

$$(x_i, y_i)_{i=1}^n$$

$$x_i \in \mathbb{R}^m$$

$x_{ij}$   $j^{\text{th}}$  feature of  $i^{\text{th}}$  point.

$$\hat{y}_i = f_w(x_i) = w^T x_i = \sum_{j=1}^m w_j x_{ij}$$

$$r_i = y_i - \hat{y}_i$$

$$l(w) = \sum_{i=1}^n (r_i)^2 = \sum_{i=1}^n (y_i - f_w(x_i))^2$$

$$= \sum_{i=1}^n \left( y_i - \sum_{j=1}^m w_j x_{ij} \right)^2$$

$$w^* \in \arg \min_w l(w)$$

sequence of model parameters  
 $w_i \in \mathbb{R}^m$

Hill Climbing.



$w_0$

$w_1$

$w_2$

$\vdots$

$\vdots$

$w_{i,j}$

which term in sequence

which weight of  $w_i$ .

"Directions"



$$w_1 = (w_{0,1} + \Delta_1, w_{0,2} + \Delta_2)$$

$$w_0 = (w_{0,1}, w_{0,2})$$

$$\Delta = (\Delta_1, \Delta_2)$$

$$w_1 = w_0 + \Delta$$



Lecture

- Gradient descent
- Basis functions
- Feature normalization.

Dir. Steepest Ascent of  $l$  at  $w_i$  is gradient:

$$\nabla l(w_i) = \left( \frac{\partial l(w_i)}{\partial w_{i,1}}, \frac{\partial l(w_i)}{\partial w_{i,2}}, \dots, \frac{\partial l(w_i)}{\partial w_{i,m}} \right)$$

$w_{i,m}$  random.  
often zero.

-  $\nabla l(w_i)$  is the dir. of steepest descent.

## Gradient Descent

$$x_i \in \mathbb{R}^m$$
$$f_{w_i}(x) = \sum_{j=1}^m w_{ij} x_j$$

Input: loss function  $l: \mathbb{R}^m \rightarrow \mathbb{R}$   
dimension  $m$  of  $w$

Output: An approximation of an element of  $\arg \min_w l(w) \Rightarrow \sum_{i=1}^n (y_i - \hat{y}_i)^2 \rightarrow f_w(x_i)$

Hyperparams:  
 $w_0$ : initial point  
 $\alpha$ : stepsize.  
Stopping criteria.

$i \leftarrow 0$

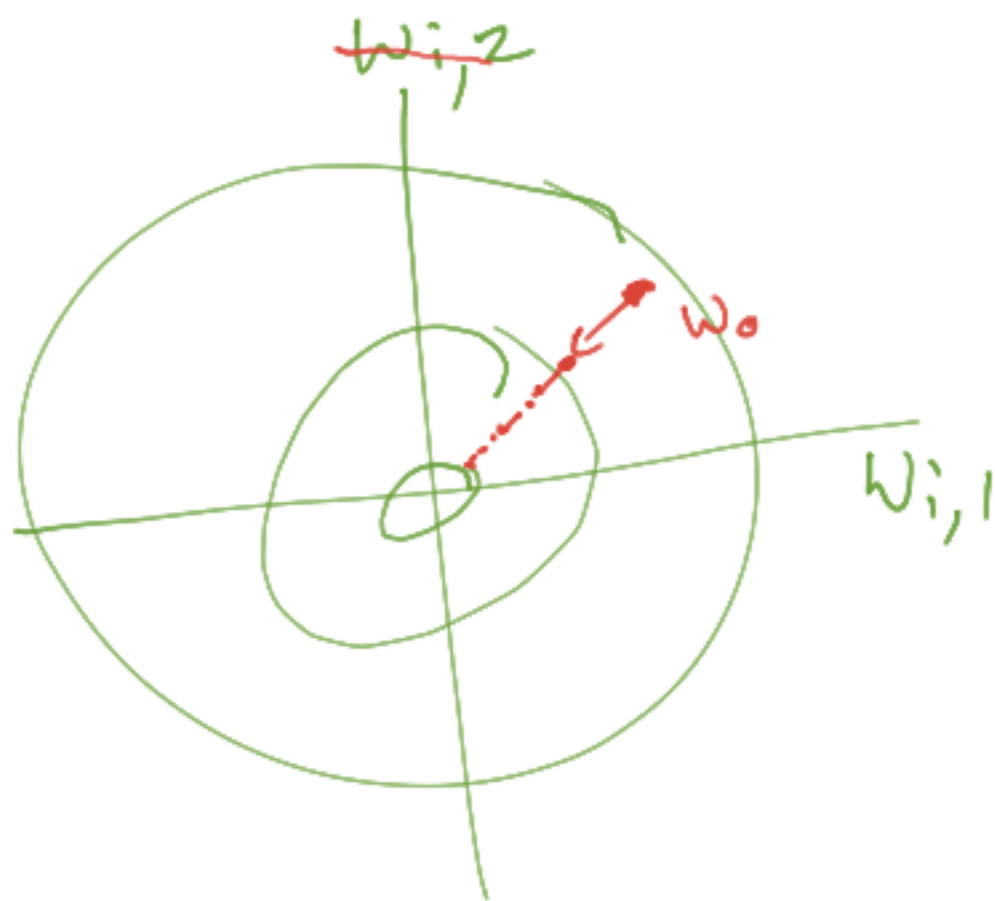
while stopping criteria not satisfied do

for  $j = 1$  to  $m$  do

$$w_{i+1,j} = w_{i,j} - \alpha \frac{\partial l(w_i)}{\partial w_{i,j}}$$

$i = i + 1$

return  $w_i$



→ After fixed time.

→ When  $l(w_i)$  sufficiently small

→ When  $\nabla l(w_i)$  is small,

→ When  $l(w_i)$  not much less than  $l(w_{i-k})$

$$l(w) = w^2$$

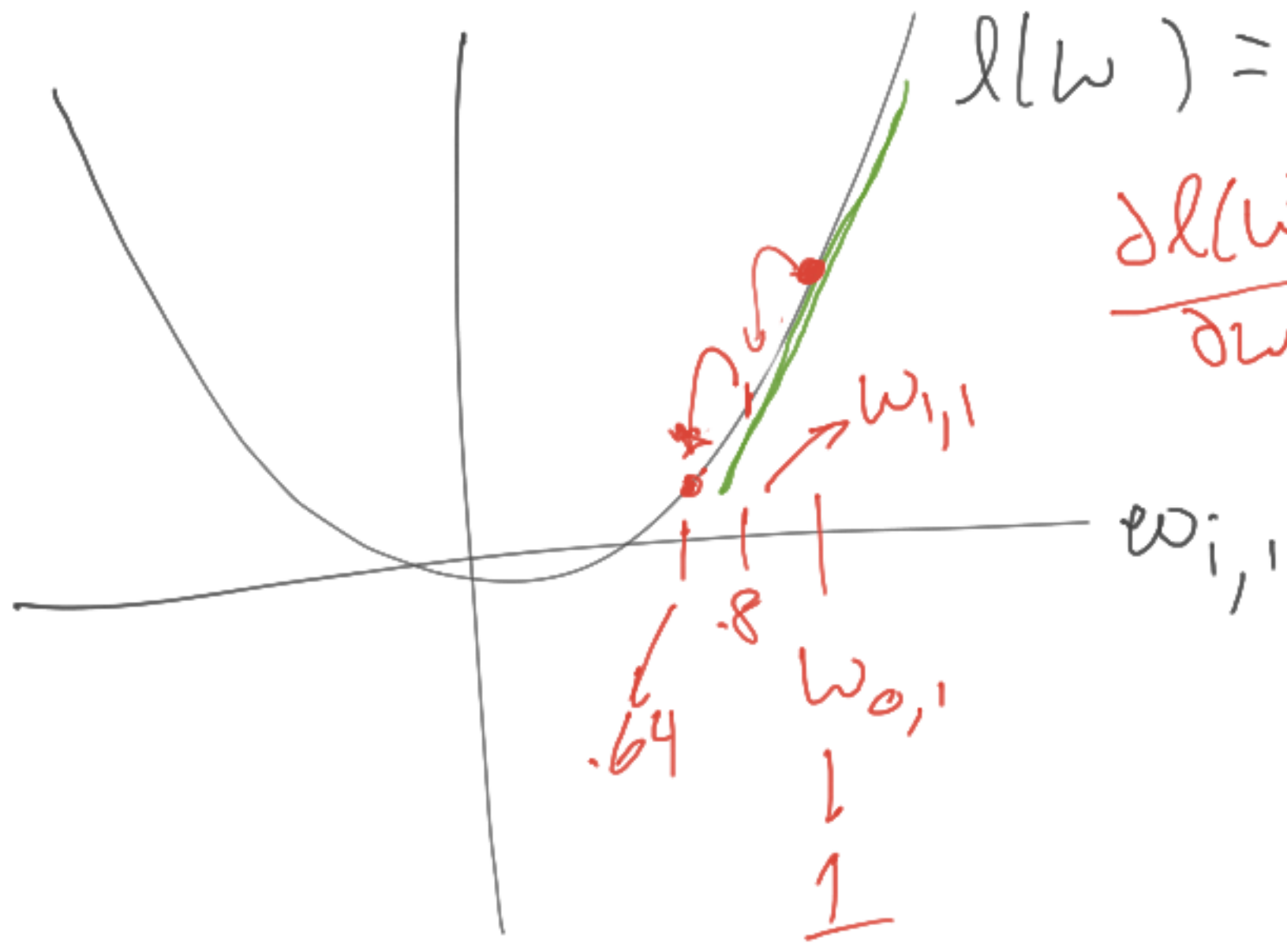
$$w \in \mathbb{R} \quad .8 - .1(1.6)$$

$$\frac{\partial l(w)}{\partial w} = 2w$$

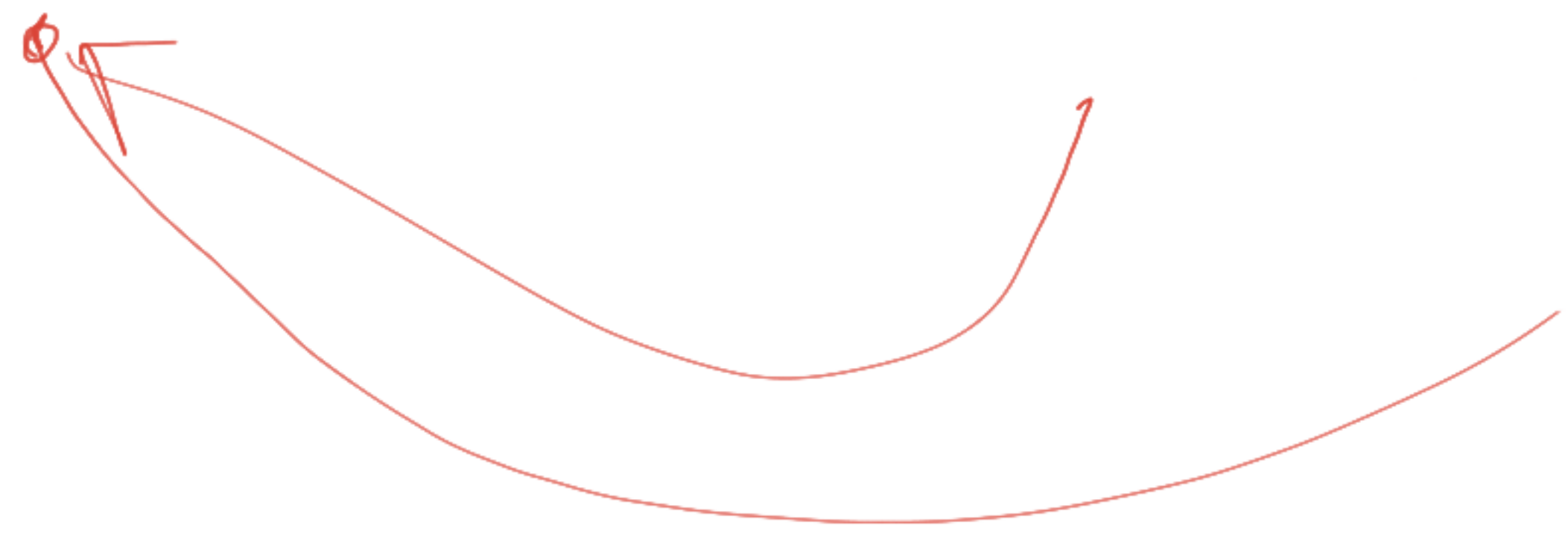
$$.8 - .1(2 \cdot .8)$$

$$w_{i+1,1} = w_{i,1} - \alpha(2w_{i,1})$$

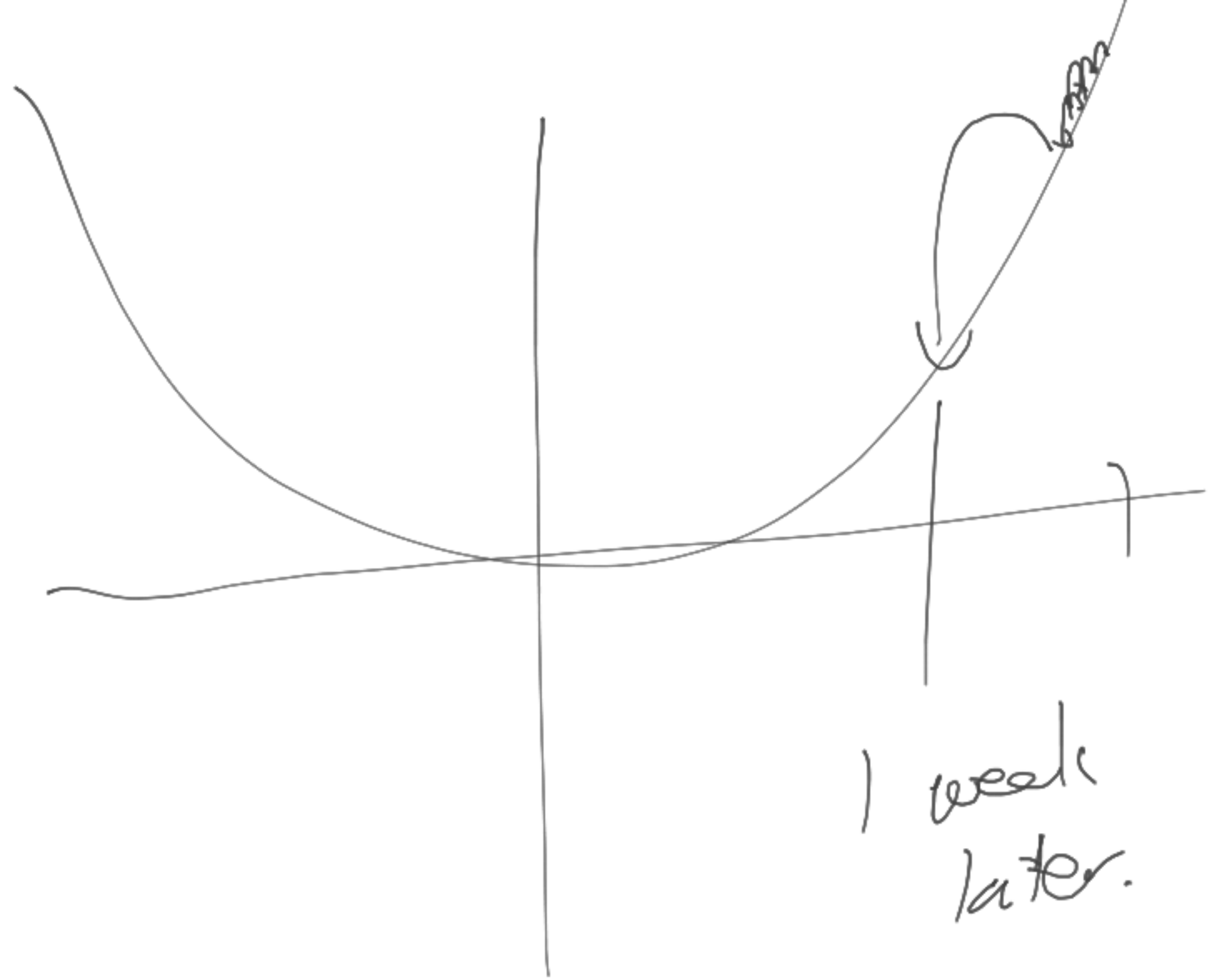
$$\alpha = 0.1$$



-1.6



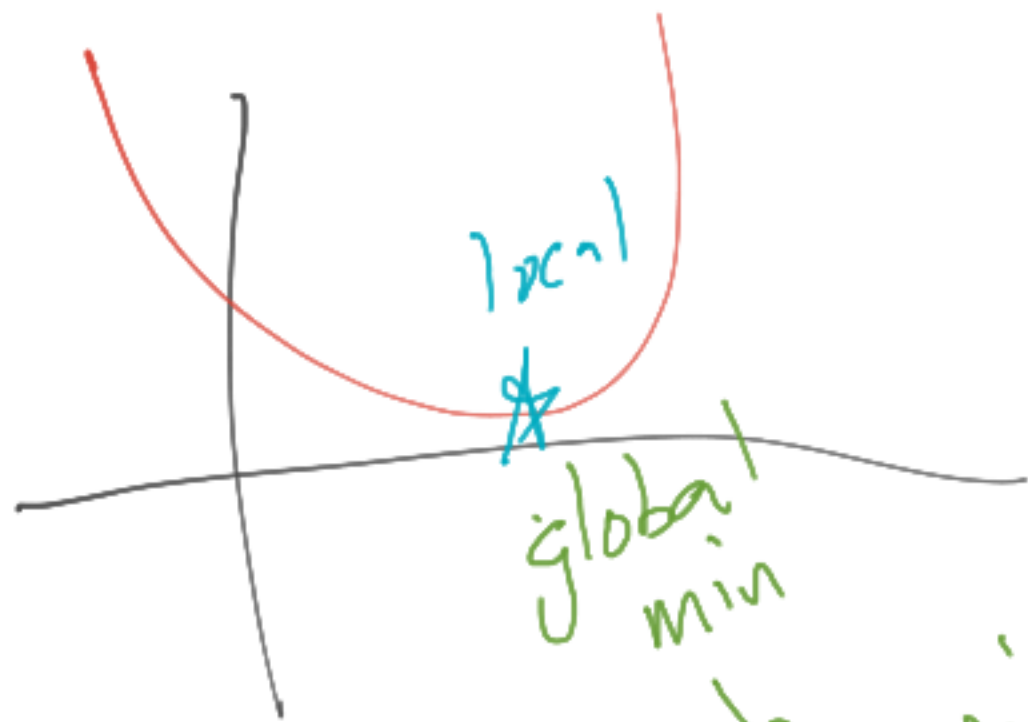
"divergence"



1 week later.

Gradient descent converges to a local minimum (given certain conditions).

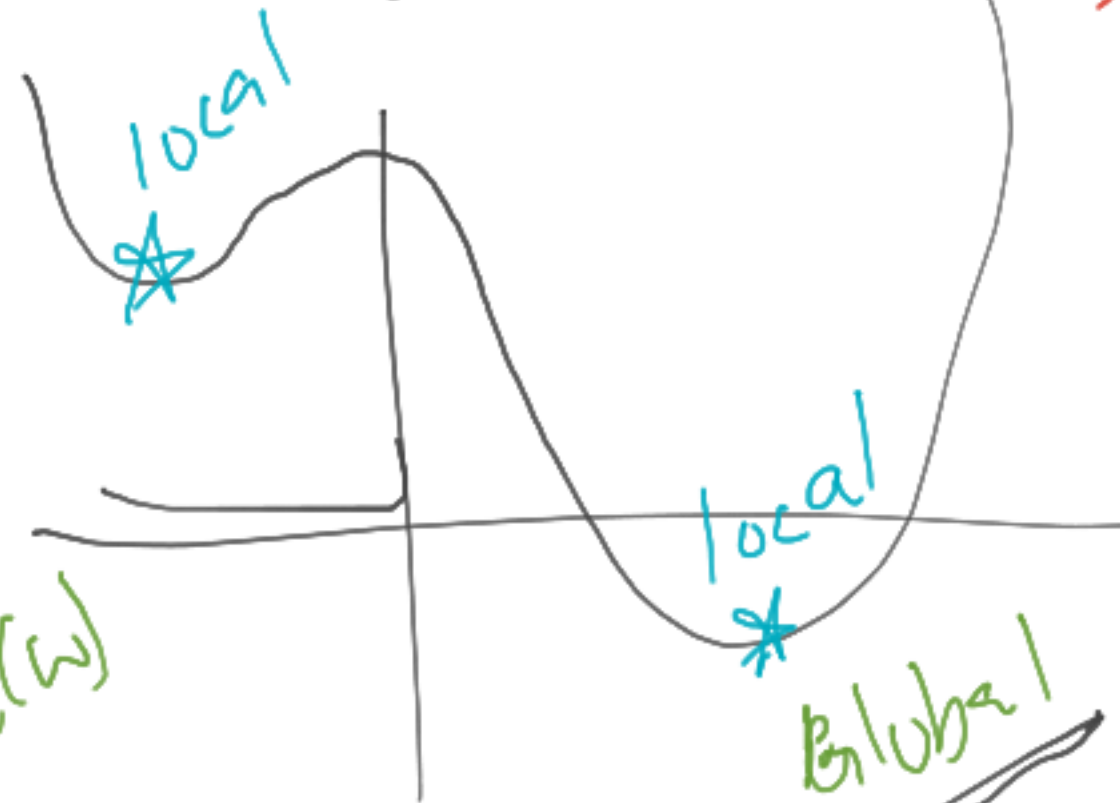
↳ no descent directions.



local

global min

↳ arg min  $f(w)$

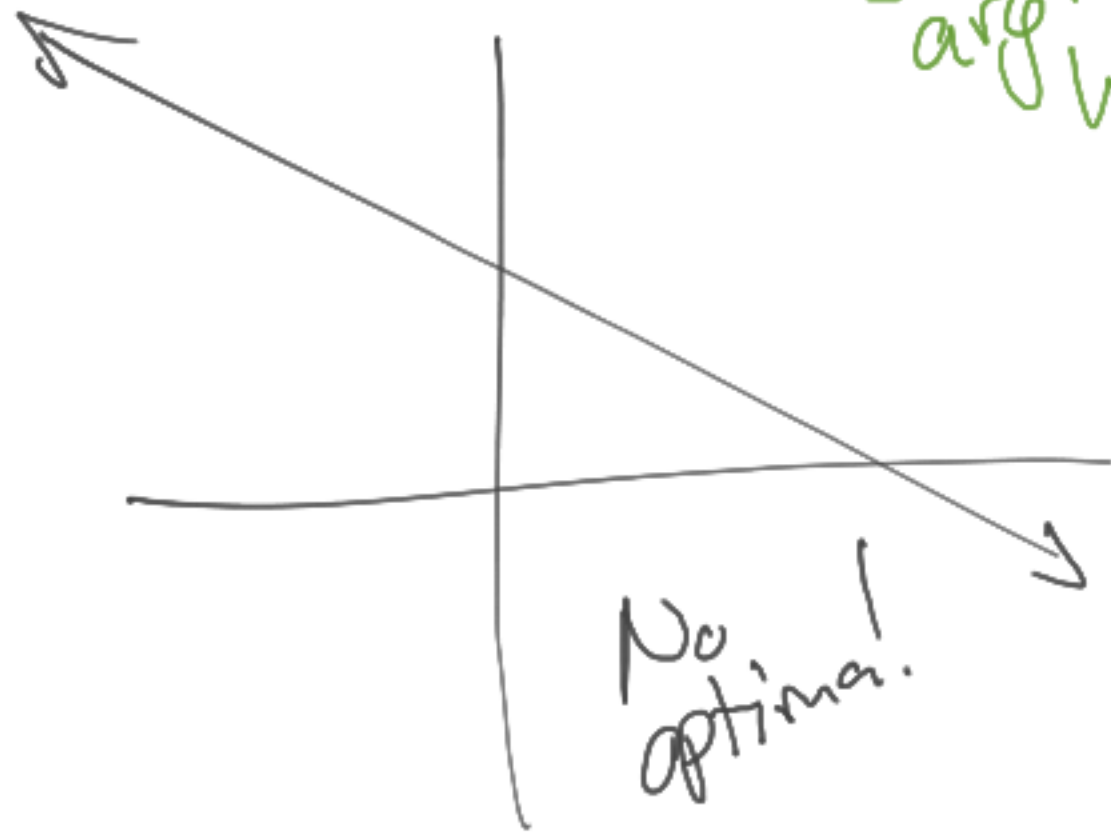


local

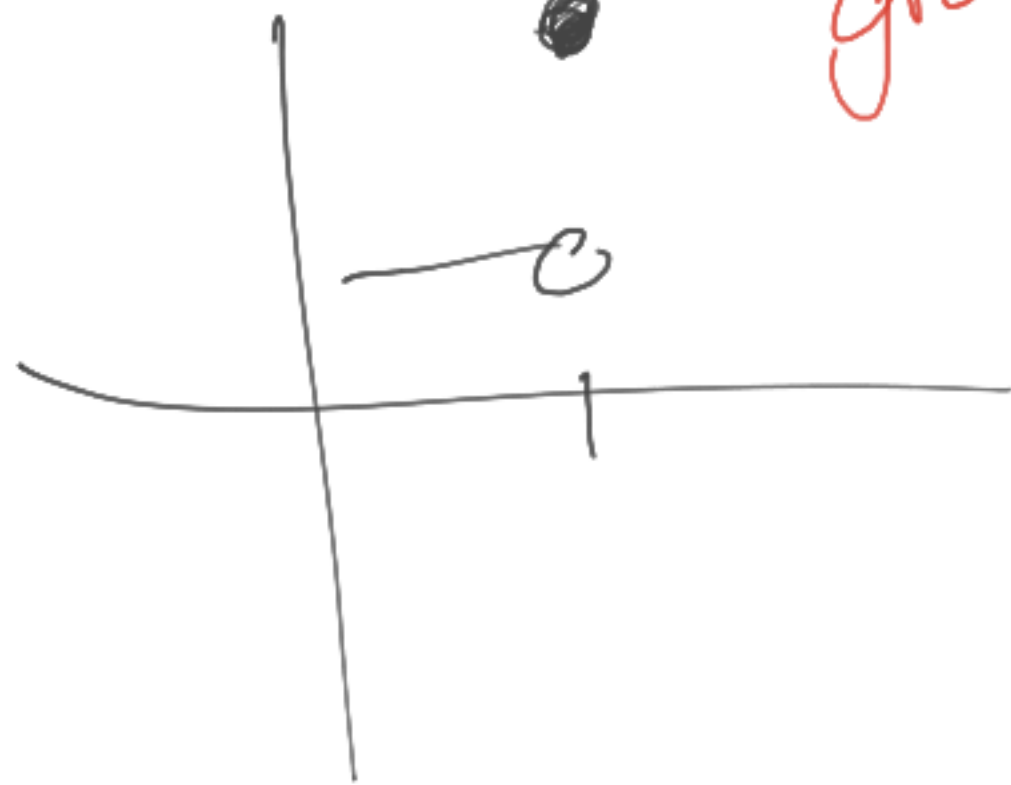
local

global

gradient not defined at discontinuity.

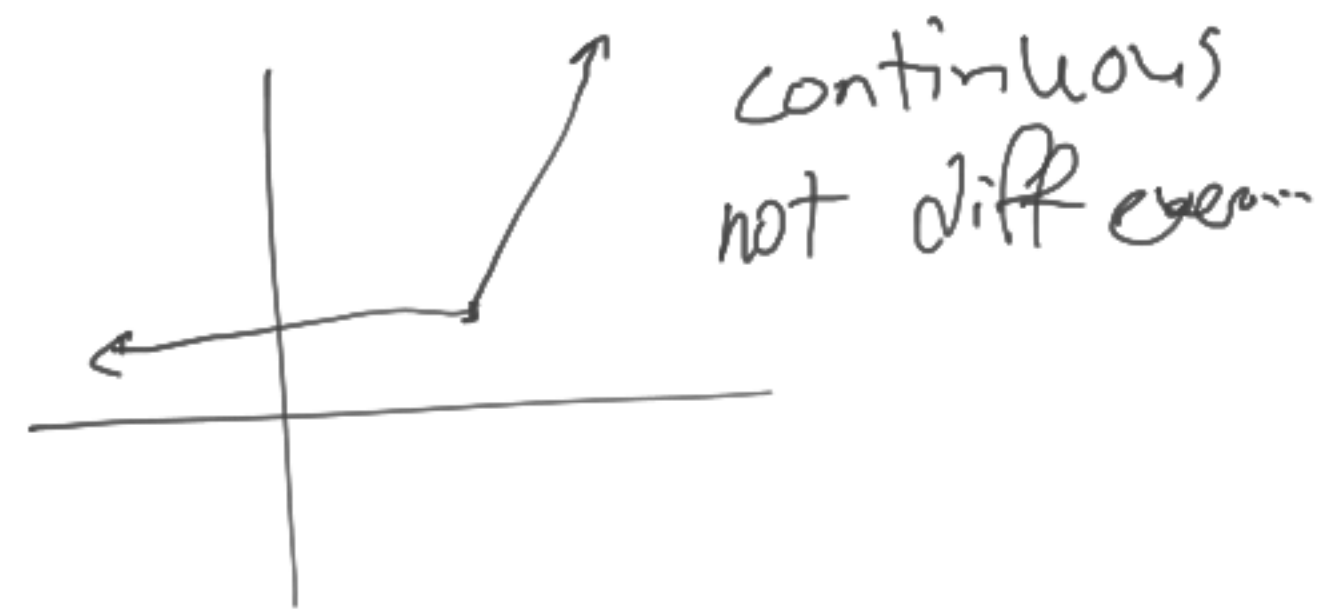


No optimal!



c

- Can diverge if step size too big.
- Requires loss function to be differentiable everywhere.



Linear least squares

$$l(w_j) = \sum_{k=1}^n \left( y_k - \sum_{j=1}^m w_{i,j} x_{k,j} \right)^2$$

$$c w^2 + c_1 w_1 + w_2 + c$$

- quadratic,

- One local min, the global min.

Common requirements for convergence:

★ -  $l$  is "smooth" (differentiable everywhere)

★ -  $l$  is Lipschitz with constant  $L$ .  $\rightarrow$  "maximum slope is  $L$ "

- Step size is ... appropriately set.

$\rightarrow$  Sufficiently small:  $\alpha < \frac{1}{L}$

★  $\rightarrow$  Decayed appropriately. (square summable but not summable).

$$\sum_{i=0}^{\infty} \alpha_i = \infty$$
$$\sum_{i=0}^{\infty} \alpha_i^2 < \infty$$

-  $\nabla l$  is Lipschitz.

-  $l$  is convex, quadratic, strongly convex.



# Types of convergence:

$$w_t \rightarrow \underset{w}{\operatorname{arg\,min}} \ell(w)$$

$$\nabla \ell(w_i) \rightarrow 0$$

$$\nabla \ell(w_i) \rightarrow 0 \quad \underline{\text{or}} \quad \ell(w_i) \rightarrow -\infty$$

$$w_{k+1, j} = w_{k, j} - \alpha \left[ \frac{\partial l(w_k)}{\partial w_{k, j}} \right]$$

$$\frac{\partial l(w_k)}{\partial w_{k, j}} = \frac{\partial}{\partial w_{k, j}} \frac{1}{2} \sum_{i=0}^n (y_i - \hat{y}_i)^2$$

$$= \sum_{i=0}^n \frac{\partial}{\partial w_{k, j}} \frac{1}{2} (y_i - \hat{y}_i)^2$$

$$= \sum_{i=0}^n \cancel{1} (y_i - \hat{y}_i) \frac{\partial}{\partial w_{k, j}} (y_i - \hat{y}_i)$$

$$= \cancel{\sum_{i=0}^n} (y_i - \hat{y}_i) \frac{\partial}{\partial w_{k, j}} \sum_{\beta=1}^3 w_{k, \beta} x_{i, \beta}$$

skipping  
steps

$$= \cancel{\sum_{i=0}^n} (y_i - \hat{y}_i) \frac{\partial}{\partial w_{k, j}} \underbrace{\sum_{\beta=1}^3 w_{k, \beta} x_{i, \beta}}_{z \cdot \text{const} = \text{const}}$$

$$= \cancel{\sum_{i=0}^n} (y_i - \hat{y}_i) x_{i, j}$$

$$l(w_k) = \frac{1}{2} \sum_{i=0}^n (y_i - f_{w_k}(x_i))^2$$

$$f_w(x_i) = \sum_{j=1}^3 w_j x_{ij}$$

$$\frac{\partial}{\partial x} x^2 \rightarrow 2x$$

$$\frac{\partial}{\partial x} f(x)^2 \rightarrow 2f(x) \frac{\partial f(x)}{\partial x}$$







