

## Linear Regression: the ML perspective

Let us now treat the regression problem from an ML perspective! We have seen that closed form solutions to the LLS exist, but we pretend that they are not available. Instead, we formulate and solve via ML. The steps we are taking and the concepts we introduce are part of (almost) all supervised learning methods, in particular Neural Networks → in fact, the algorithm we show here will result in the innermost part of a NN, a linear neuron.

Recap:  $D = \{(x_i, y_i)\}_{i=1}^N$  : data

$f: X \rightarrow Y$ : generating function

$h: X \rightarrow Y$  hypothesis

$H_0$ : model (a specific hypothesis for parameters  $\theta$ )

Supervised learning: given  $D$ , find optimal  $\hat{\Theta}$   
for your choice of  $h$

For an  $(x_j, \underline{y_j}) \in P \cap D = \emptyset$ ,  $h_{\hat{\theta}}(x_j) = \hat{y}_j$

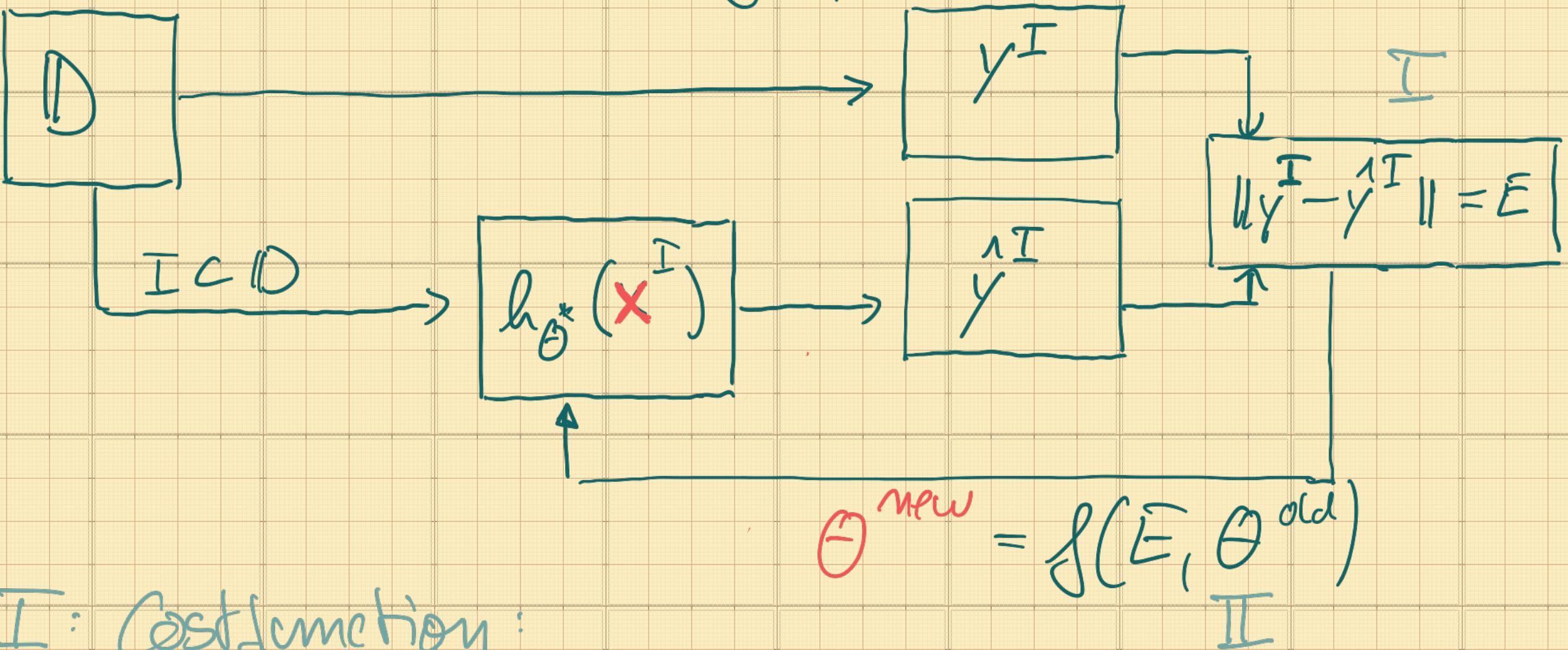
↳ "inference"  
"prediction stage"

?

"training stage"

validation data ( $y_j$  known)  
test data ( $y_j = ?$ )

# The Supervised Learning Cycle: Fail better!

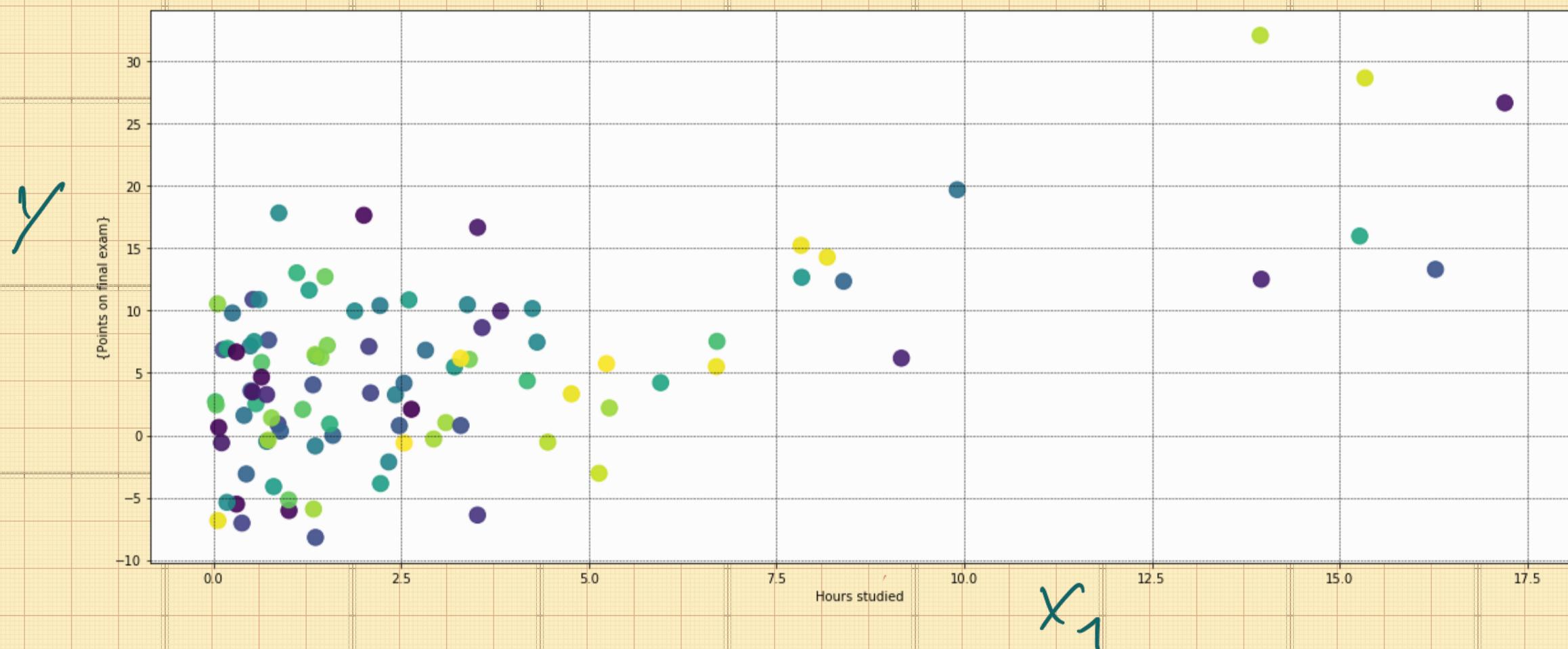


I : Cost function :

How to compute the error

II : how to use knowledge from I to improve  $\theta$

## Univariate case:

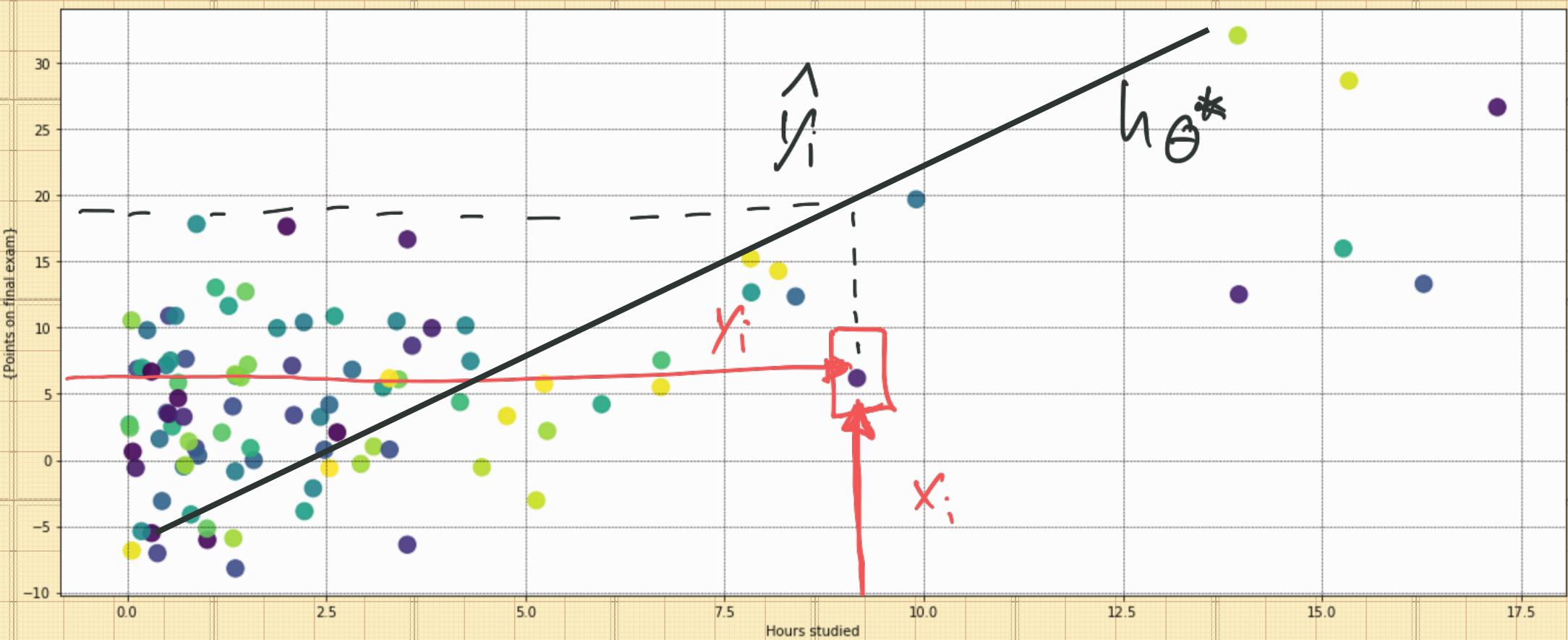


h: univariate linear model

$$h_0: y = x^T w = w_1 x_1 + w_0 x_0 \quad \xrightarrow{=} 1$$

initial guess for  $w_0, w_1$ : random!

for a given sample  $\{x_i, y_i\} \in D$ , the forward pass  
 with current  $\theta^*$  gives :  $h_{\theta^*}(x_i) = \hat{y}_i$



Measure of failure:  $L_2$  / Square error :

for sample i :  $C_i = (\hat{y}_i - y_i)^2 = j_i$

$$C = \frac{1}{2N} \sum_{i=1}^N (\hat{y}_i - y_i)^2$$

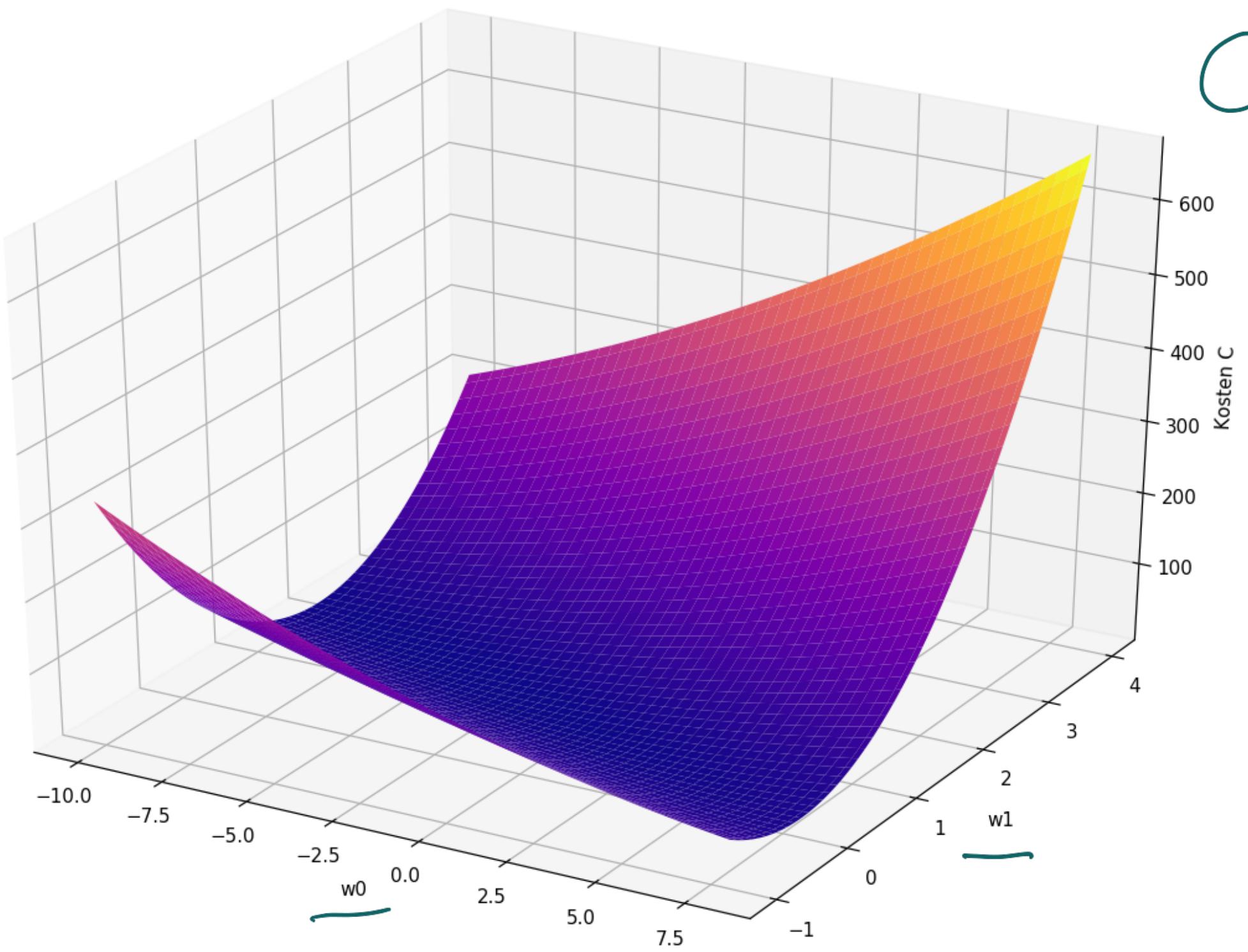
$$\frac{d x^2}{d x} = 2x$$

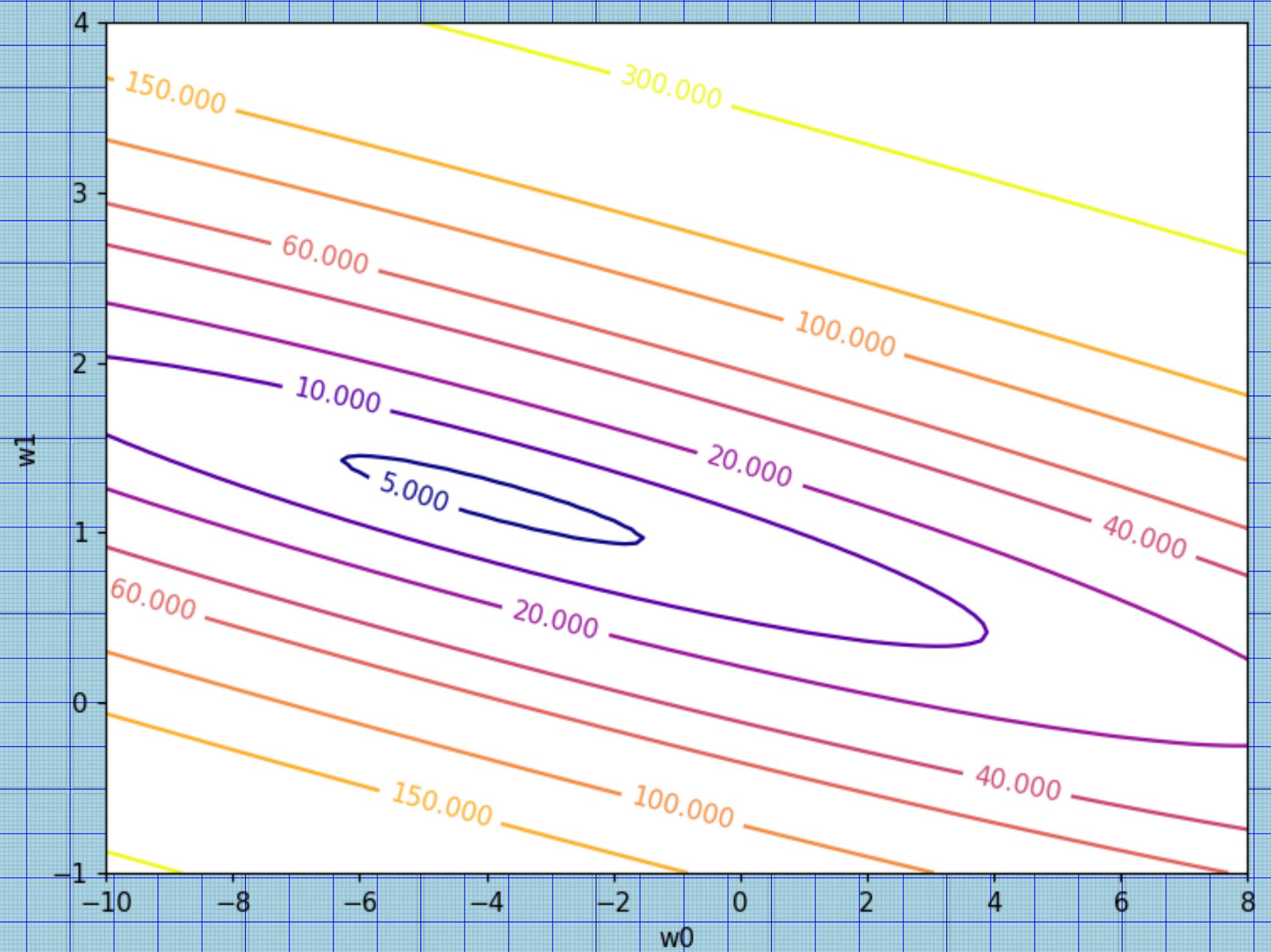
No. Samples

$$\text{so } C = \frac{1}{2N} \sum_{i=1}^N (h_{\theta^*}(x_i) - y_i)^2 = C(\theta) = C((w_0, w_1))$$

Cost function , sampled of  $\theta$

C





How does the cost function help us to achieve goal II:  $\rightarrow$  Gradient Descent

$$\rightarrow \text{Compute } \nabla_{\theta} C = \left( \frac{\partial C}{\partial w_0}, \frac{\partial C}{\partial w_1} \right)^T$$

$$\frac{\partial C}{\partial w_j} = \frac{1}{2N} \sum_{i=1}^N \frac{\partial}{\partial w_j} \left[ \underbrace{\left[ \sum_k w_k x_k(x_i) - y_i \right]}_{h_{\theta}(x_i)} \right]^2$$

linear model:  $\frac{\partial w_i x_i}{\partial w_j} = \begin{cases} x_i & : i=j \\ 0 & : i \neq j \end{cases}$

k: Basis functions / features of linear model

j: which parameters / which gradient?

i: sample

$$\frac{\partial C}{\partial w_j} = - \frac{1}{N} \sum_{i=1}^N (\gamma_i - \hat{\gamma}_i) x_j(x_i)$$

j-th basis evaluated  
 at i-th sample

$C$ : Costfunction, often  
 also called  $J$

Note: this is an average gradient over all samples!

$N = |D|$  : Batch Gradient Descent

$N < |D|$  : Mini-Batch GD

$N = 1$  : Stochastic GD

→ Simple optimization of  $C$  w.r.t.  $\Theta$ :  
iterative gradient descent:

$$w_0^{\text{next}} = w_0^{\text{old}} - \alpha \frac{\partial C}{\partial w_0}$$

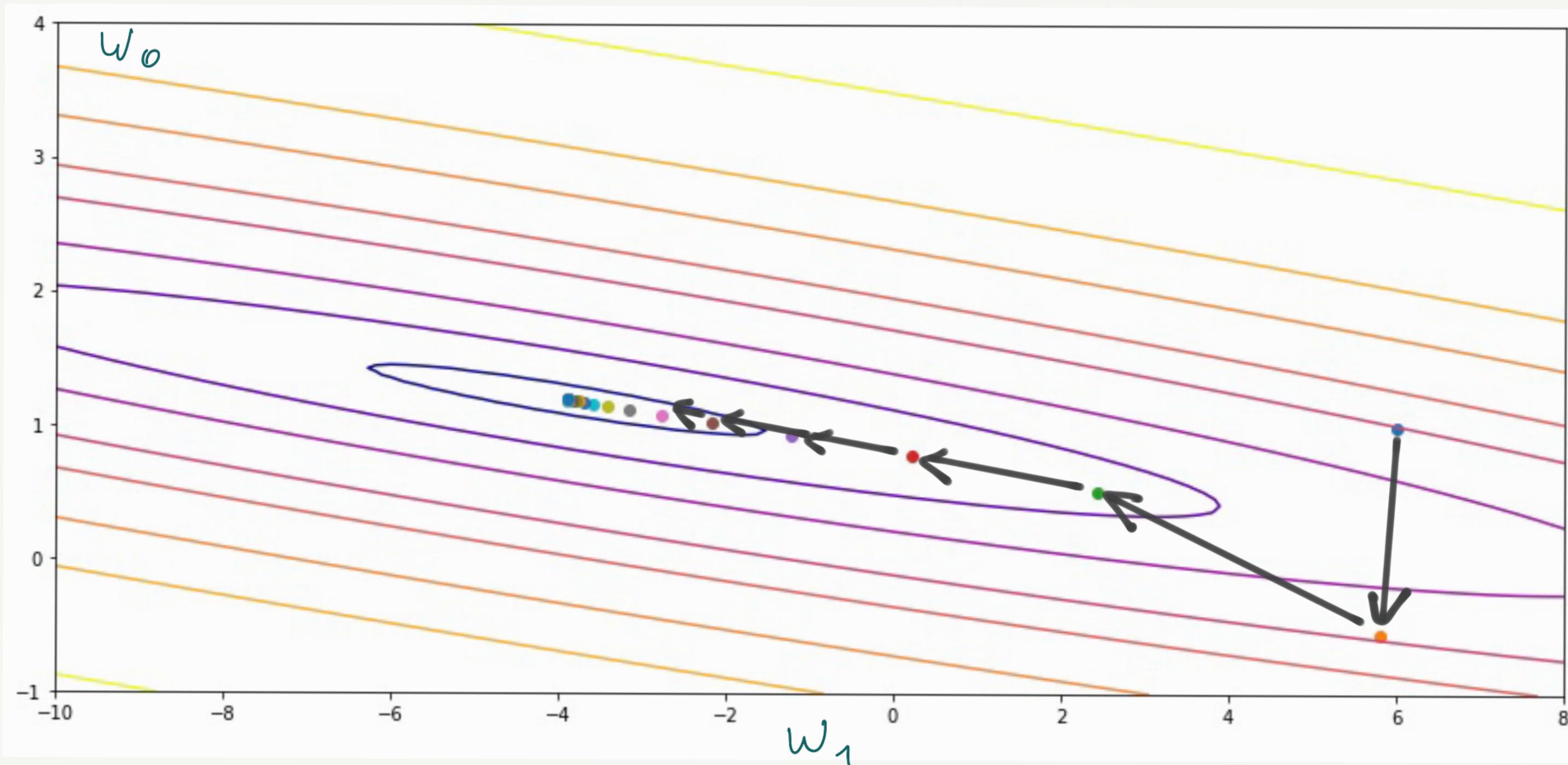
$$w_1^{\text{next}} = w_1^{\text{old}} - \alpha \frac{\partial C}{\partial w_1}$$



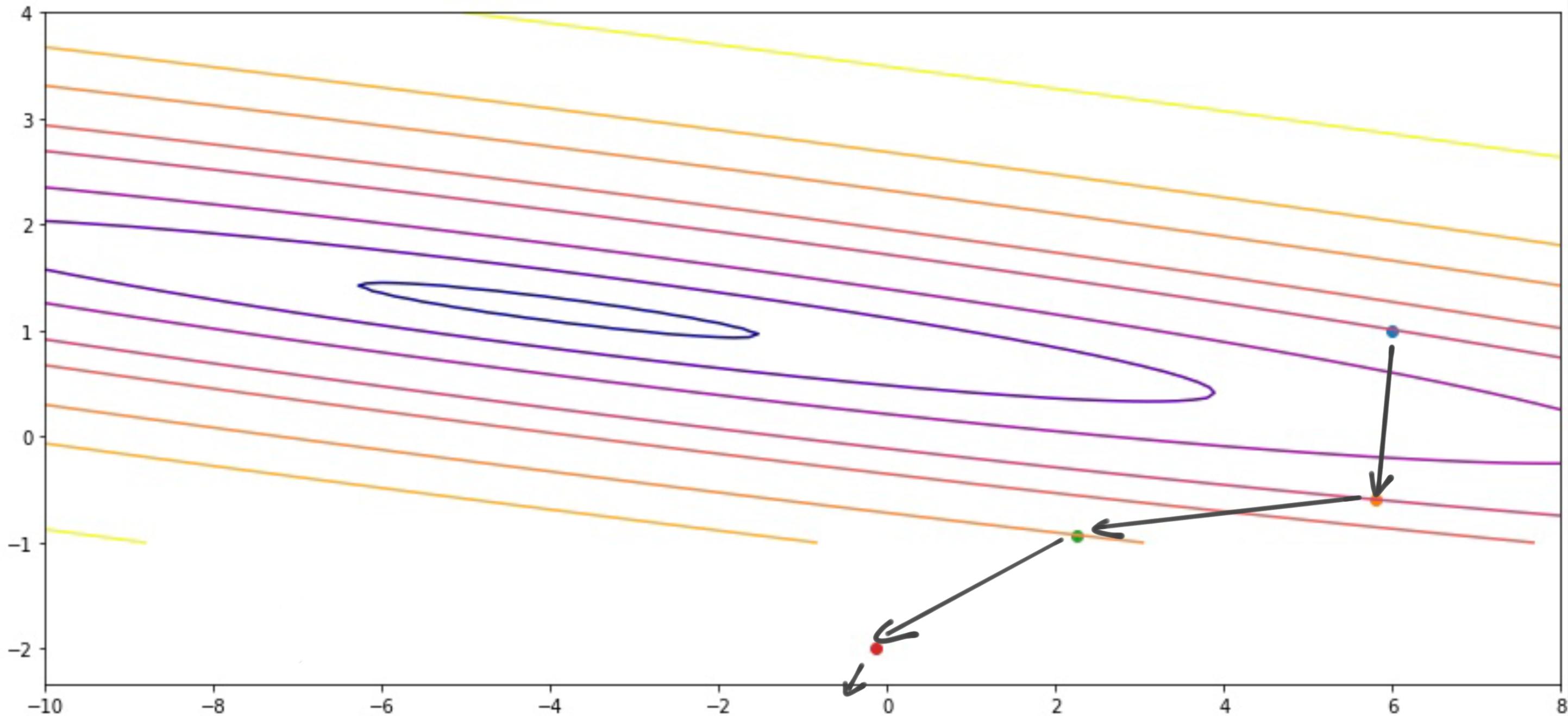
step size: "learning rate"  
 $\rightarrow$  hyperparameter

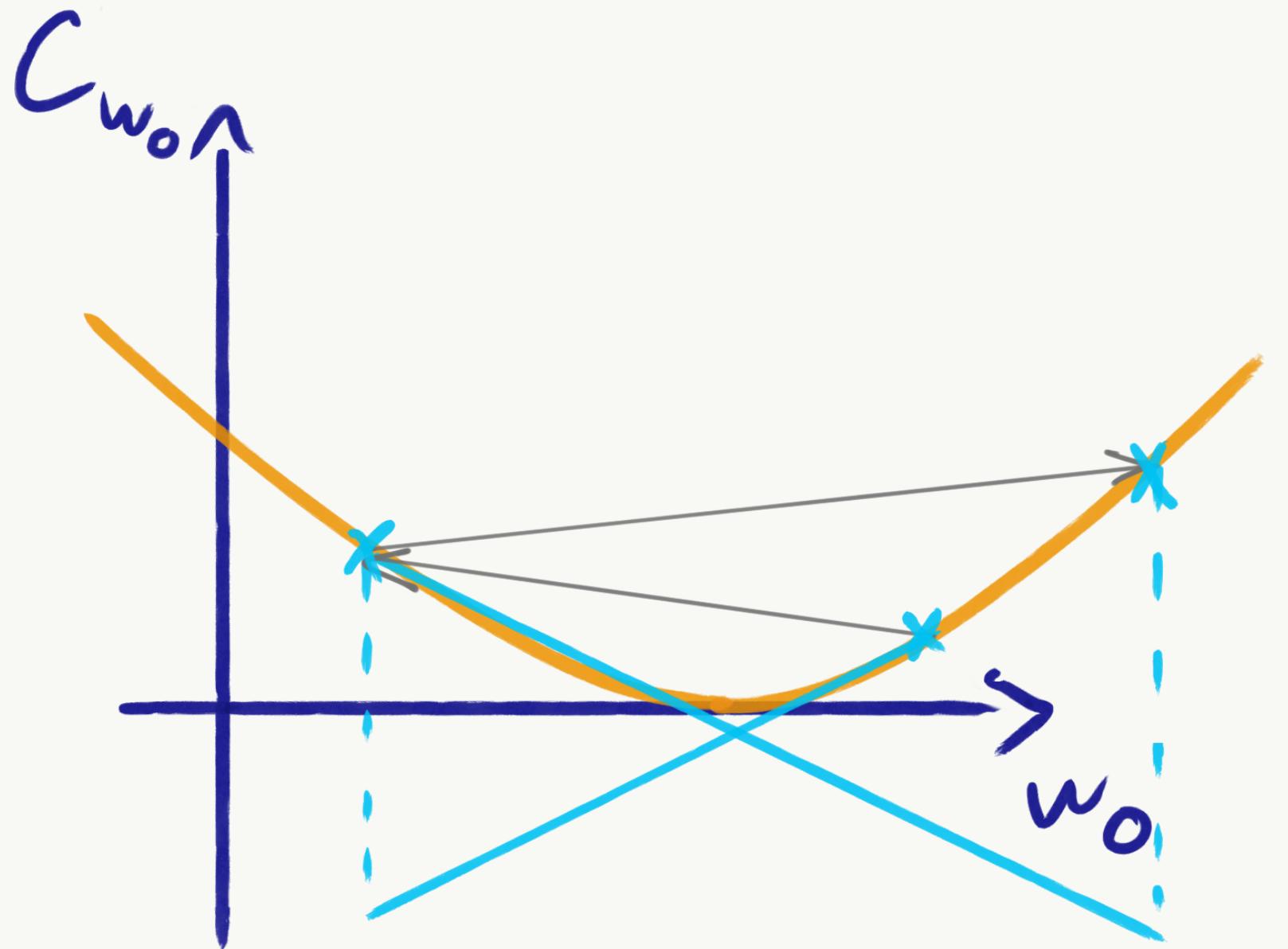
iterate and update!

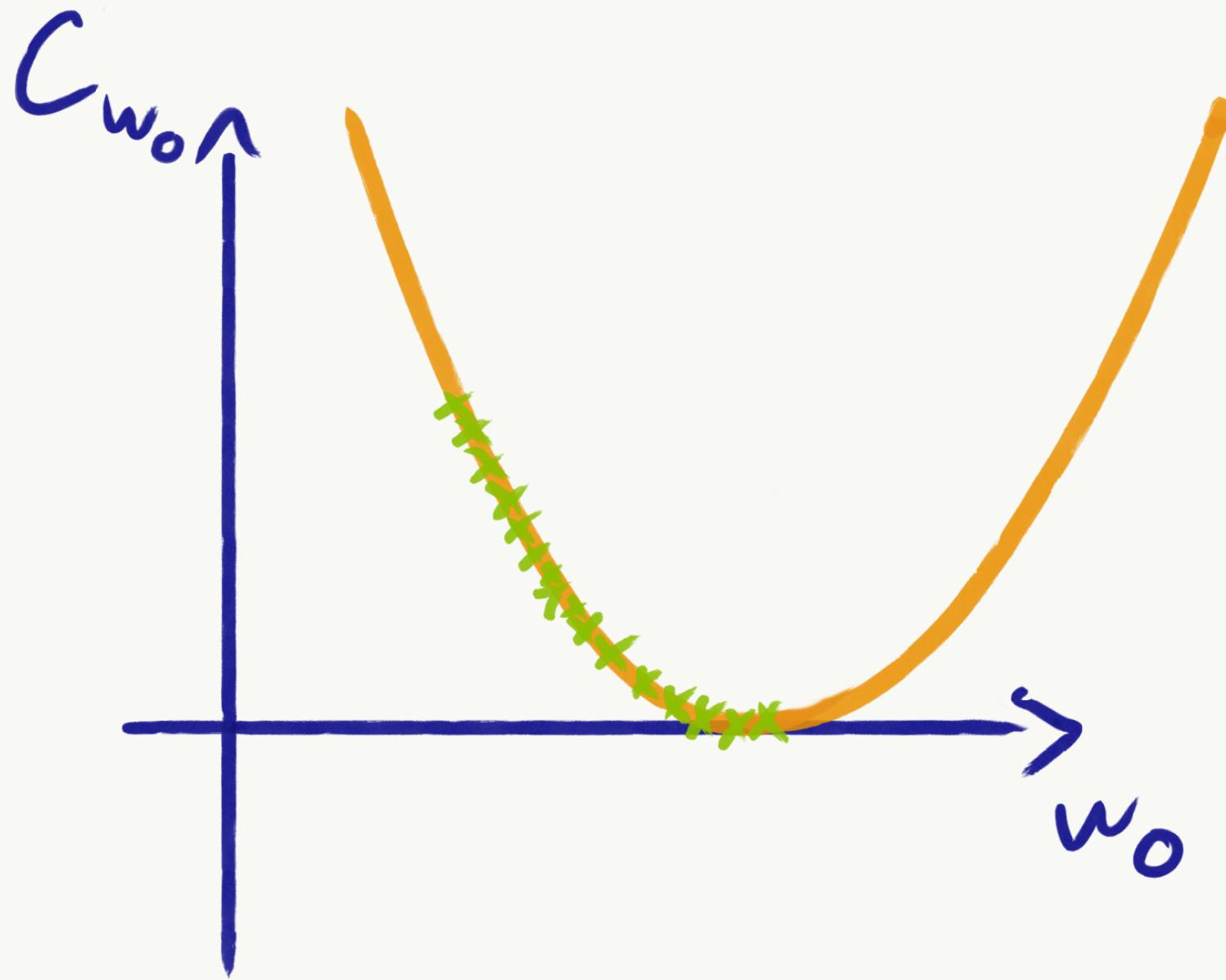
# Successful optimization

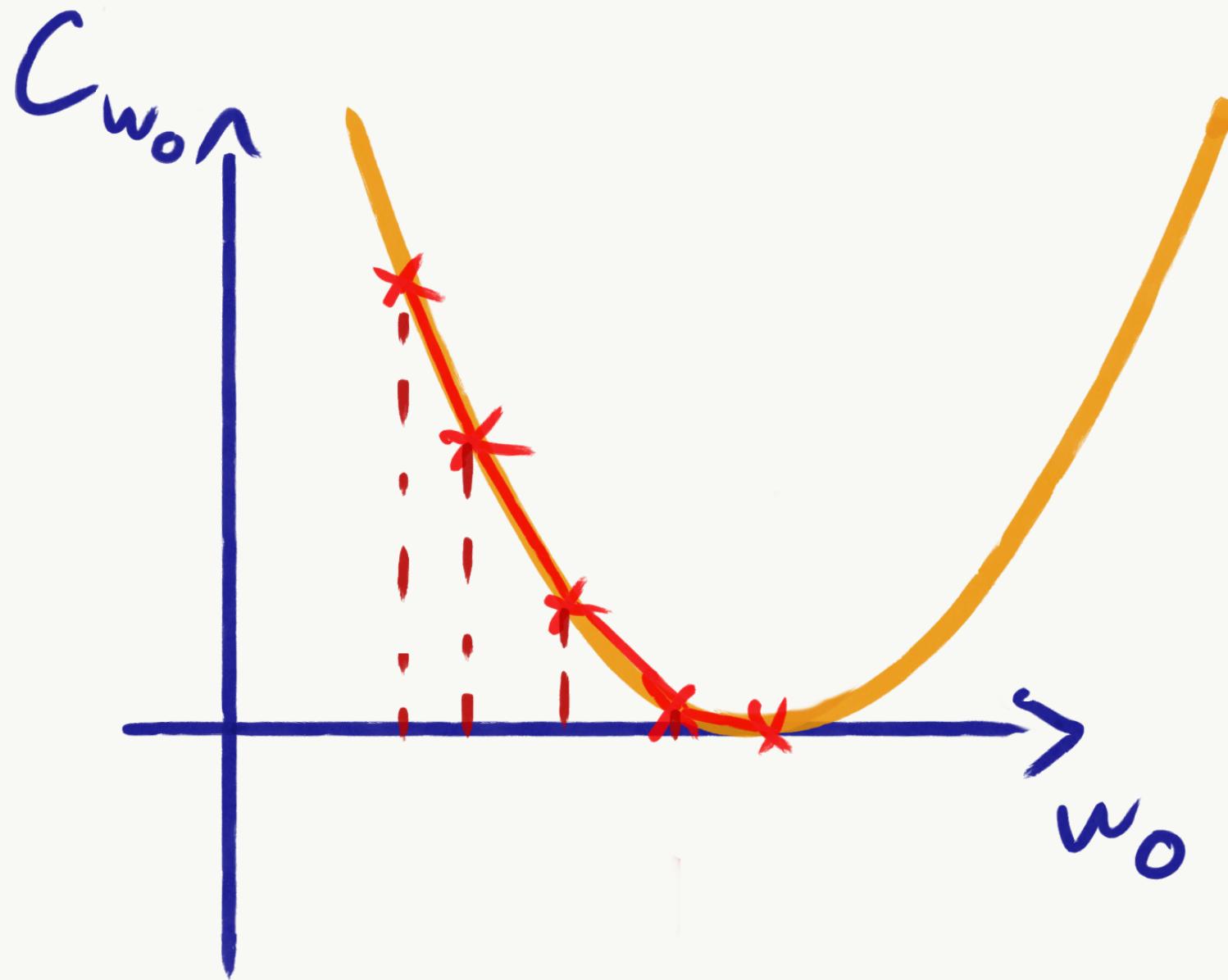


Learning Rate too large  $\rightarrow$  diverging









- Note that if done right, this method converges to the LS we discussed before!
- We can increase the expressiveness by
  - more features:  $w_0x_0 + w_1x_1 + w_2x_2 + \dots$
  - high order ansatz:  $x^2, x^3, x_1 \cdot x_2, \dots$
- Multivariate LS: nothing new!
- Linear Neuron as the building block of NNs:

