

# Least Square Gradient Descent

- first lets compute gradient wrt single param  $\theta_i$  & single observation  $(x, y)$  & then vectorize so<sup>l</sup> to update param simultaneously
- one obs at a time? → iterative / online gradient learning.

for single obs, error  $\Rightarrow J = (y - x\theta)^2$

$$\Rightarrow \frac{d}{d\theta_i} (y^2 - 2yx\theta + (x\theta)^2)$$

$$= -2yx_i + 2x\theta \cdot x_i$$

$$\frac{d}{d\theta_i} = 2(x\theta - y)x_i$$

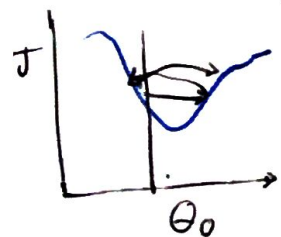
$$\frac{dJ}{d\theta} = \begin{bmatrix} 2(x\theta - y)x_0 \\ 2(x\theta - y)x_1 \\ \vdots \\ 2(x\theta - y)x_n \end{bmatrix} = x^T$$

$$\frac{dJ}{d\theta} = 2x^T(x\theta - y)$$

$$\theta = \theta - 2x^T(x\theta - y)$$

Learning Rate:

$$\theta = \theta - \frac{dJ}{d\theta}$$



if we move params according to  $\frac{dJ}{d\theta}$

we can go too slowly → too much time

go too fast → overshoot desired minimum/maxim

(basically step size)

∴ we use  $\eta$  called learning rate that controls how much to go in direction of gradient.

(maxing step-size dynamic sorta)

→ standardize the data or  $\eta$  goes crazy!!