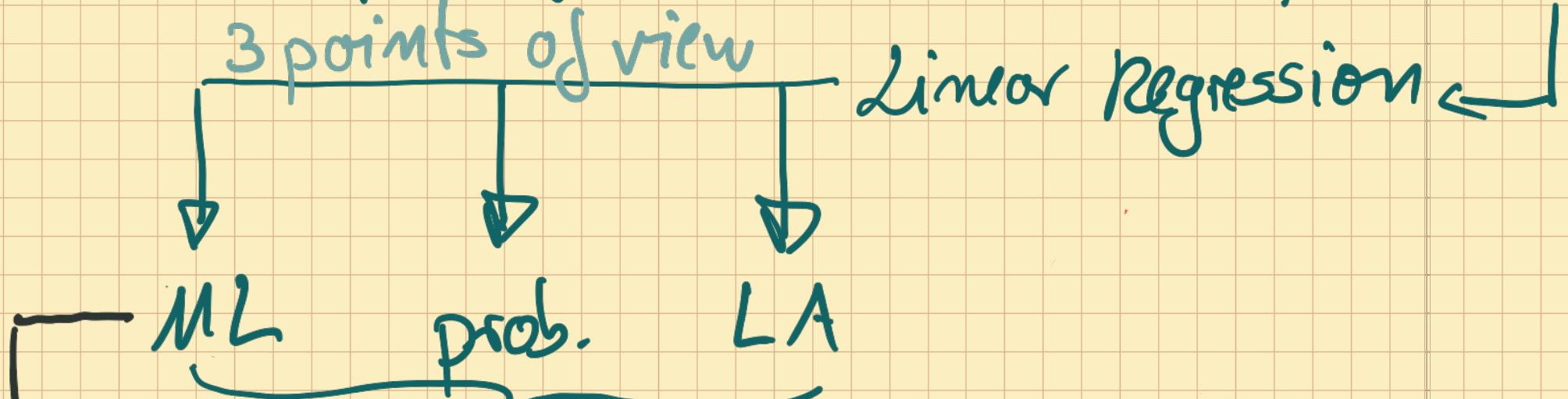


Welcome to shallow ML!

Types of Learning → Supervised Learning

3 points of view



↳ Linear Neuron → Classification

Logistic Neuron ← Logistic Regression ←

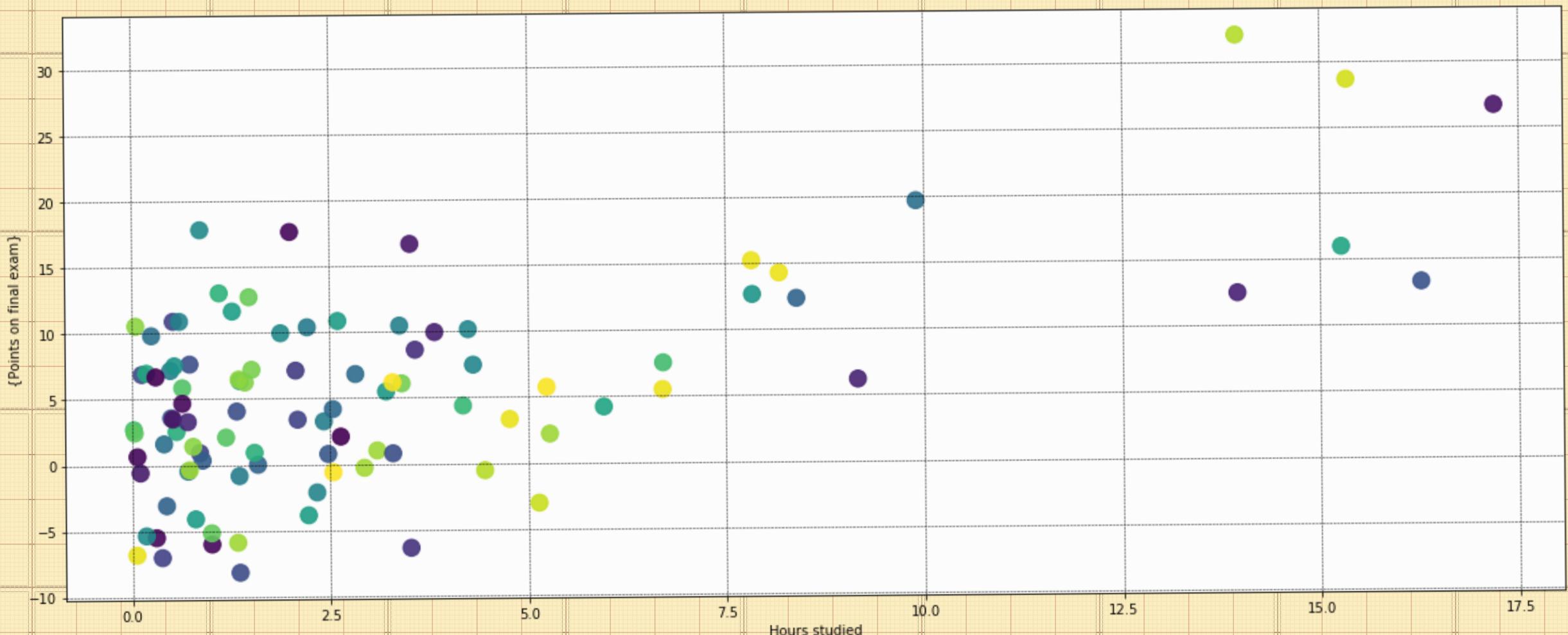
→ Assignment IV

Types of Learning

- 1.) unsupervised learning: descriptive model
unlabeled data; find lower dimensional representation /model of data $P(x|D)$
- 2.) Supervised learning: predictive model
Mappings from input to output given data D
 $\{(x_i, y_i)\}_{i=1}^N := D$
 $P(y|x, D)$ or $f(x, D) \rightarrow$ model as function
↳ model as a probability distribution
- 3.) Reinforcement, semi-supervised ...

Supervised Learning: Nomenclature

A sample problem: Grade on final test



$x_i \in X$: Input : hours spent studying

$y_i \in Y$: Output : points on test

$f: X \rightarrow Y$: Generating function : the truth

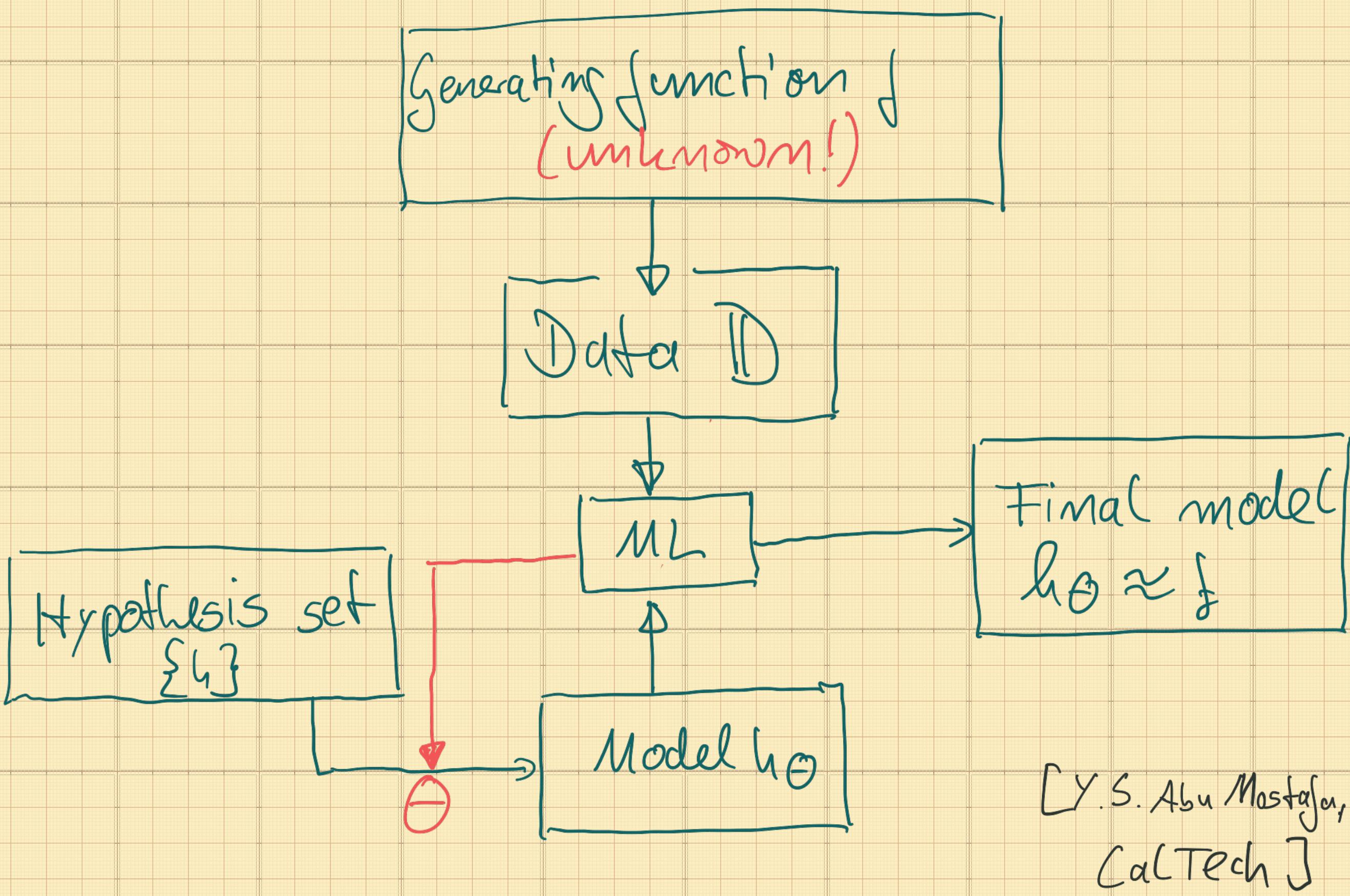
$\{(x_1, y_1), (x_2, y_2), \dots, (x_N, y_N)\} =: D$: data

$h: X \rightarrow Y$: hypothesis set

Θ : parameters : trainable

$h_\theta: X \xrightarrow{\Theta} Y$: Model : a specific choice from $\{h\}$

$h_\theta(x^*)$: Prediction



Linear Regression : 3 Approaches

↳ continuous functions \leftrightarrow classification

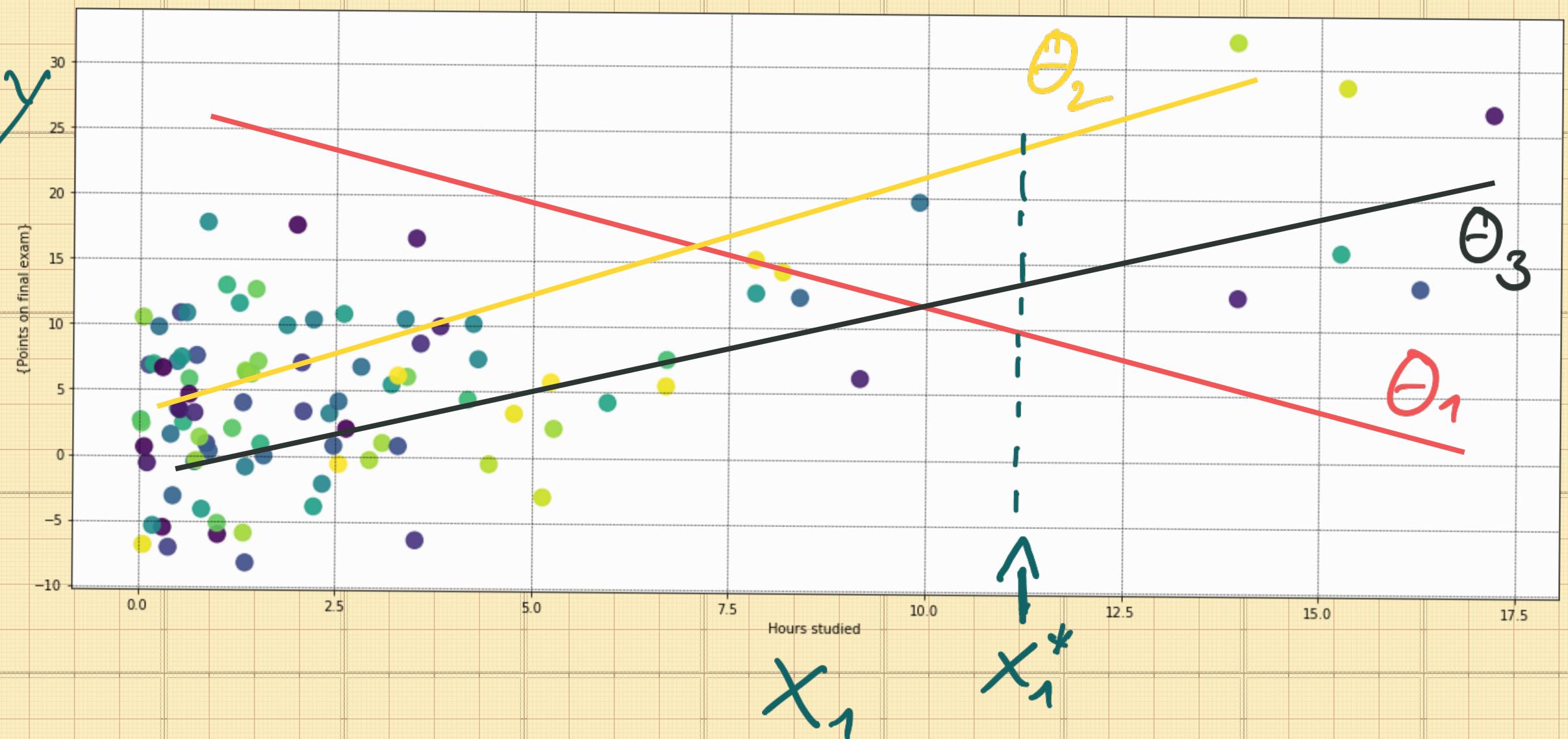
$$f: x \rightarrow y$$
$$h: y = x^T \theta = \sum_{i=1}^n w_i x_i$$

We consider the univariate case here where we have a single feature x_1 :

$$h = w_1 x_1 + w_0 x_0 = w_1 x_1 + \underbrace{w_0 \cdot 1}_{\text{bias term}}$$

Linear : refers to the parameters !

$h = w_1 x_1^2 + w_0 x_0$ is still LR (with feature engineering)



Approach 1: Classical Linear Algebra
→ Linear Least Squares → Normal form

Approach 2: Max. Likelihood Estimate
→ Gaussian Noise → Normal form

Approach 3: Iterative Learning
→ Square Loss → approx. LLS

→ Most famous shallow ML algorithm!

The linear algebra perspective

(heavily inspired by Pavel Grinstejn)

A) Quadratic Form Minimization :

A function $f: \mathbb{R}^N \rightarrow \mathbb{R}$ of the form

$$f(\underline{x}) = \frac{1}{2} \underline{x}^T \underline{A} \underline{x} - \underline{x}^T \underline{b}$$
 with \underline{A} square, positive definite

is called the "quadratic form" (with linear extension).

$$\text{Then, } \min_{\underline{x}} f(\underline{x}) : \underline{A} \underline{x} = \underline{b}$$

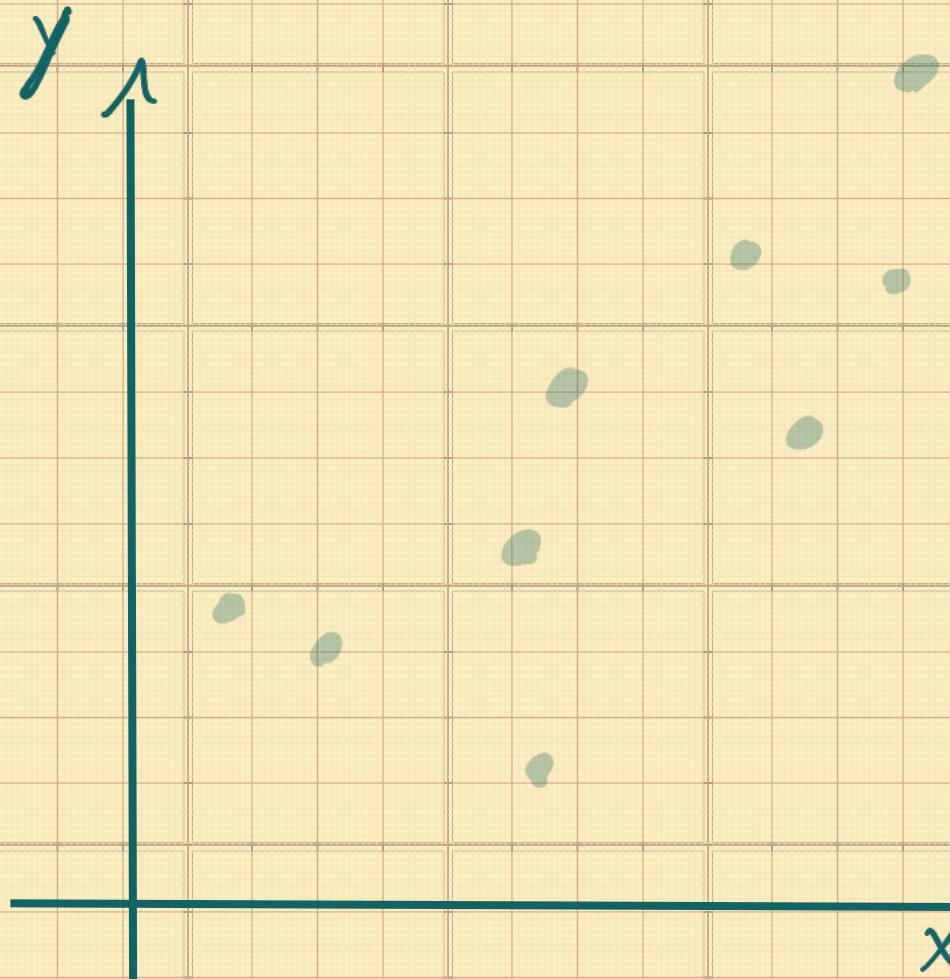
One way to show this: $\underline{x} \in \mathbb{R}^3$, compute $\nabla f = 0$

→ solving the linear system $\boxed{\underline{A}\underline{x} = \underline{b}}$
minimizes the associate quadratic

form $\frac{1}{2} \underline{x}^T \underline{A} \underline{x} - \underline{x}^T \underline{b}$.

→ Scalar case: $f(x) = \frac{1}{2} ax^2 - bx$,
 $f'_x = ax - b \stackrel{!}{=} 0$
⇒ $\boxed{ax = b}$

B) Linear least squares



→ Find a best linear fit
to the data $\{(x_i, y_i)\}_{i=1}^n = D$

$$\rightarrow h_{\theta}(x) = w_1 x + w_0 \cdot 1$$

$$\Theta = (w_1, w_0)^T$$

Basis functions

→ how to fit the model

$h_{\theta}(x)$ on the data D ?

→ Interpolation leads to an overdetermined system:

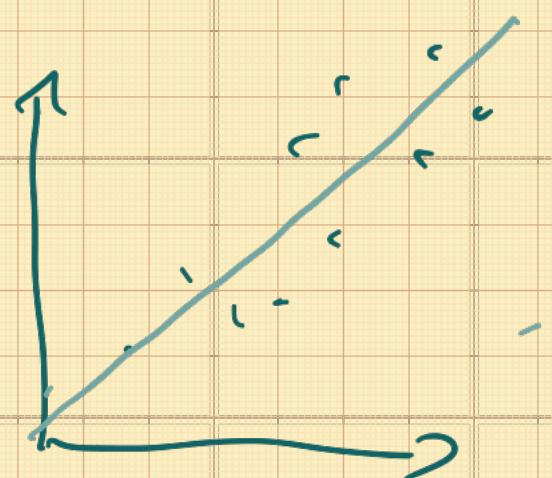
$$y_i = w_1 x_i + w_0 \cdot 1 \quad i = 1, \dots, N$$

2 unknowns, N constraints

in matrix vector form with $x = \phi_1, 1 = \phi_0$:

$$\begin{pmatrix} \phi_1(x_1) & \phi_0(x_1) \\ \phi_1(x_2) & \vdots \\ \phi_1(x_3) & \vdots \\ \vdots & \vdots \\ \phi_1(x_N) & \phi_0(x_N) \end{pmatrix} \begin{pmatrix} w_1 \\ w_0 \end{pmatrix} = \begin{pmatrix} y_1 \\ \vdots \\ \vdots \\ y_N \end{pmatrix}$$

↳ Vandermonde - Matrix



A

x = b

→ cannot be inverted exactly!

→ What can we say about the error?

$$\underline{r} = \underline{b} - \underline{\underline{A}} \underline{x}$$

→ need norm for $\underline{r} \rightarrow$ inner product!

$$\begin{aligned}\underline{r}^T \underline{r} &= (\underline{b} - \underline{\underline{A}} \underline{x})^T (\underline{b} - \underline{\underline{A}} \underline{x}) \\ &= \underline{x}^T \underline{\underline{A}}^T \underline{\underline{A}} \underline{x} - \underline{x}^T \underline{\underline{A}}^T \underline{b} - \underline{b}^T \underline{\underline{A}} \underline{x} + \underline{b}^T \underline{b}\end{aligned}$$

$$\downarrow \quad \underline{b}^T (\underline{x}^T \underline{\underline{A}}^T)^T = \underline{b}^T (\underline{\underline{A}} \underline{x})$$

$$\begin{aligned}&= 2 \left(\frac{1}{2} \underline{x}^T \underline{\underline{A}}^T \underline{\underline{A}} \underline{x} - \underline{x}^T \underline{\underline{A}}^T \underline{b} \right) + \underline{b}^T \underline{b} \\&\quad \underbrace{\qquad\qquad}_{\text{"A" }} \qquad \underbrace{\qquad\qquad}_{\text{"b" }} \qquad \underbrace{\qquad\qquad}_{\text{=Const.}}\end{aligned}$$

$$= 2 \cdot \left(\frac{1}{2} \underline{x}^T \underline{\underline{\tilde{A}}} \underline{x} - \underline{x}^T \underline{\tilde{b}} \right) + \text{const.}$$

→ The term in brackets is a quadratic form!
 (the position of the minimum does not change through addition of a constant or scaling)

→ This error becomes minimal if we thus

solve $\tilde{\underline{A}} \underline{x} = \tilde{\underline{y}}$

or $\underline{\underline{A}}^T \underline{\underline{A}} \underline{x} = \underline{\underline{A}}^T \underline{b}$

$$\Rightarrow \underline{x} = \frac{\underline{\underline{A}}^T \underline{b}}{\underline{\underline{A}}^T \underline{\underline{A}}}$$

$$\begin{pmatrix} w_1 \\ w_0 \end{pmatrix} = \underline{\theta}^T$$

→ The famous normal equation of linear least squares!

What did we learn from this?

- Solving a linear system minimizes the associated quadratic form
- Thus, for a linear hypothesis, a quadratic error is natural choice
- A closed form solution to the linear least squares problem is given by the normal form
- This is thus a nice test case for learning algorithms!
- The normal form corresponds to a discrete L_2 -projection onto P_1 with standard inner product:

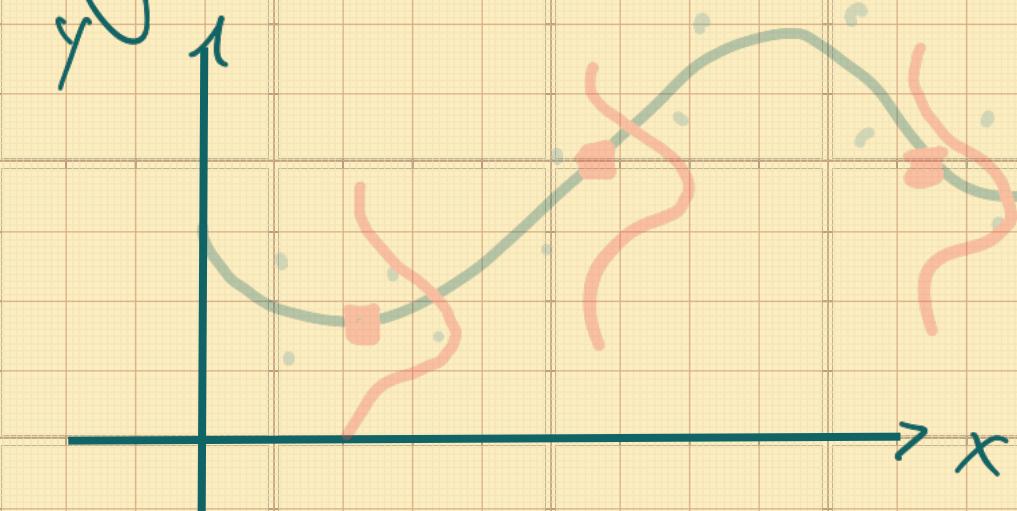
Linear Least Squares

Probabilistic P.D.V.: Signal/data contains noise, e.g. measurement error

→ regression with likelihood:

$$p(y|x) = \mathcal{N}(y | f(x), \sigma^2)$$

here, x is again the input, y is the output/target and we have $y = f(x) + \epsilon$, where ϵ is an i.i.d. Gaussian noise with zero μ and variance σ^2



- Note: $f(x)$ is the hidden generating function, that we can only observe through sampling
- Note: we assume that σ^2 is known and constant
- Linear regression means linear in the parameters Θ : (features can be combined non-linearly):

$$f(x) \approx w_1 x + w_0 \cdot 1 = x^T \Theta$$

$$= \begin{pmatrix} x \\ 1 \end{pmatrix}^T \begin{pmatrix} w_1 \\ w_0 \end{pmatrix}$$

↙
features

↘
parameters

$$\text{so: } p(y|x, \theta) = \mathcal{N}(y | x^T \theta, \sigma^2), *$$

$$y = f(x) + \epsilon = x^T \theta + \epsilon, \epsilon \sim \mathcal{N}(0, \sigma^2)$$

↑
linear in θ

Remark : iff $\sigma^2 \rightarrow 0$: $p(y|x, \theta) \rightarrow \delta$

How is (*) useful in finding a model?

⇒ classic estimator: MLE (see lecture statistics II)

maximizing the likelihood / prediction

of the training data D given the unknown, but fixed parameters θ :

$$\theta_{\text{MLE}} := \arg \max_{\theta} p(D, \theta)$$

so more specifically:

$$\theta_{MLE} = \arg \max_{\theta} P(\underline{Y} | \underline{\bar{X}}, \theta)$$

feature
1
↑

set of N
training samples

$$\left(\begin{array}{c} (y_1, y_2, \dots, y_N)^T \\ \vdots \\ (x_1^0, x_1^1, x_2^0, x_2^1, \dots, x_N^0, x_N^1, \dots) \end{array} \right)$$

Since we assume the samples to be i.i.d., we can factorize the likelihood as:

$$P(\underline{Y} | \underline{\bar{X}}, \theta) = \prod_{m=1}^N P(y_m | x_m, \theta)$$

$$= \prod_{m=1}^N \mathcal{N}(y_m | x_m^T \theta, \sigma^2)$$

**

Remark: $P(\underline{Y} | \underline{X}, \Theta)$ is not a pdf in Θ , but it is one in \underline{Y}

Instead of maximizing (**), it is numerically easier to minimize the log-likelihood.

Since $\log(x)$ is strictly monotone, the extrema of x are identical to those of $\log(x)$!

[MMLB, §.2]

$$\begin{aligned} NLL: NLL(\Theta) &:= -\log P(\underline{Y} | \underline{X}, \Theta) \\ &= -\log \prod_{n=1}^N p(y_n | x_n, \Theta) \\ &= -\sum_{n=1}^N \log p(y_n | x_n, \Theta) \end{aligned}$$

Since we have $p(y_n | x_n, \theta) = \mathcal{N}(y_n | x_n^T \theta, \sigma^2)$:

$$\begin{aligned} NLL(\theta) &= -\sum_{n=1}^N \log \mathcal{N}(y_n | x_n^T \theta, \sigma^2) \\ &= -\sum_{n=1}^N \log \frac{1}{\sqrt{2\pi\sigma^2}} \exp -\frac{(y_n - x_n^T \theta)^2}{2\sigma^2} \end{aligned}$$

$$= -\sum_{n=1}^N \log \exp \left(\frac{-(y_n - x_n^T \theta)^2}{2\sigma^2} \right) + \text{const.}$$

$$= \frac{1}{2\sigma^2} \sum_{n=1}^N \frac{(y_n - x_n^T \theta)^2}{2\sigma^2} + \text{const'}$$

\downarrow
no influence
on minimum
location!

so, we need to find the minimum of

$$NLL^*(\theta) = \sum_{n=1}^N (y_n - \mathbf{x}_n^T \theta)^2$$

This is simply an L_2 -error / quadratic form minimization as we have seen before!

To make the connection clearer, rewrite NLL^* in vector form:

$$NLL^*(\theta) = (\mathbf{Y} - \mathbf{X}\theta)^T (\mathbf{Y} - \mathbf{X}\theta)$$

We have seen before that this can be seen as a quadratic form, for which the minimum is found by solving the associated linear system, which leads again to the

normal equation of LLS:

$$\hat{\theta}_{MLE} = \frac{\underline{\underline{X}}^T \underline{\underline{Y}}}{\underline{\underline{X}}^T \underline{\underline{X}}} \quad \rightarrow$$

see talk on
normal form.

→ More on MLE and MAP estimator +
regularized least squares : Bishop Sec 3.1.

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 θ)

Supervised learning: given D , find optimal $\hat{\Theta}$

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

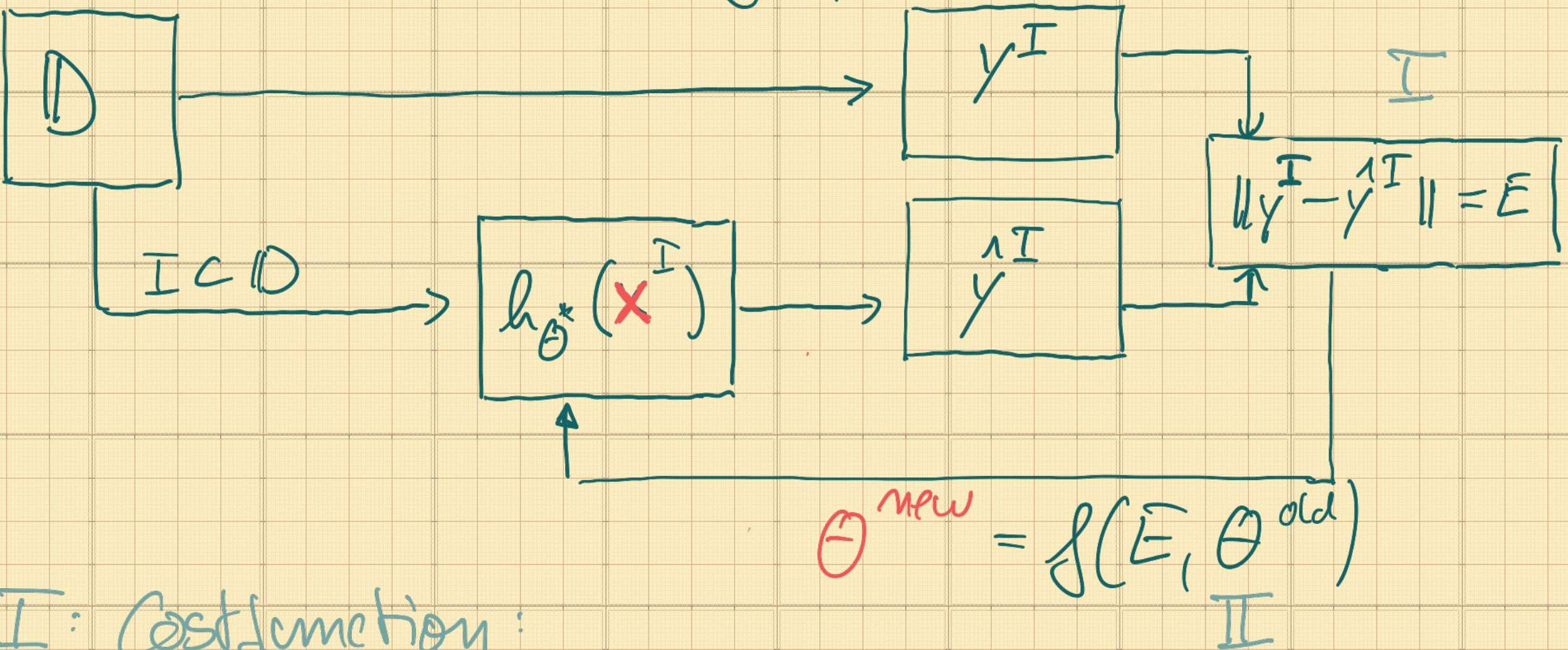
↳ "inference"
"prediction stage"

?

"training stage"

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

The Supervised Learning Cycle: Fail better!

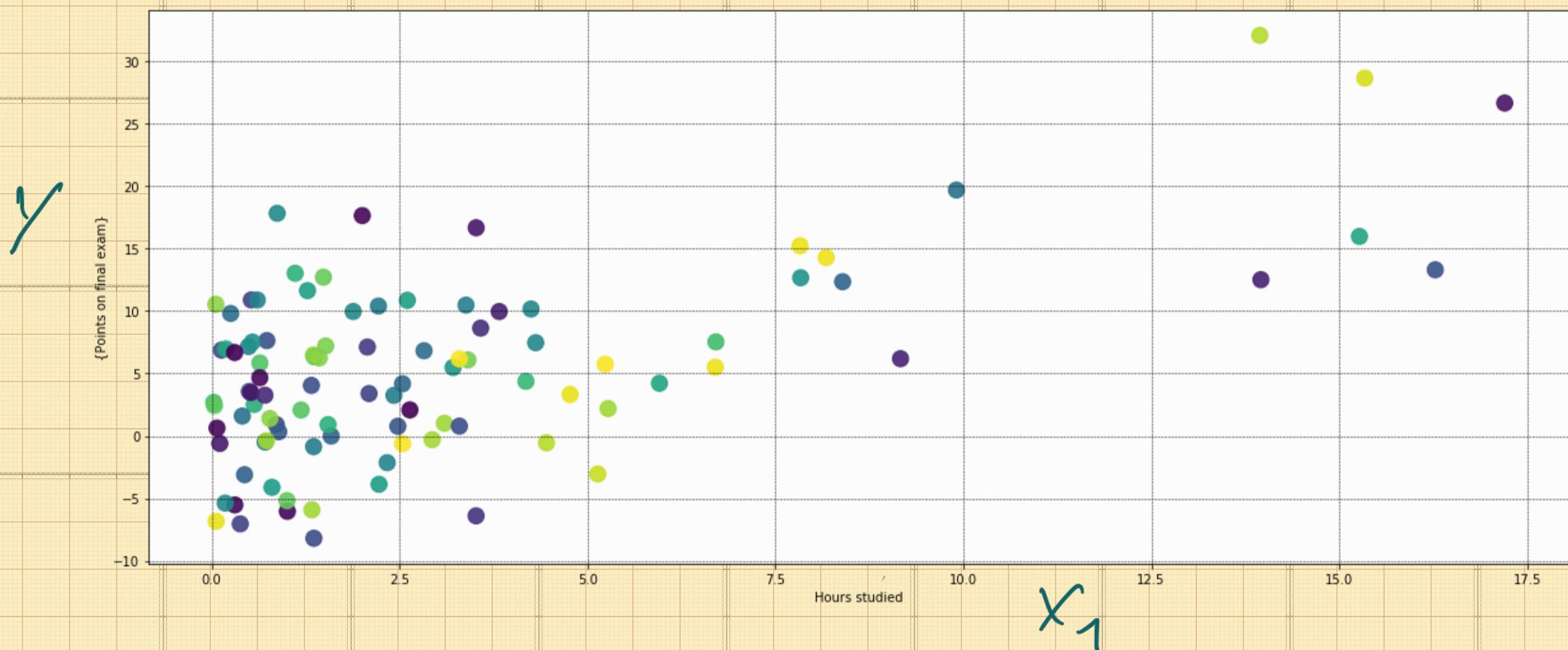


I : Cost function :

How to compute the error

II : how to use knowledge from I to improve θ

Univariate case:

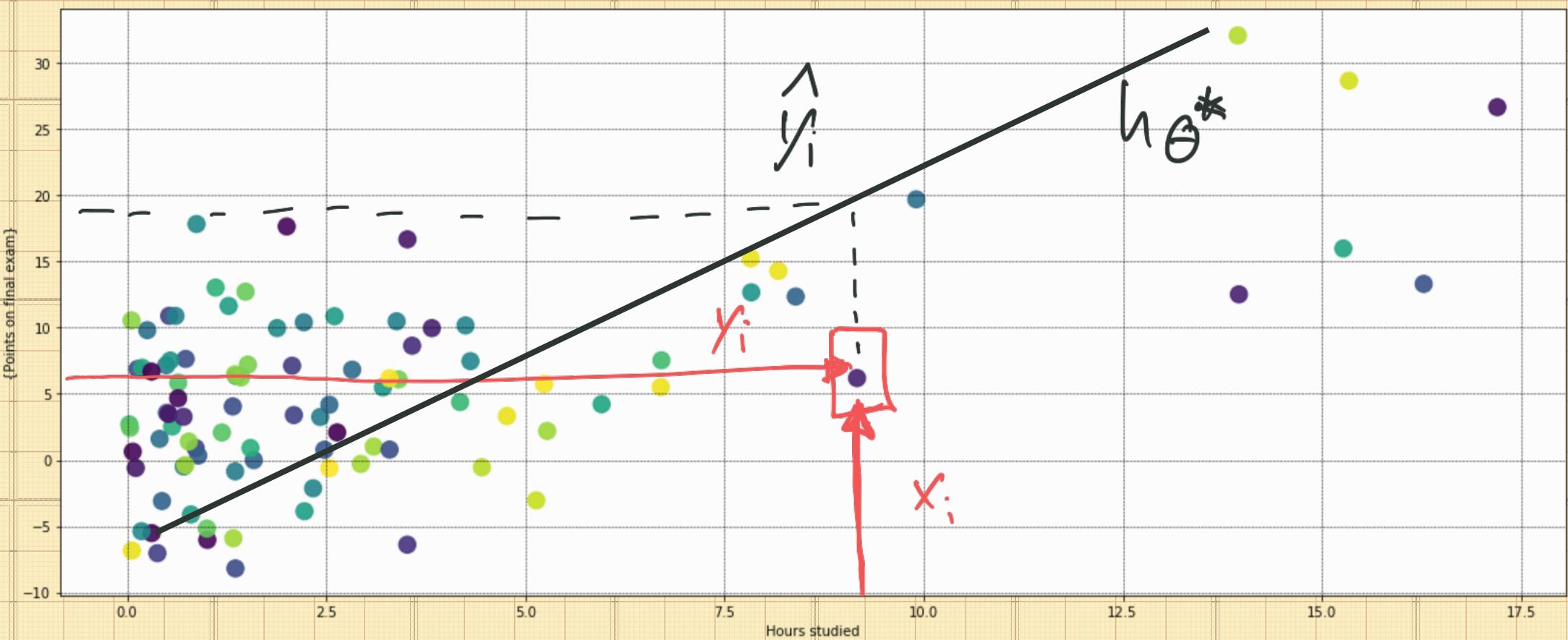


h: univariate linear model

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

initial guess for w_0, w_1 : random!

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



Measure of failure: L_2 / Square error :

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

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

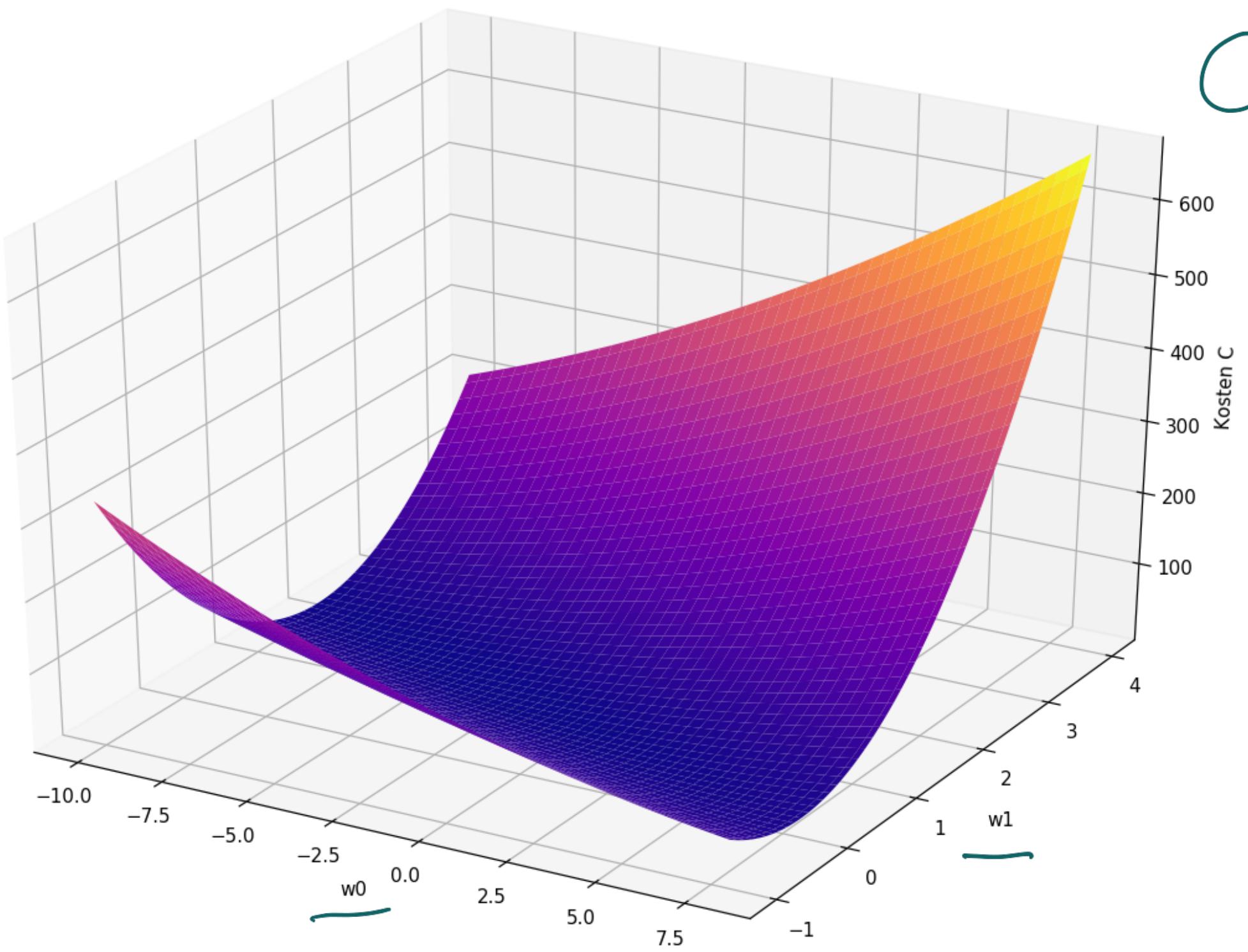
$$\frac{d x^2}{d x} = 2x$$

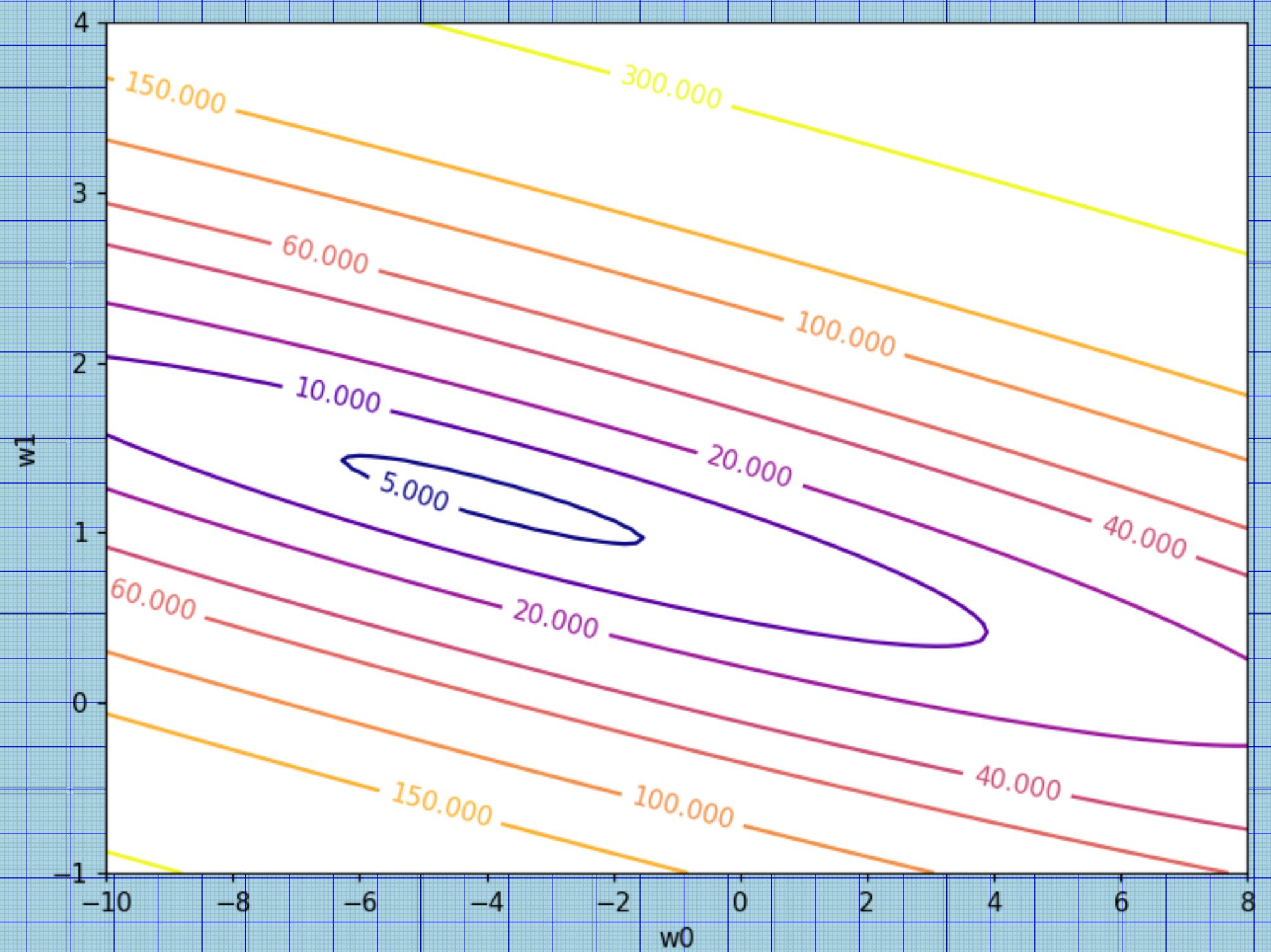
No. Samples

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

Cost function , sampled of θ

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. Θ :
iterative gradient descent:

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

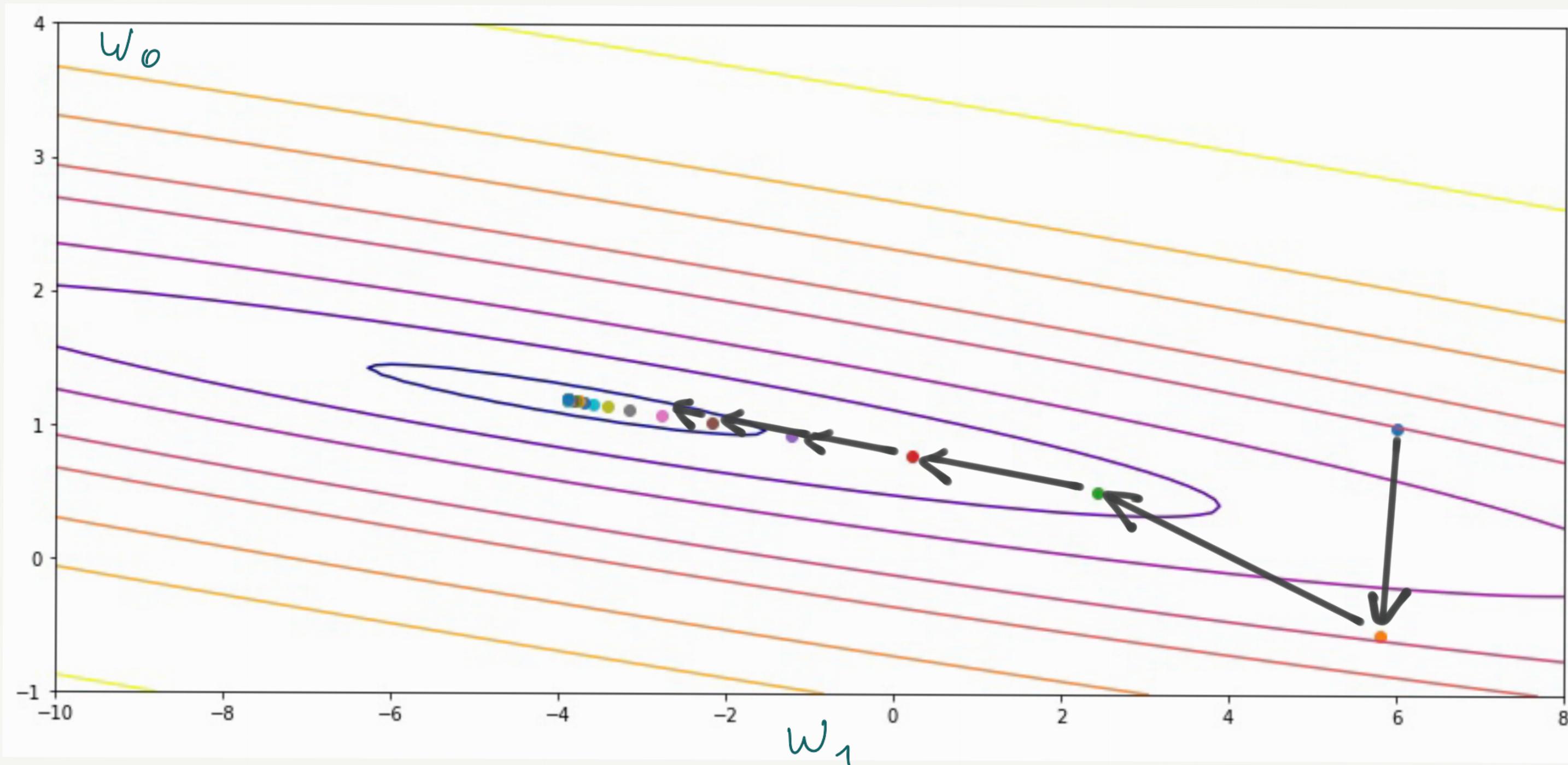
$$w_1^{\text{next}} = w_1^{\text{old}} - \alpha \frac{\partial C}{\partial w_1}$$



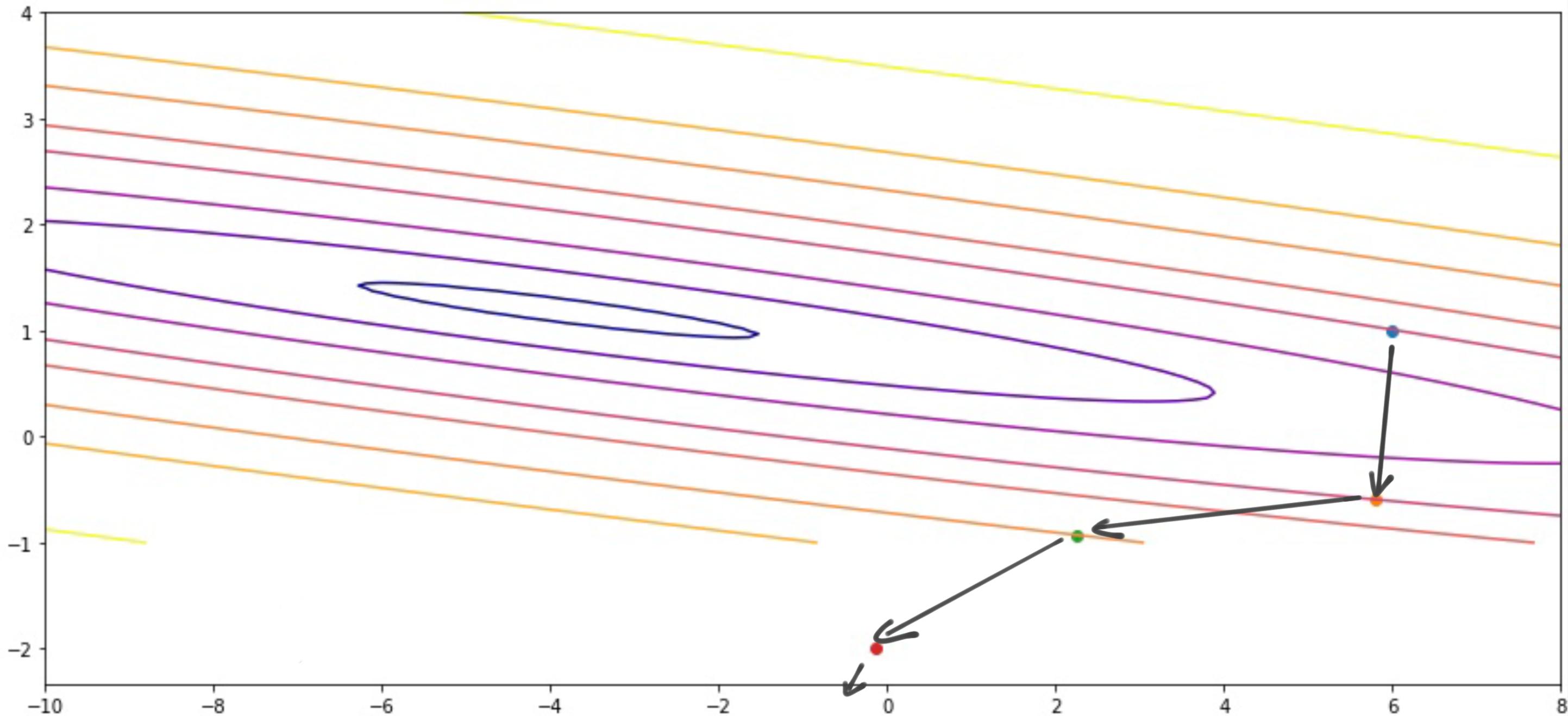
step size: "learning rate"
 \rightarrow hyperparameter

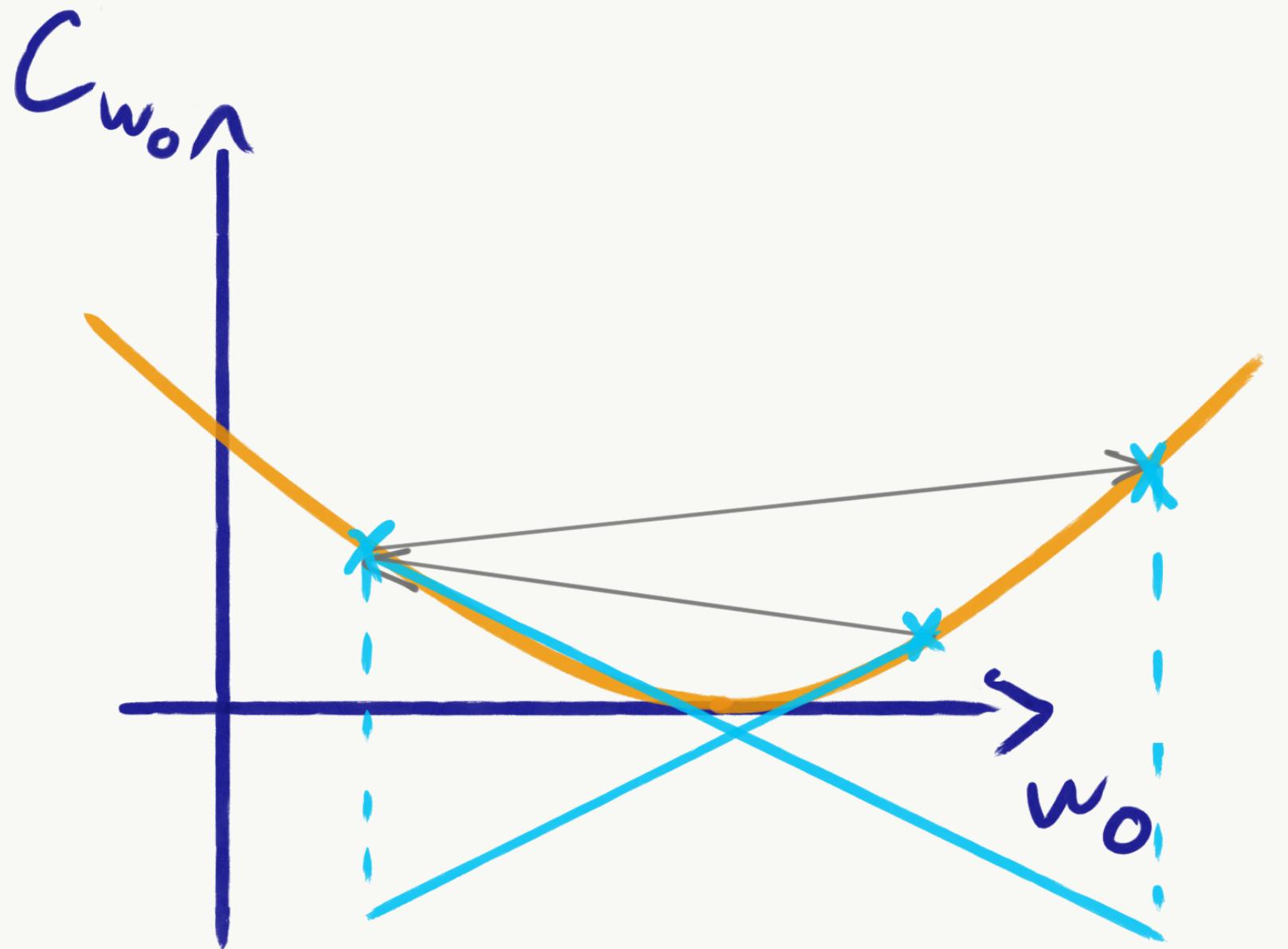
iterate and update!

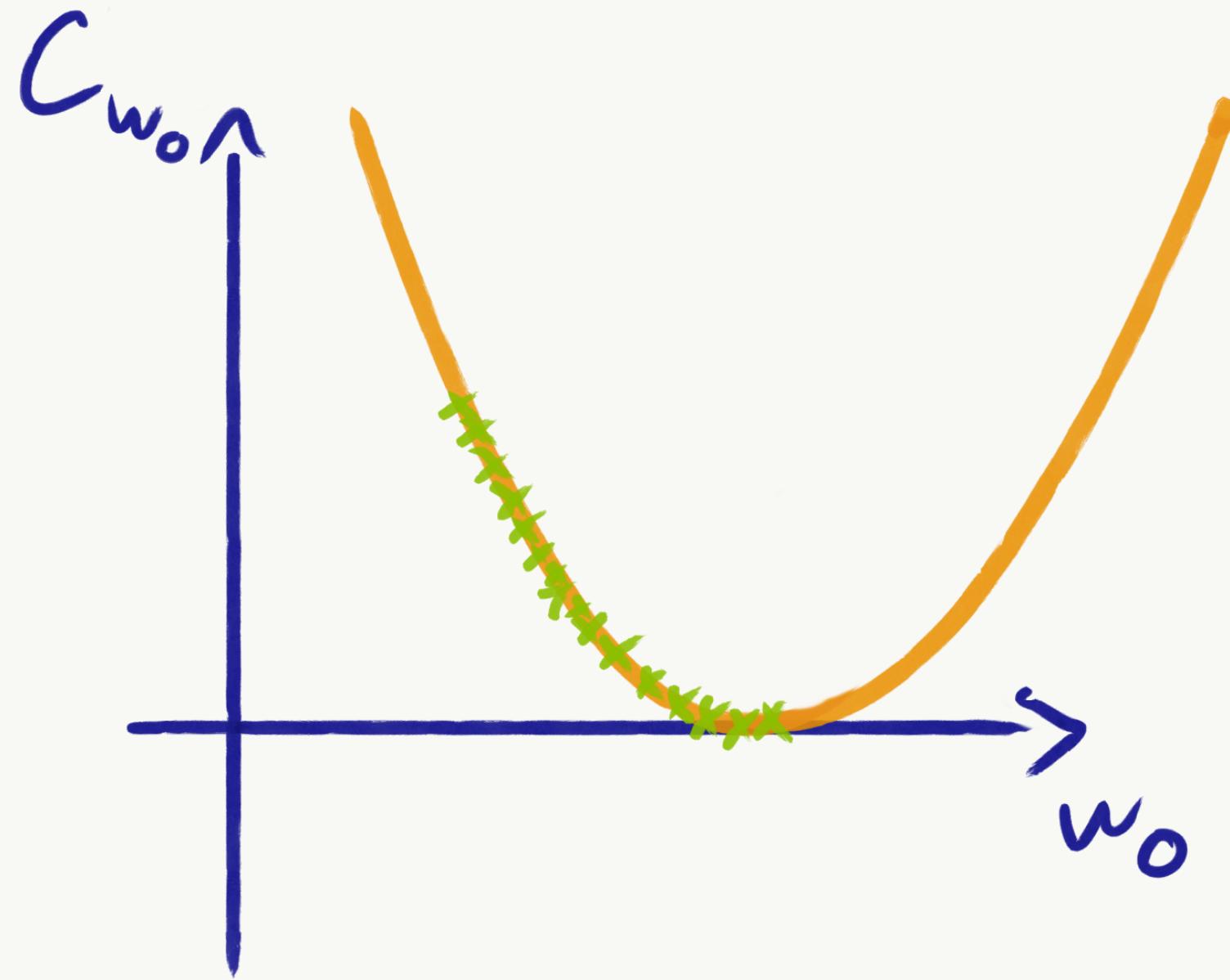
Successful optimization

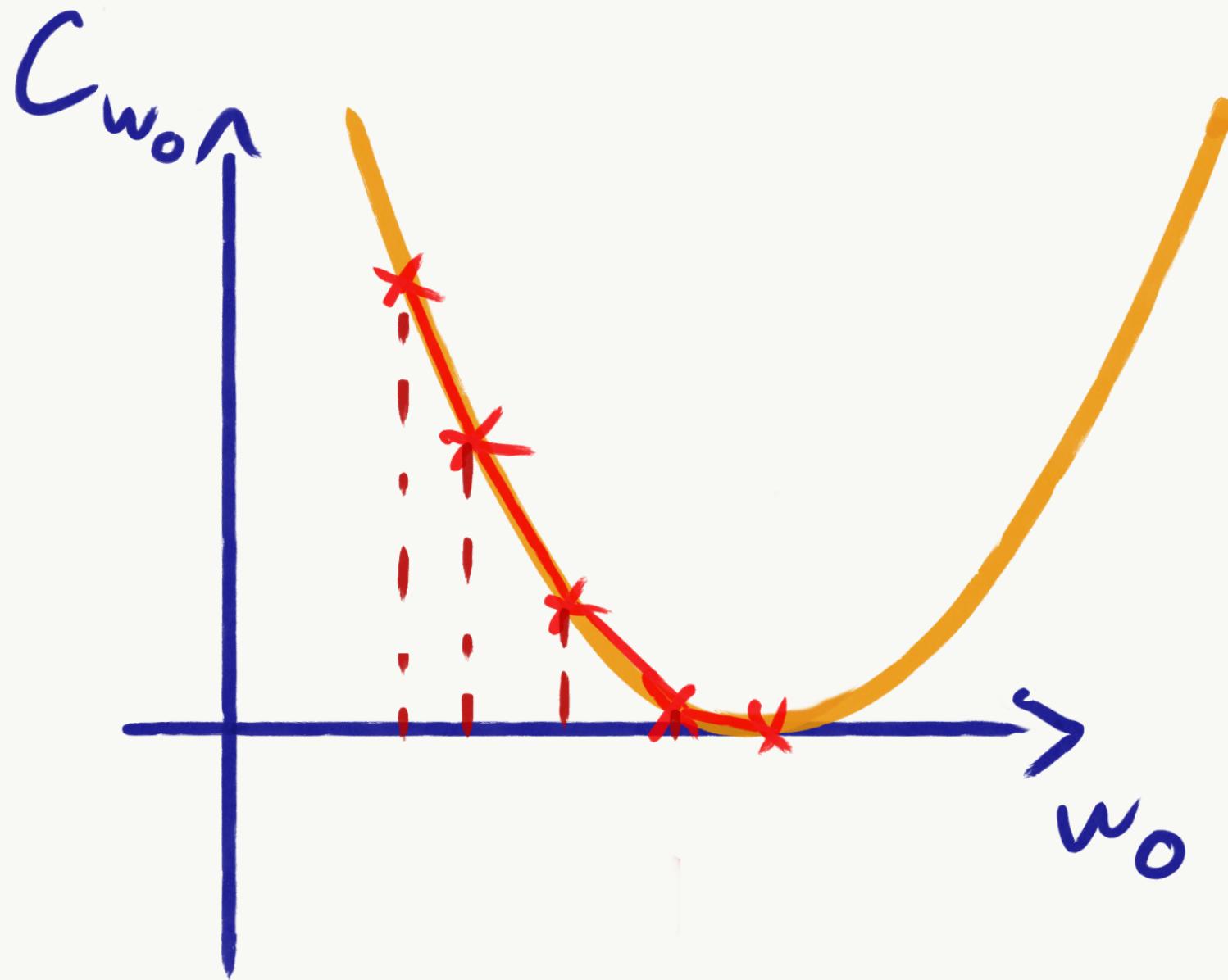


Learning Rate too large \rightarrow diverging

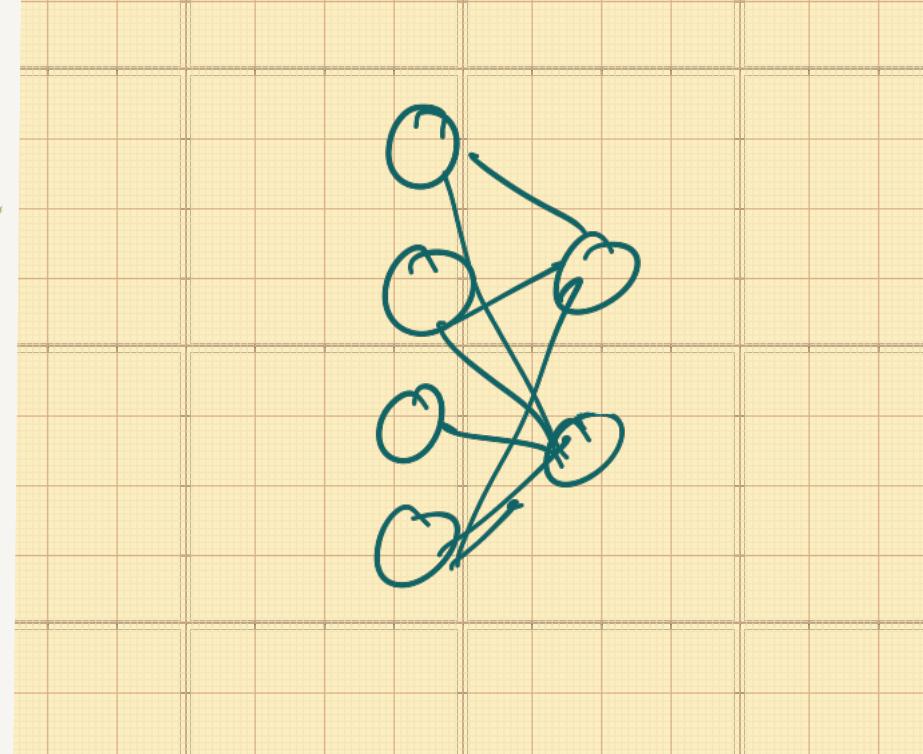
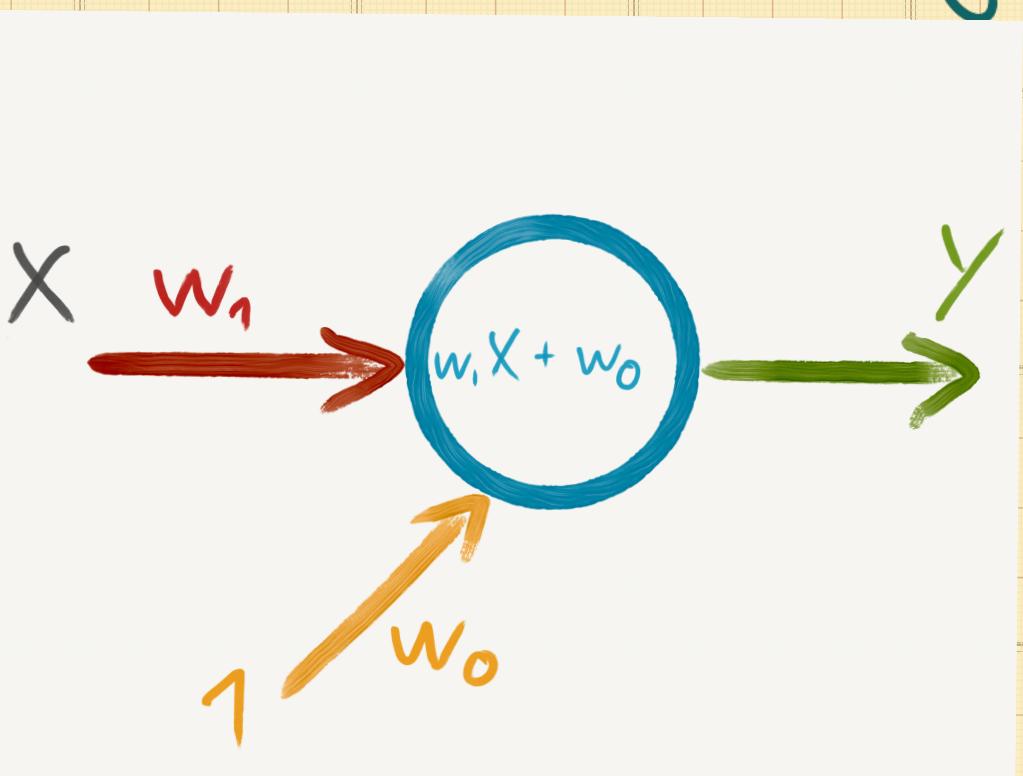








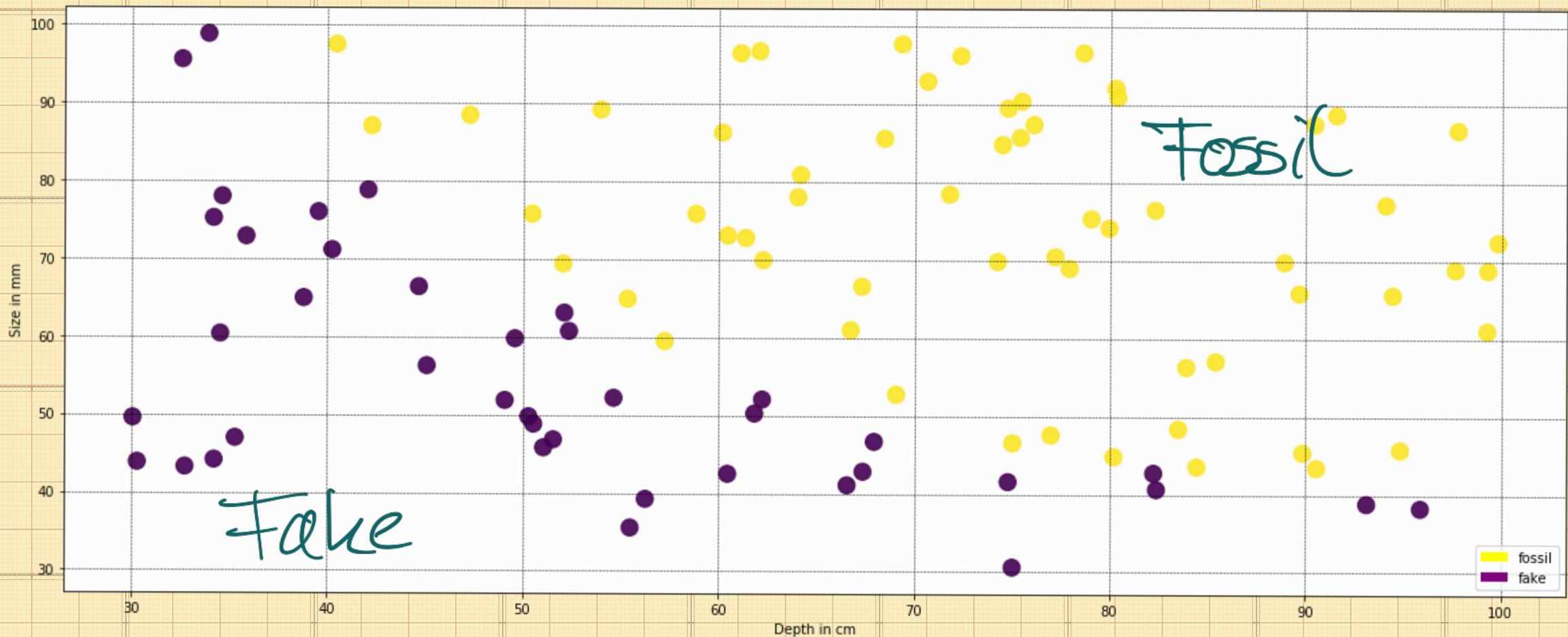
- Note that if done right, this method converges to the LS we discussed before!
- We can increase the expressiveness by
 - more features: $w_0x_0 + w_1x_1 + w_2x_2 + \dots$
 - high order ansatz: $x^2, x^3, x_1 \cdot x_2, \dots$
- Multivariate LS: nothing new!
- Linear Neuron as the building block of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



Logistic Regression (a Classifier!)

lin. Reg : regression ; log. Reg : Classification

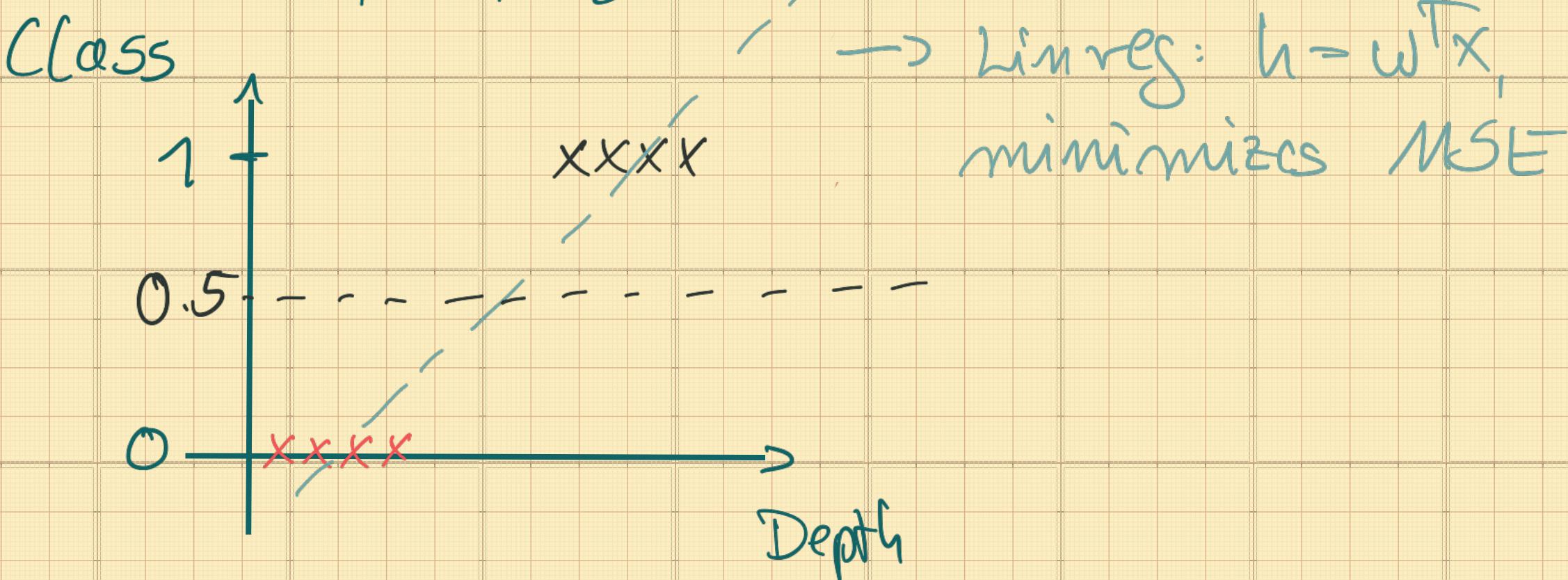
→ We consider binary classification: pass/fail;
spam/no spam; $y = \{0, 1\}$

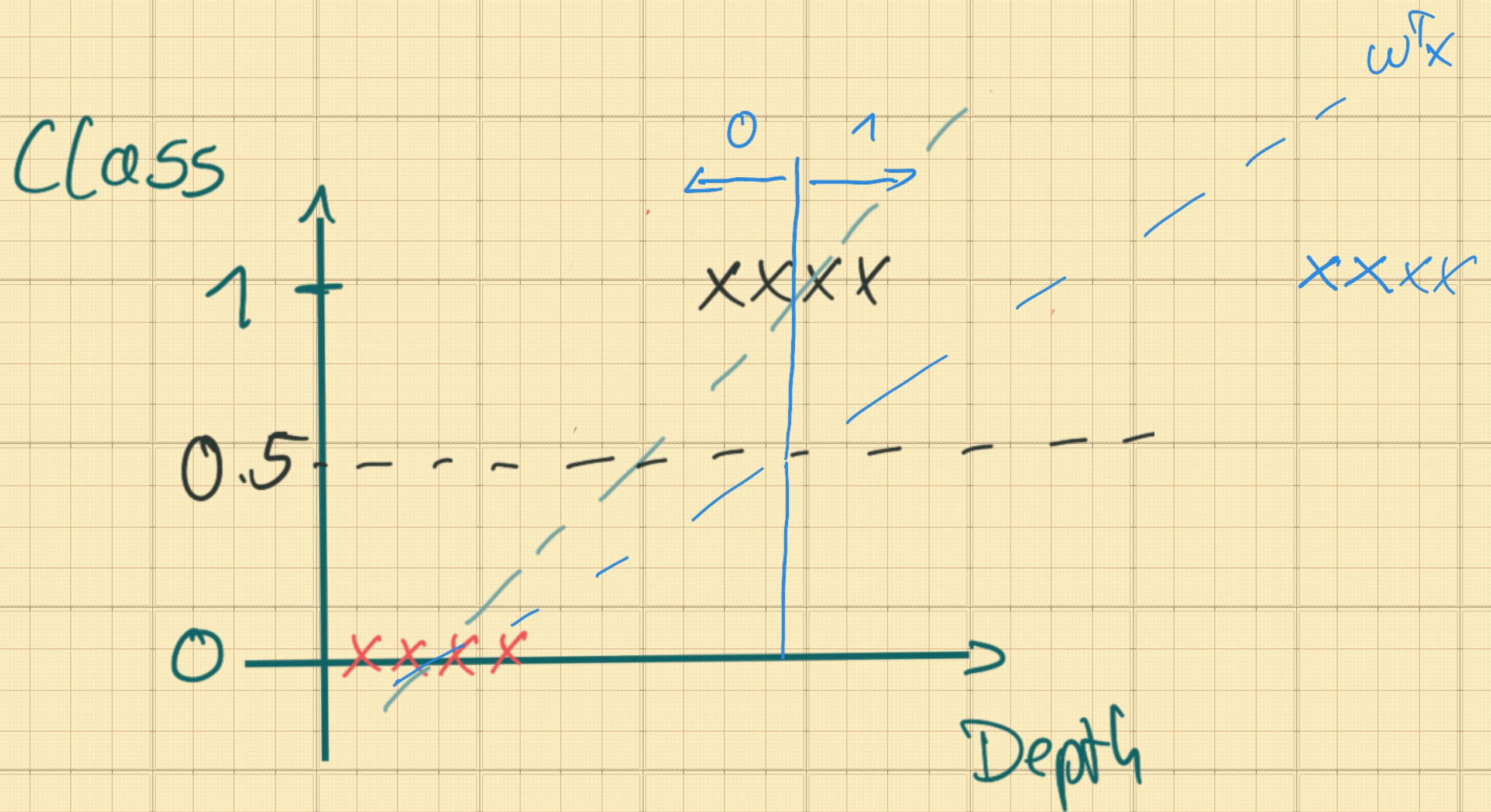


→ train an ML method with inputs x_1 : depth,
 x_2 : size to predict 0/1 : fake / fossil

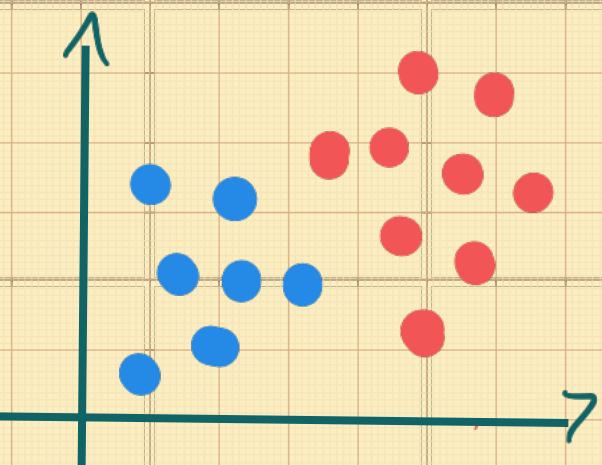
→ let us try the linear regression:

$$h_{\theta} = w_0 + w_1 x_1 + w_2 x_2$$

$$\Theta = (w_0, w_1, w_2)^T$$




- linear model not well suited for classification!
- outliers move the optimum strongly



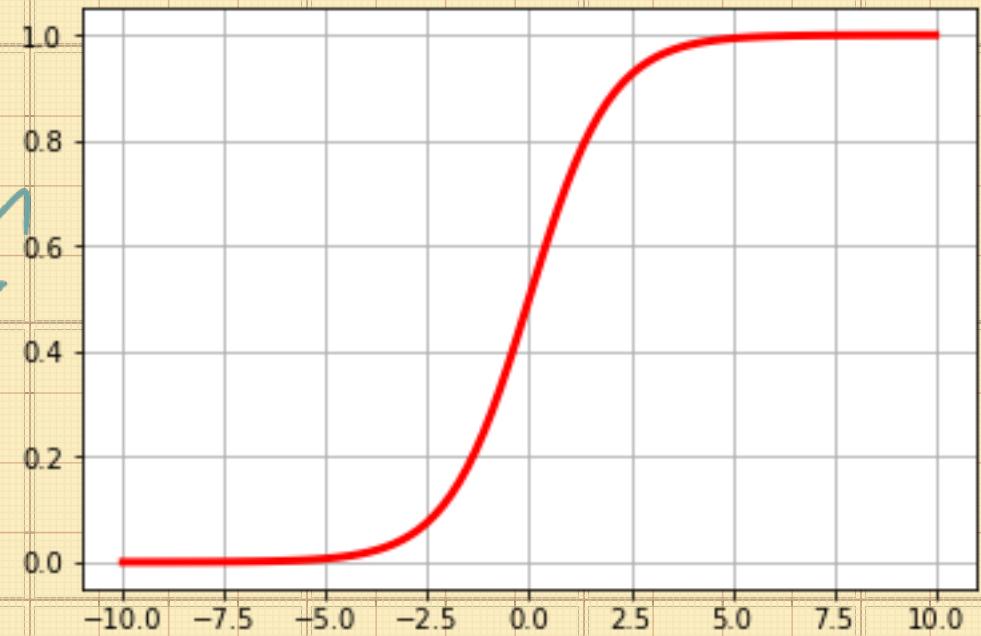
- output of $x^T w$ is unbounded, not $0 \leq x^T w \leq 1$
- Probabilistic P.O.V : $p(y | x, \theta) = \text{Ber}(y | \mu(x))$
- MLE → cost function: "cross entropy"
- non-linear w.r.t. θ .
- optimization more difficult
- practical approach:

→ extend linear hypothesis by
non-linear activation function

→ "sigmoid" function:

↓
 induced by

Bernoulli distr.



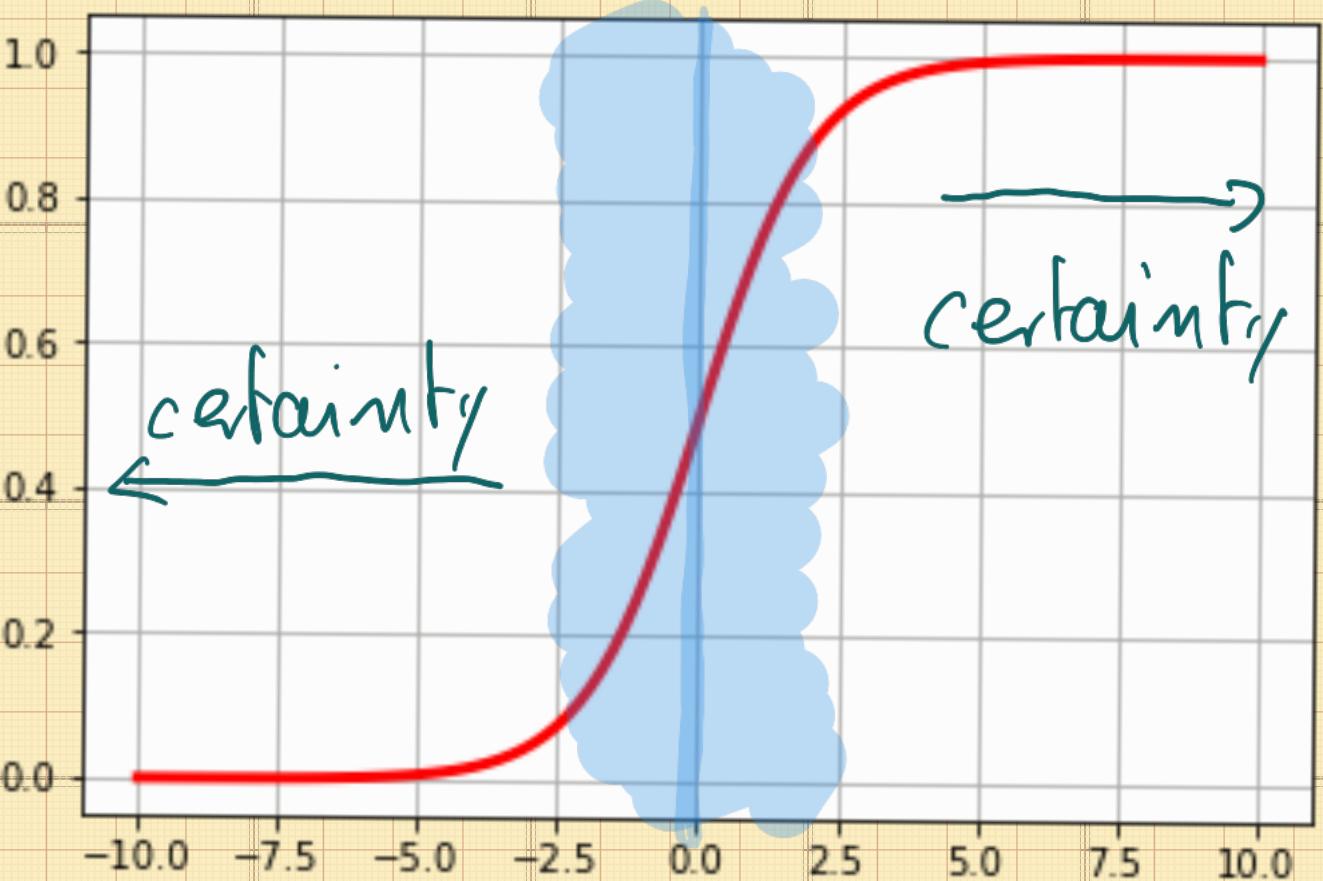
$$\hat{G}(z) = \frac{1}{1 + \exp(-z)}$$

$$z \rightarrow \infty : \hat{G}(z) \rightarrow 1$$

$$z \rightarrow -\infty : \hat{G}(z) \rightarrow 0$$

$$\rightarrow h_{\theta} = \hat{G}(h_{\theta}(x)) = \hat{G}(w_0 + w_1 x_1 + w_2 x_2) \epsilon_{(0,1)}$$

→ We can now interpret $h_{\theta}^{\log}(x)$ as the probabilities
 $P(y=1|x)$: "How likely is class 1?"

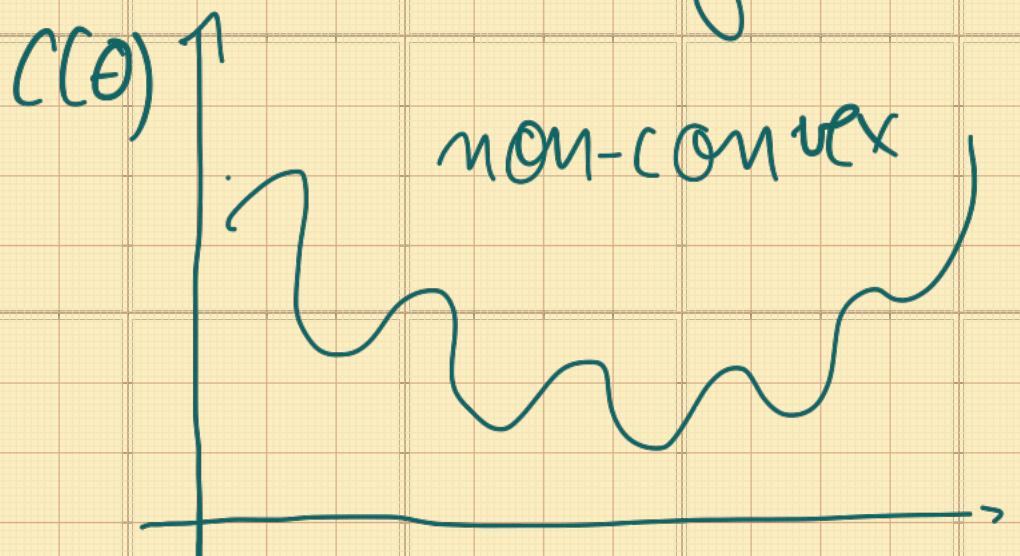


→ note that now all points are equally important
in the cost function (outliers do not skew the
cost)

$h_{\theta}^{\text{log}}(x) = 0.5$: decision boundary

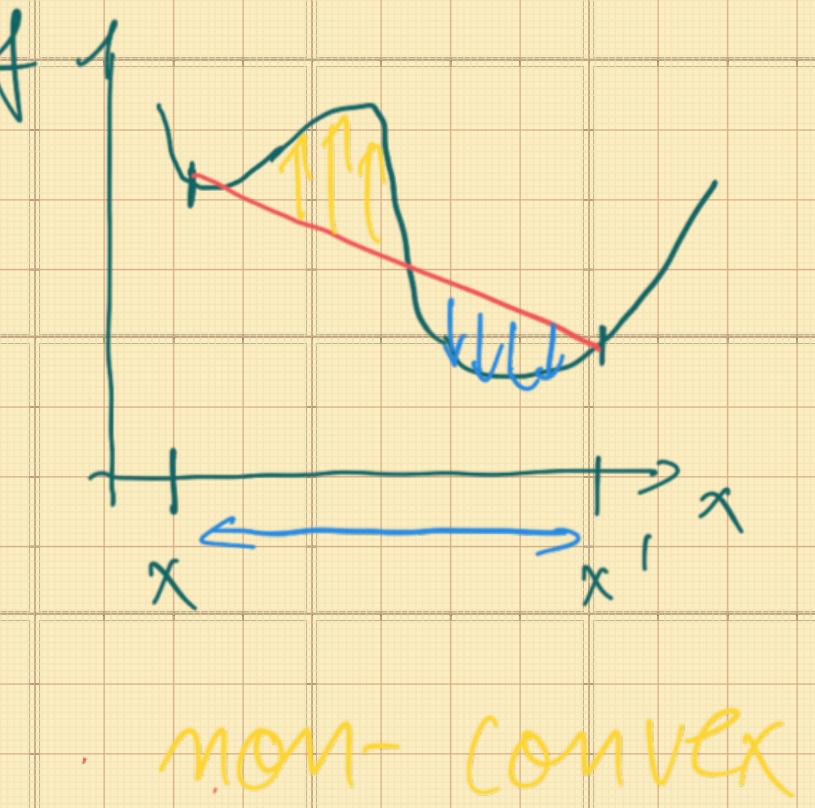
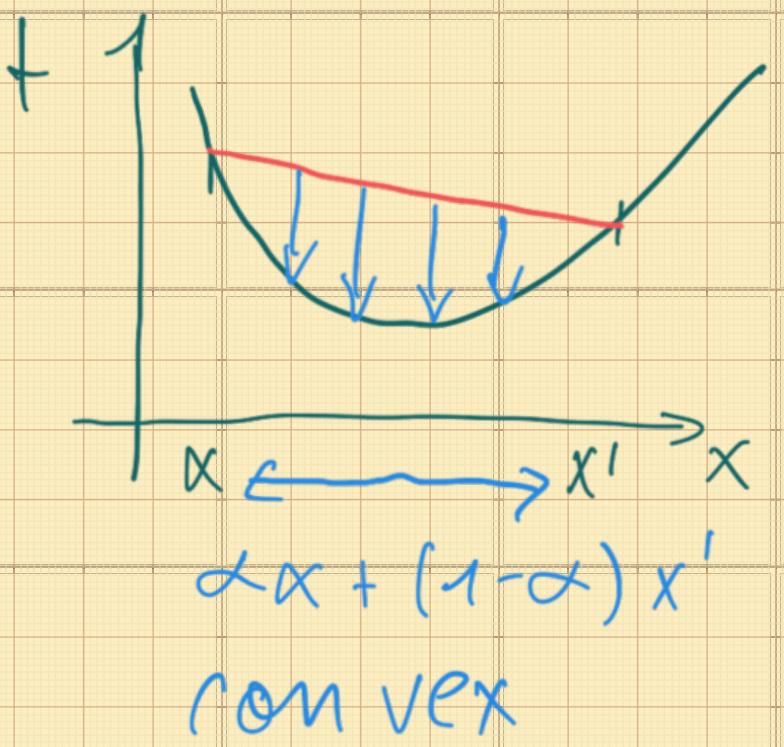
$$C = \frac{-1}{N} \sum_{n=1}^N (y_i \log(h_{\theta}(x_i)) + (1-y_i) \log(1-h_{\theta}(x_i)))$$

(a MSE cost function is not convex w. r. t. θ)

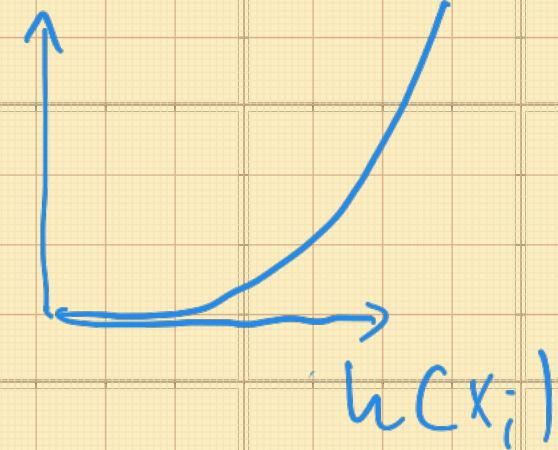
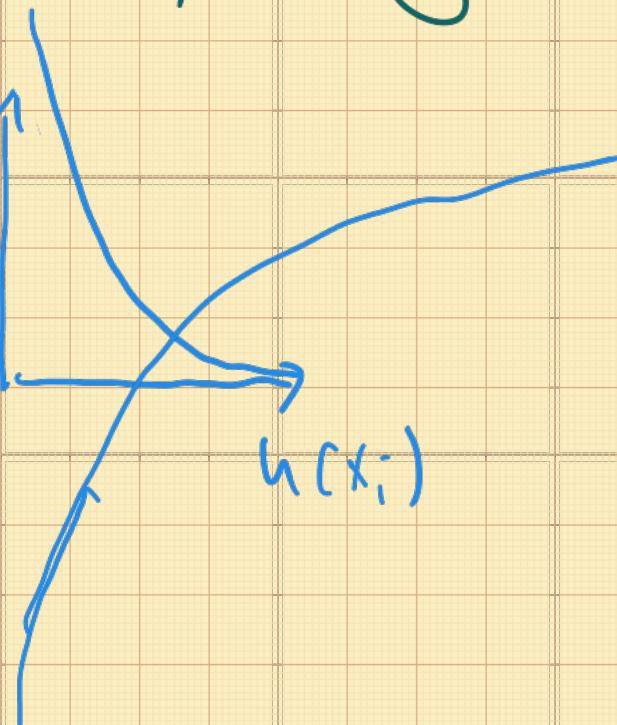


Def.: $0 \leq \alpha \leq 1$: if $f(\alpha x + (1-\alpha)x') \leq \underbrace{\alpha f(x) + (1-\alpha)f(x')}_{\text{linear blending}}$

linear blending



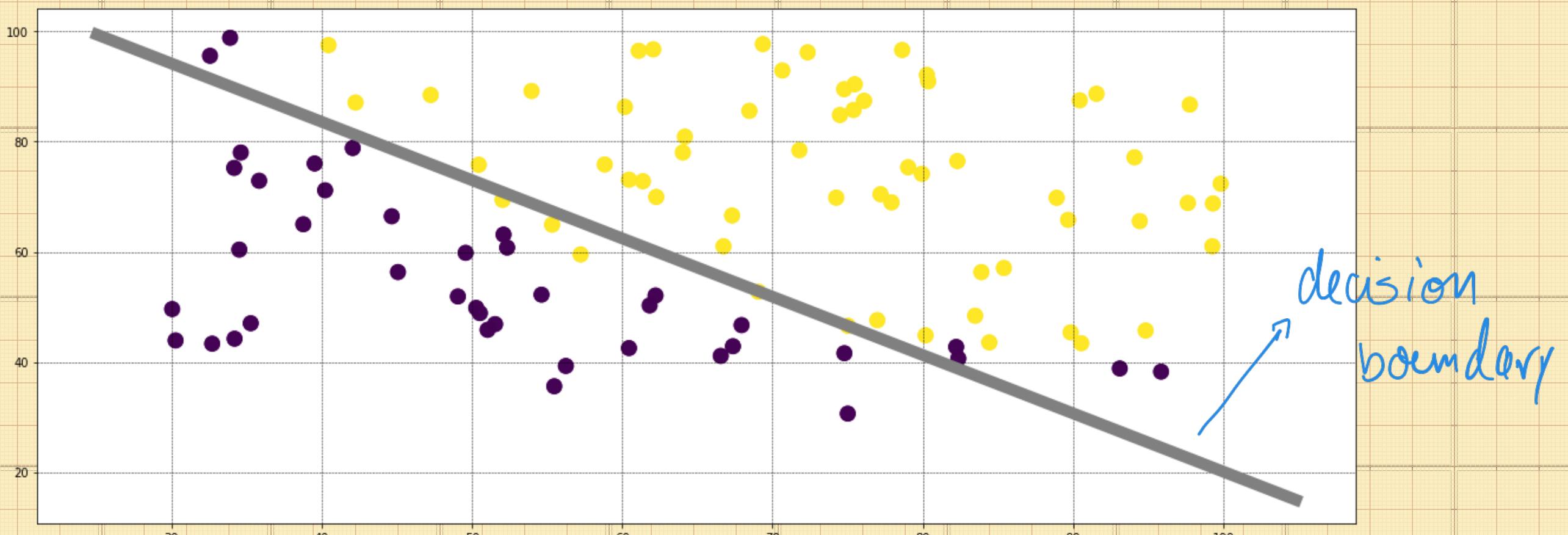
$$C^i = -y_i \log(h_i) - (1-y_i) \log(1-h_i)$$



$$\nabla_{\theta} C^i = \frac{\partial C^i}{\partial w_j} = (y_i - \hat{y}_i) x_{ij}$$

j: feature index
 $(w_0x_0 + w_1x_1 + w_2x_2)$
i: which sample

$$\text{e.g. } w_1^{\text{new}} = w_1^{\text{old}} - \frac{LR}{N} \sum_{n=1}^N \frac{\partial C^n}{\partial w_1}$$



Problem : linearly separable data

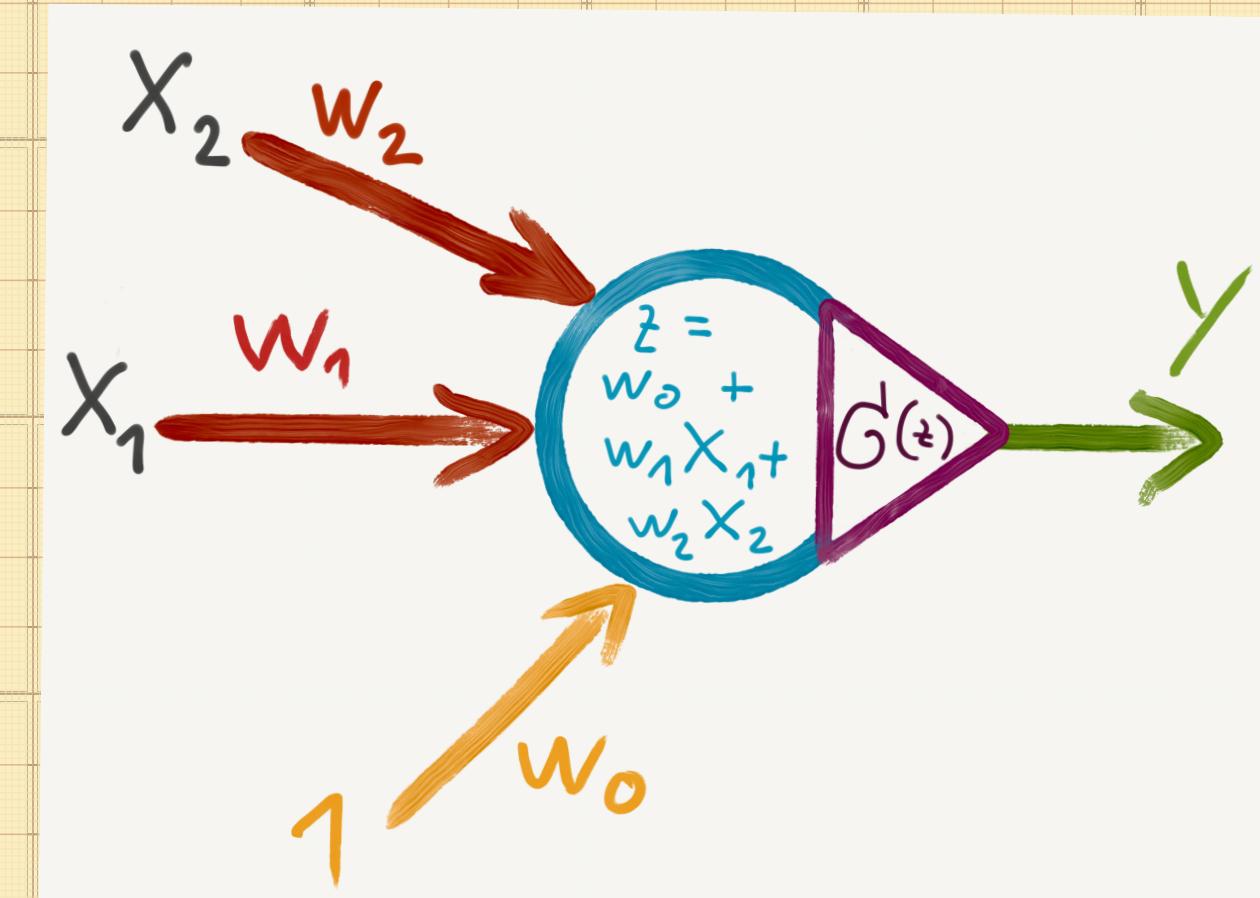
Logistic neuron :
 linear neuron
 → activation function

Sigmoid, tanh, ReLU, ...

"Perceptron": Rosenblatt ~~1956~~
 1962

→ Book recommendation:

Pattern recognition & machine learning, Ch. Bishop
 Springer



$$\hat{G}(z) = \begin{cases} 1 & z \geq 0 \\ -1 & z < 0 \end{cases}$$