\cdot \beta^* = \underset{\beta}{\text{Arg Min}} \quad F(\beta)$$

$$\cdot F(\beta^{(t)}) - F^* = O\left(\frac{1}{t}\right)$$

$$F(\beta^{(t)}) - F^* = \varepsilon$$

$$\frac{2L\|\beta^{(0)} - \beta^*\|}{\varepsilon} \leq \varepsilon \quad \text{---}$$

$$t \geq \frac{2L\|\beta^{(0)} - \beta^*\|}{\varepsilon}$$

$$t \geq \frac{2LD}{\varepsilon}$$

$$T = \left\lceil \frac{2LD}{\varepsilon} \right\rceil + 1$$

$$\underset{\beta \in \mathbb{R}^d}{\text{Min}} \quad F(\beta) := \frac{1}{n} \sum_{i=1}^n (y_i - x_i^\top \beta)^2 + \lambda \|\beta\|_2^2$$

$$\|\beta\|_2^2 = \sum_{j=1}^d \beta_j^2$$

Algo . Gradient Descent

$$\begin{array}{ll}\text{Input:} & L \\ \text{Init:} & \beta^{(0)}\end{array}$$

Iterate: for  $t = 1, \dots, T$

$$\beta^{(t+1)} = \beta^{(t)} - \frac{1}{L} \nabla_{\beta} F(\beta^{(t)})$$

$$\text{Return } \beta^{(T+1)}$$

## Simplest case

$$n=1 \quad d=1 \quad \beta \in \mathbb{R}^1$$

$$F(\beta) = \frac{1}{1} \sum_{i=1}^1 (y_i - x_i \beta)^2 + \lambda \sum_{j=1}^1 \beta_j^2$$

$$F(\beta) = (y - x\beta)^2 + \lambda \beta^2$$

•  $F$  convex ✓ sum of two quadratics  
 $f'' \geq 0$

•  $F$  smooth ✓

$$F'(\beta) = \cdot$$

Chain rule

$$f(u) = g(h(u))$$

$$f'(u) = h'(u) \cdot g'(h(u))$$

$$f(u) = au^2 \quad f'(u) = 2au$$

$$F(\beta) = 2 \cdot (-x) \cdot (y - x\beta) + 2\lambda\beta$$

$$\underline{n \geq 1} \quad \underline{d \geq 1}$$

$$F(\beta) = \frac{1}{n} \sum \dots \quad \beta = \begin{pmatrix} \beta_1 \\ \vdots \\ \beta_d \end{pmatrix}$$

$$\bar{X} \quad \beta \quad \bar{Y}$$

$$\bar{X} = \left( \dots \begin{array}{|c|} \hline x_i \\ \hline \dots \end{array} \dots \right) \quad \bar{Y} = \begin{pmatrix} y_1 \\ \vdots \\ y_n \end{pmatrix}$$

$n$

$$\sum_{i=1}^n (y_i - x_i^\top \beta)^2$$

$$\tilde{X}^T \beta = \begin{pmatrix} x_1^T \\ \vdots \\ x_n^T \end{pmatrix} \begin{pmatrix} \beta_1 \\ \vdots \\ \beta_d \end{pmatrix}$$

$$\begin{matrix} n \times d & & d \times 1 \\ \curvearrowright & & \\ & n \times 1 & \end{matrix}$$

$$\underline{\underline{X}}^T \beta$$

$$\bar{Y} = \begin{pmatrix} y_1 \\ \vdots \\ y_n \end{pmatrix}$$

$$\bar{Y} - \bar{X}\beta$$

$$= \begin{pmatrix} y_1 \\ \vdots \\ y_n \end{pmatrix} - \underbrace{\bar{X}^T \beta}_{n \times 1}$$

$n \times 1$

$n \times 1$

(OK)

$$\| \bar{Y} - \bar{X}\beta \|_2^2$$

$$\| v \|_p = \left( \sum_{i=1}^n |v_i|^p \right)^{\frac{1}{p}}$$

$p \geq 1$

$$\begin{aligned} L_1\text{-norm} \\ \|\mathbf{v}\|_1 &= \left( \sum_{i=1}^n |v_i|^1 \right)^{\frac{1}{1}} \\ &= \sum_{i=1}^n |v_i| \end{aligned}$$

$$\begin{aligned} L_2\text{-norm} \\ \|\mathbf{v}\|_2 &= \left( \sum_{i=1}^n |v_i|^2 \right)^{\frac{1}{2}} \\ &= \left( \sum_{i=1}^n v_i^2 \right)^{\frac{1}{2}} \end{aligned}$$

$$\|\mathbf{v}\|_2^2 = \sum_{i=1}^n v_i^2$$

$$\frac{1}{n} \sum_{i=1}^n (y_i - x_i^\top \beta)^2$$

$$= \frac{1}{n} \| Y - X^\top \beta \|_2^2$$

$$F(\beta) = \underbrace{\frac{1}{n} \| Y - X^\top \beta \|_2^2}_g + \lambda \underbrace{\| \beta \|_2^2}_h$$

$\frac{\partial h}{\partial \beta} = 2 \lambda \beta \quad \dots$

useful general identity

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

$$A \leftarrow I_d \quad ; \quad x \leftarrow \beta$$

$$g(\beta) = \frac{1}{n} \| y - X\beta \|_2^2$$

$$g(\beta) = \frac{1}{n} \left\{ (y - X\beta)^T (y - X\beta) \right\}$$

$$= \frac{1}{n} \left\{ Y^T Y - Y^T X^T \beta - \beta^T X Y \right.$$

$$\left. + \beta^T X X^T \beta \right\}$$

$$(X^T \beta)^T = \beta^T (X^T)^T = \beta^T X$$

$$\frac{\partial g(\beta)}{\partial \beta}$$

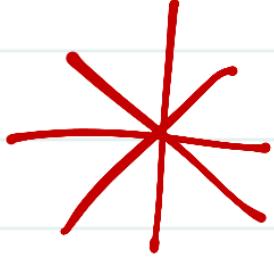
$$= \frac{1}{n} \sum_{i=1}^n - \frac{\partial}{\partial \beta} \left( (x_i^T y_i)^\top \beta \right)^T * \\ - \frac{\partial}{\partial \beta} \left( \beta^T (x_i^T y_i) \right) * \\ + \frac{\partial}{\partial \beta} \left( \beta^T (x x^T) \beta \right)$$

Useful formula

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

$$\frac{\partial}{\partial \beta} \left( ((x^y)^T \beta) \right) =$$

$X^Y$



$$\frac{\partial}{\partial \beta} \left( \beta^T (x^y) \right) = x^y .$$

$$\frac{\partial}{\partial \beta} \left( \beta^T [BLAH] \beta \right)$$

$$(AB)^T = B^T A^T$$

$$= (BLAH + BLAH^T) \beta$$

$$\frac{\partial}{\partial \beta} \left( \beta^T (\underbrace{xx^T}_{}) \beta \right) .$$

$$\begin{aligned} (xx^T)^T &= (x^T)^T x^T \\ &= xx^T \end{aligned}$$

$$= (xx^T + (xx^T)^T) \beta$$

$$= (xx^T + xx^T) \beta = 2xx^T \beta$$

Collecting everything

$$\frac{\partial g}{\partial \beta} = \frac{1}{n} \left\{ -2XY + 2XX^T\beta + 0 \right\}$$
$$= -\frac{2}{n} X(Y - X^T\beta)$$

$$\frac{\partial F}{\partial \beta} = \frac{\partial g}{\partial \beta} + \frac{\partial h}{\partial \beta}$$
$$= \boxed{-\frac{2}{n} X(Y - X^T\beta) + 2\lambda\beta}$$

Algo    Init     $\beta^{(0)} = 0_{IR^d}$   
Iterate    for  $t = 1, \dots, T$

$$\beta^{(t+1)} = \beta^{(t)} - \frac{1}{L} \boxed{\nabla_{\beta} F(\beta^{(t)})}$$

Return :  $\beta^{(T+1)}$

.  $\angle$  constant

$$\boxed{h = 1 \quad d = 1}$$

$$F(\beta) = (y - x\beta)^2 + 2\lambda\beta^2$$

$$|\nabla F(\beta) - \nabla F(\beta')| \leq L |\beta - \beta'|$$

$$\nabla F(\beta) = -2x(y - x\beta) + 2\lambda\beta$$

$$\nabla F(\beta) - \nabla F(\beta')$$

$$= -2x(y - x\beta) + 2\lambda\beta$$

$$+ 2x(y - x\beta') - 2\lambda\beta'$$

$$= \cancel{-2xy} + 2x^2\beta + 2\lambda\beta$$
  
~~+ 2xy~~  $- 2x^2\beta' - 2\lambda\beta'$

$$\begin{aligned}
 & \nabla F(\beta) - \nabla F(\beta') \\
 = & 2x^2 (\beta - \beta') + 2\gamma (\beta - \beta') \\
 = & 2(x^2 + \gamma)(\beta - \beta')
 \end{aligned}$$

$$|\nabla F(\beta) - \nabla F(\beta')| \leq L |\beta - \beta'|$$

$$L = 2(x^2 + \gamma)$$