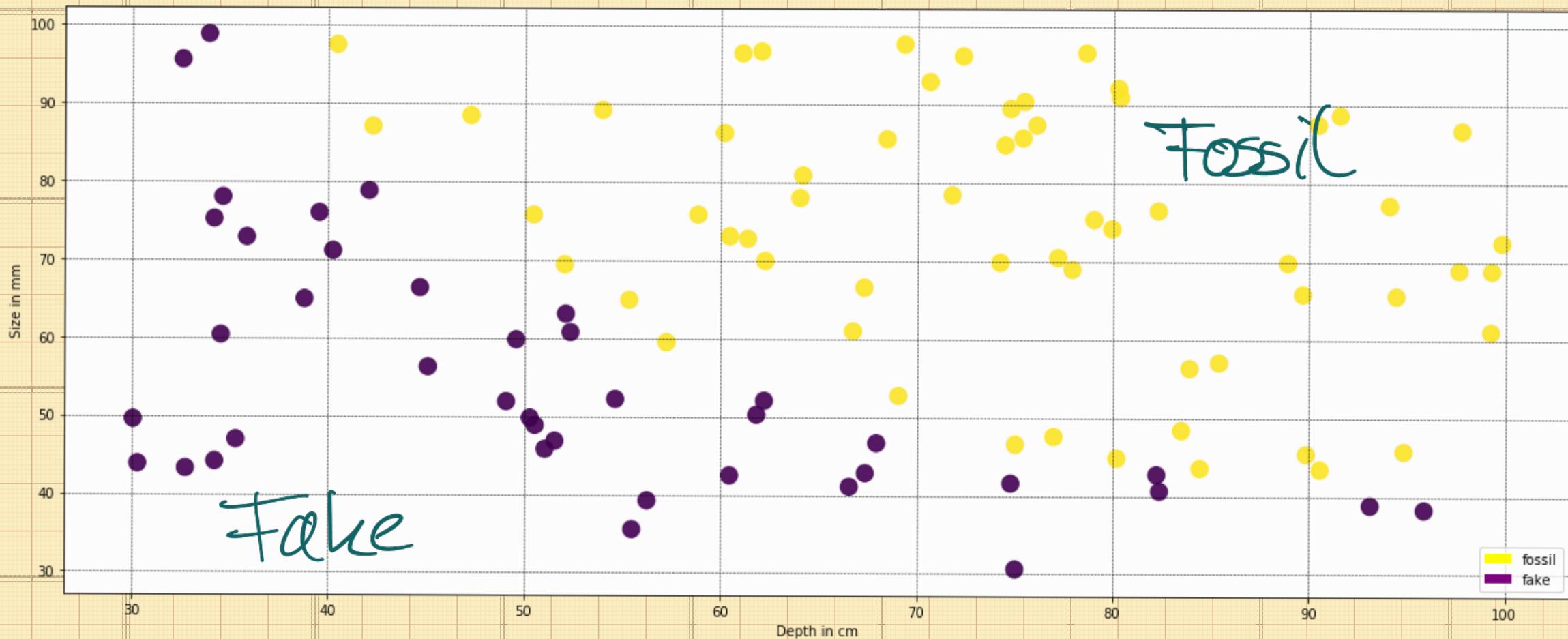


Logistic Regression (a Classifier!)

lin. Reg : regression ; log. Reg : Classification

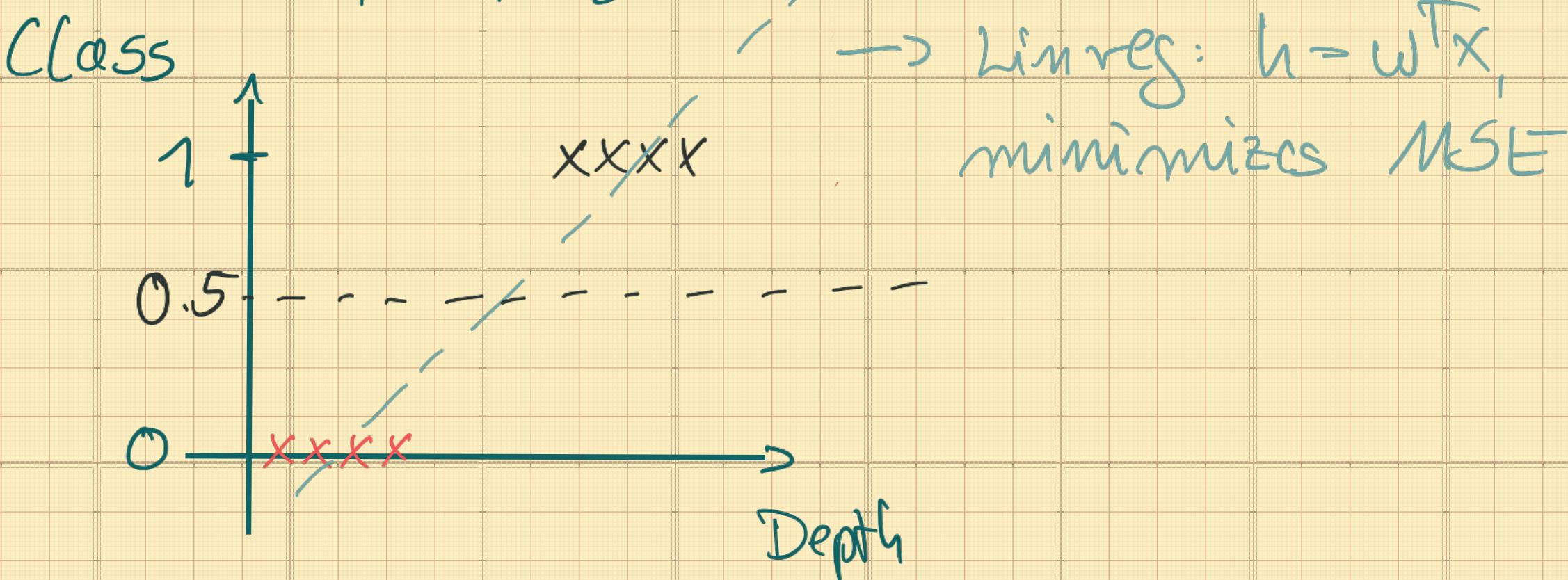
→ We consider binary classification: pass/fail;
spam/no spam; $y = \{0, 1\}$

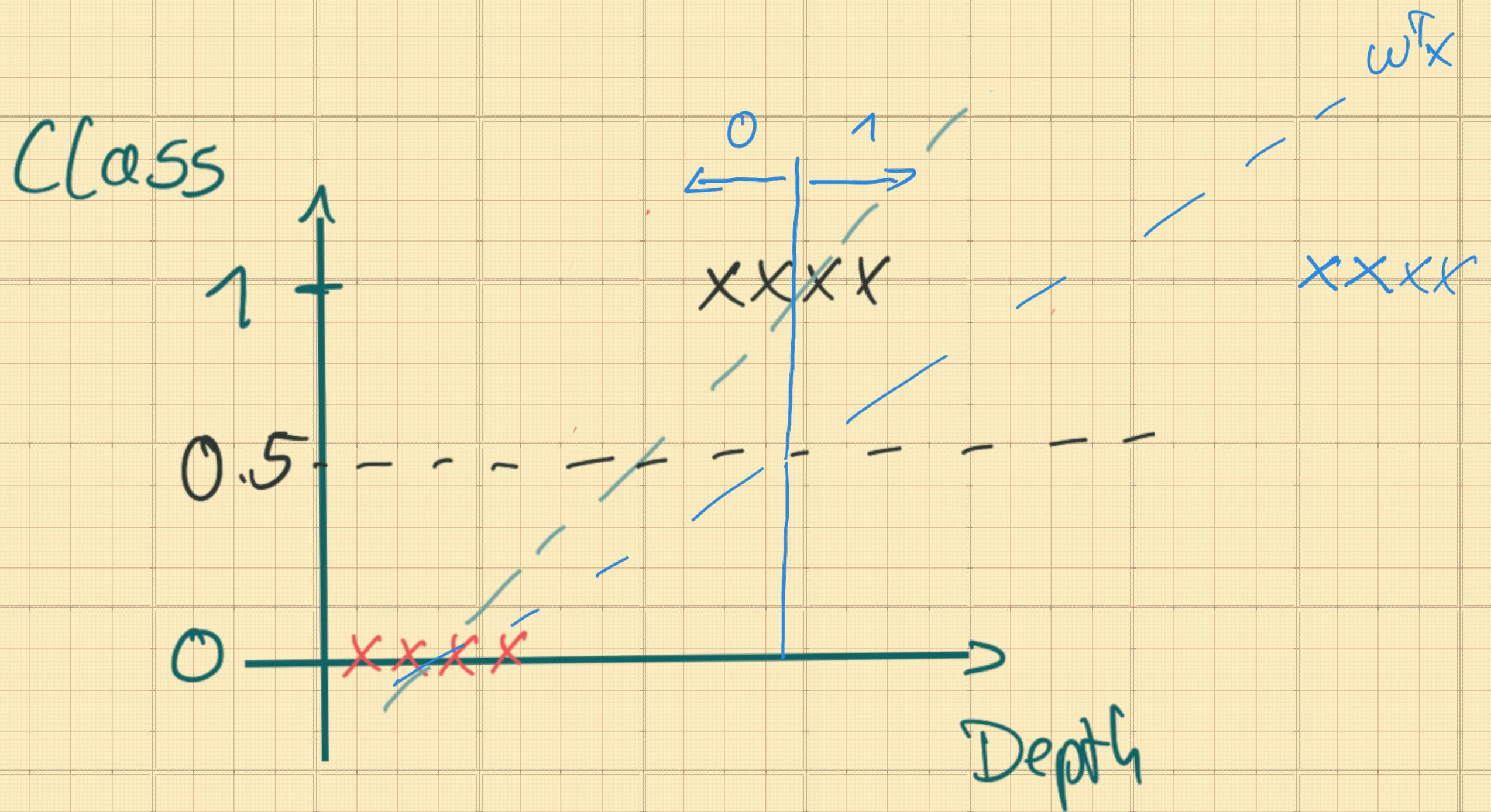


→ train an ML method with inputs x_1 : depth,
 x_2 : size to predict 0/1 : fake / fossil

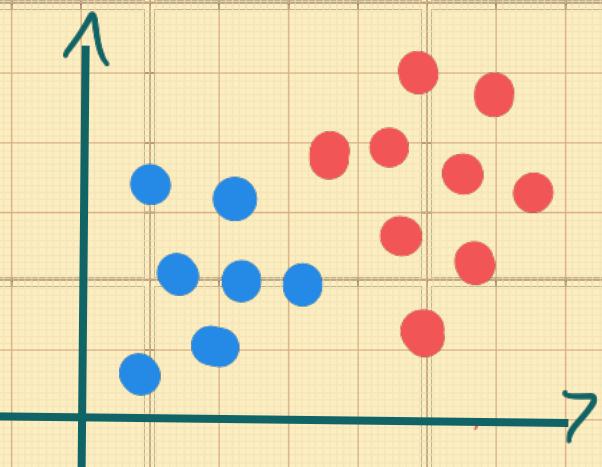
→ let us try the linear regression:

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

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




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



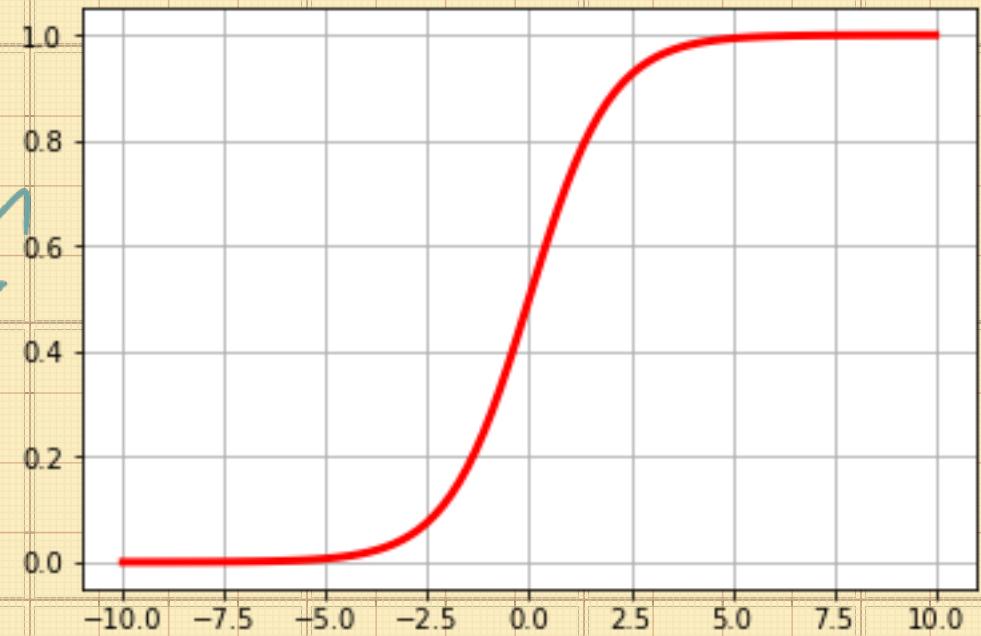
- output of $x^T w$ is unbounded, not $0 \leq x^T w \leq 1$
- Probabilistic P.O.V : $p(y | x, \theta) = \text{Ber}(y | \mu(x))$
- MLE → cost function: "cross entropy"
- non-linear w.r.t. θ
- optimization more difficult
- practical approach:

→ extend linear hypothesis by
non-linear activation function

→ "sigmoid" function:

↓
 induced by

Bernoulli distr.



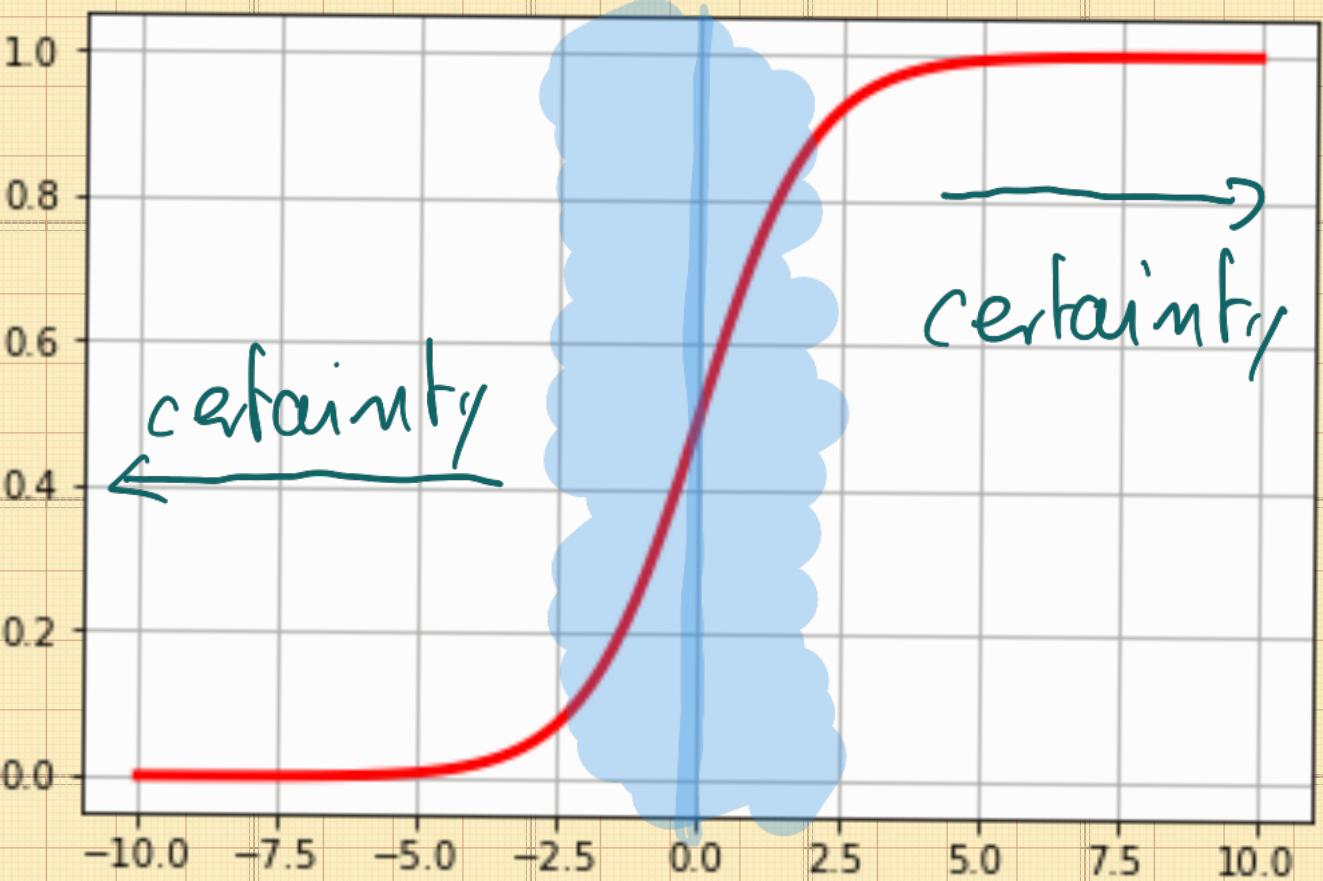
$$\hat{G}(z) = \frac{1}{1 + \exp(-z)}$$

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

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

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

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

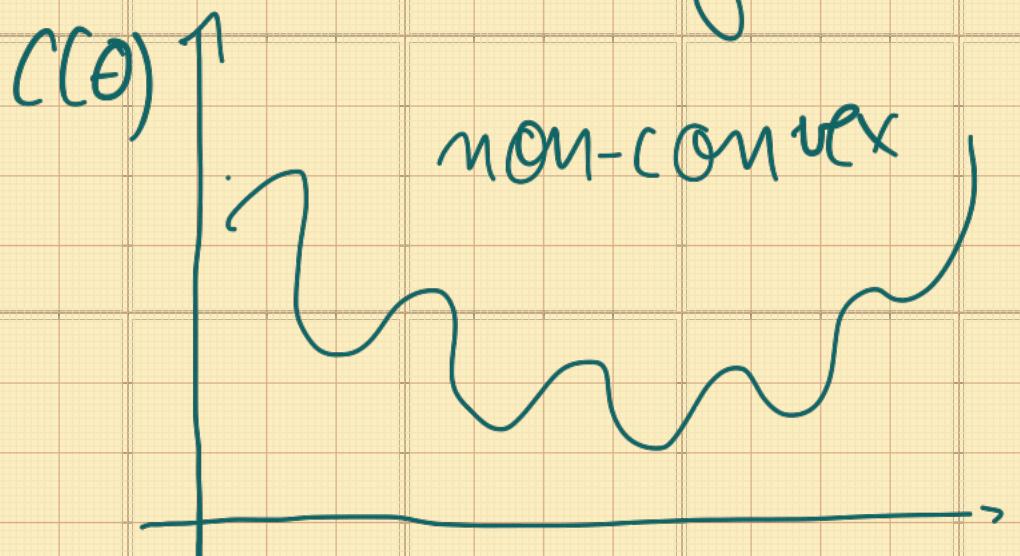


→ note that now all points are equally important
in the cost function (outliers do not skew the
cost)

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

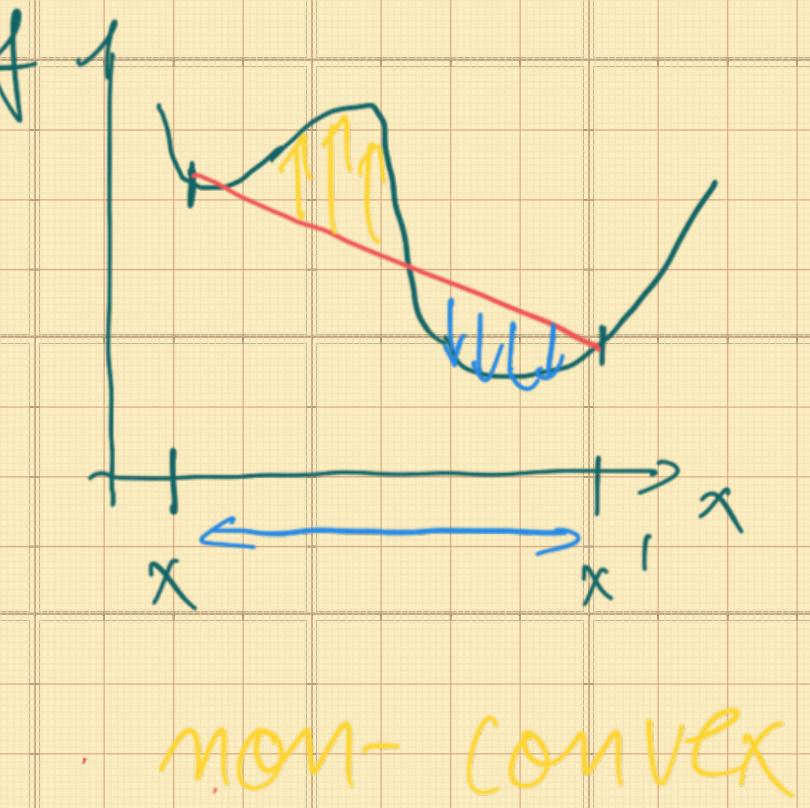
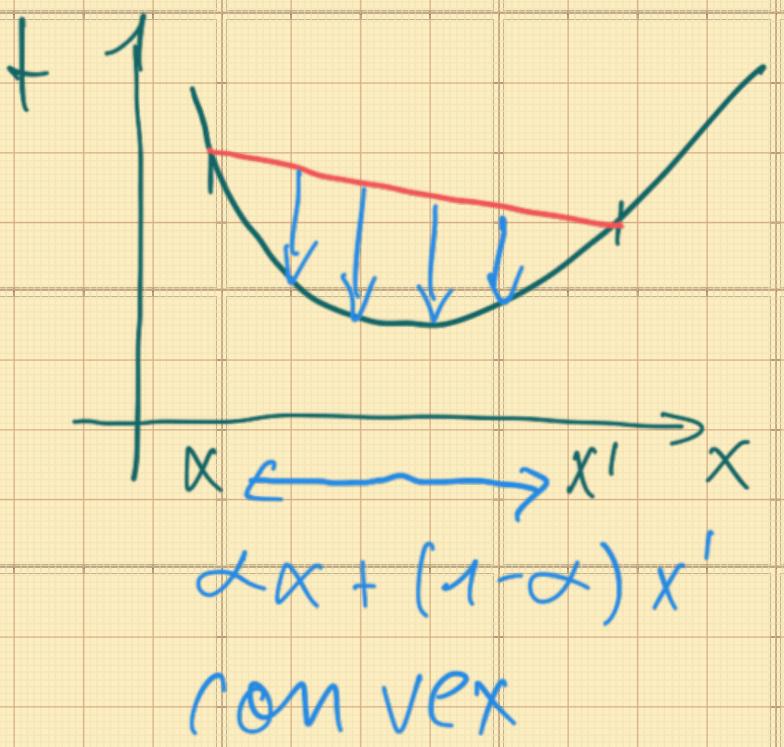
$$C = \frac{-1}{N} \sum_{n=1}^N (y_i \log(h_{\theta}(x_i)) + (1-y_i) \log(1-h_{\theta}(x_i)))$$

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

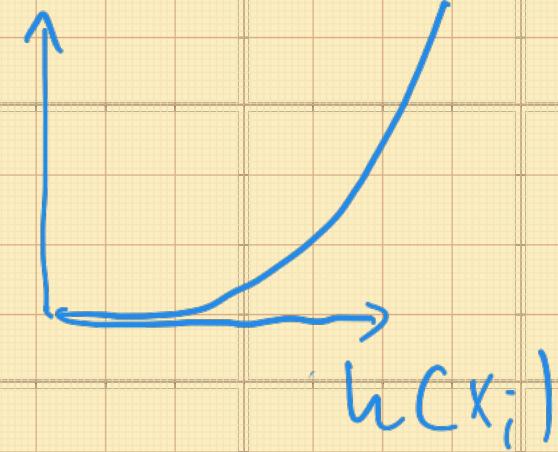
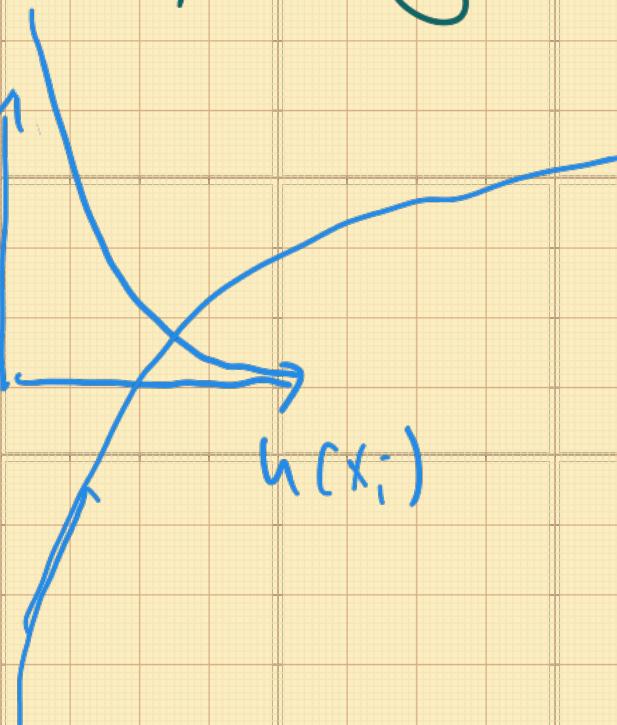


Def.: $0 \leq \alpha \leq 1$: if $f(\alpha x + (1-\alpha)x') \leq \underbrace{\alpha f(x) + (1-\alpha)f(x')}_{\text{linear blending}}$

linear blending



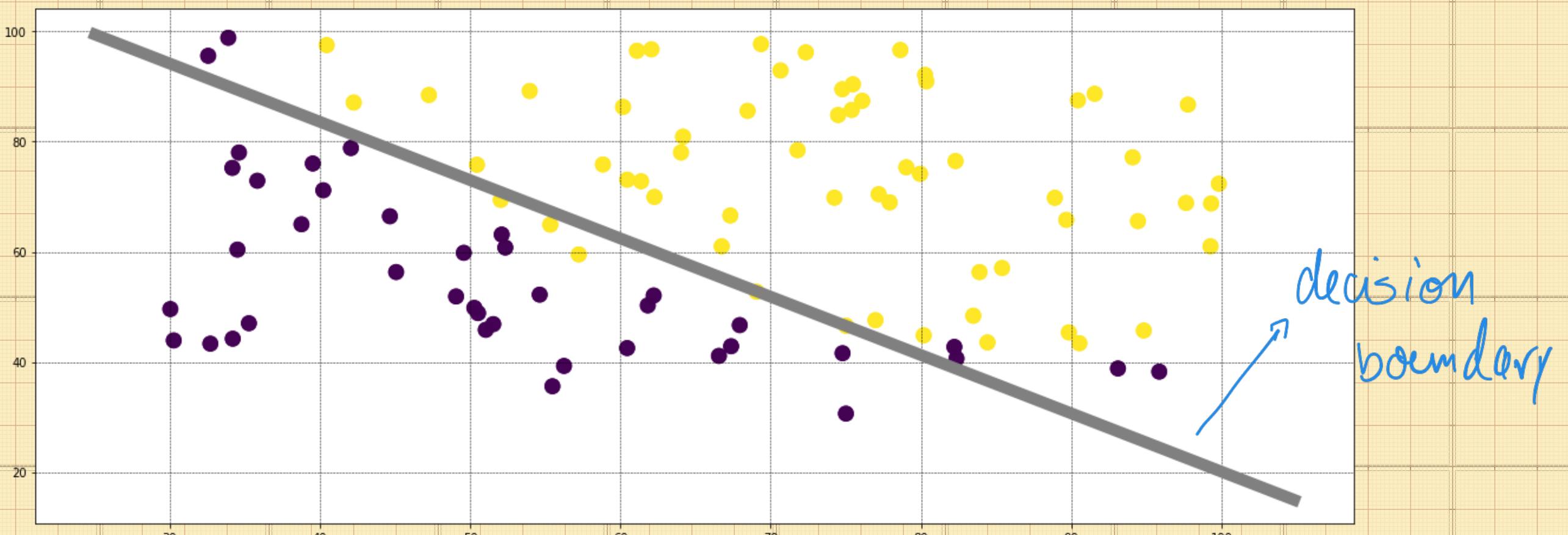
$$C^i = -y_i \log(h_i) - (1-y_i) \log(1-h_i)$$



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

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

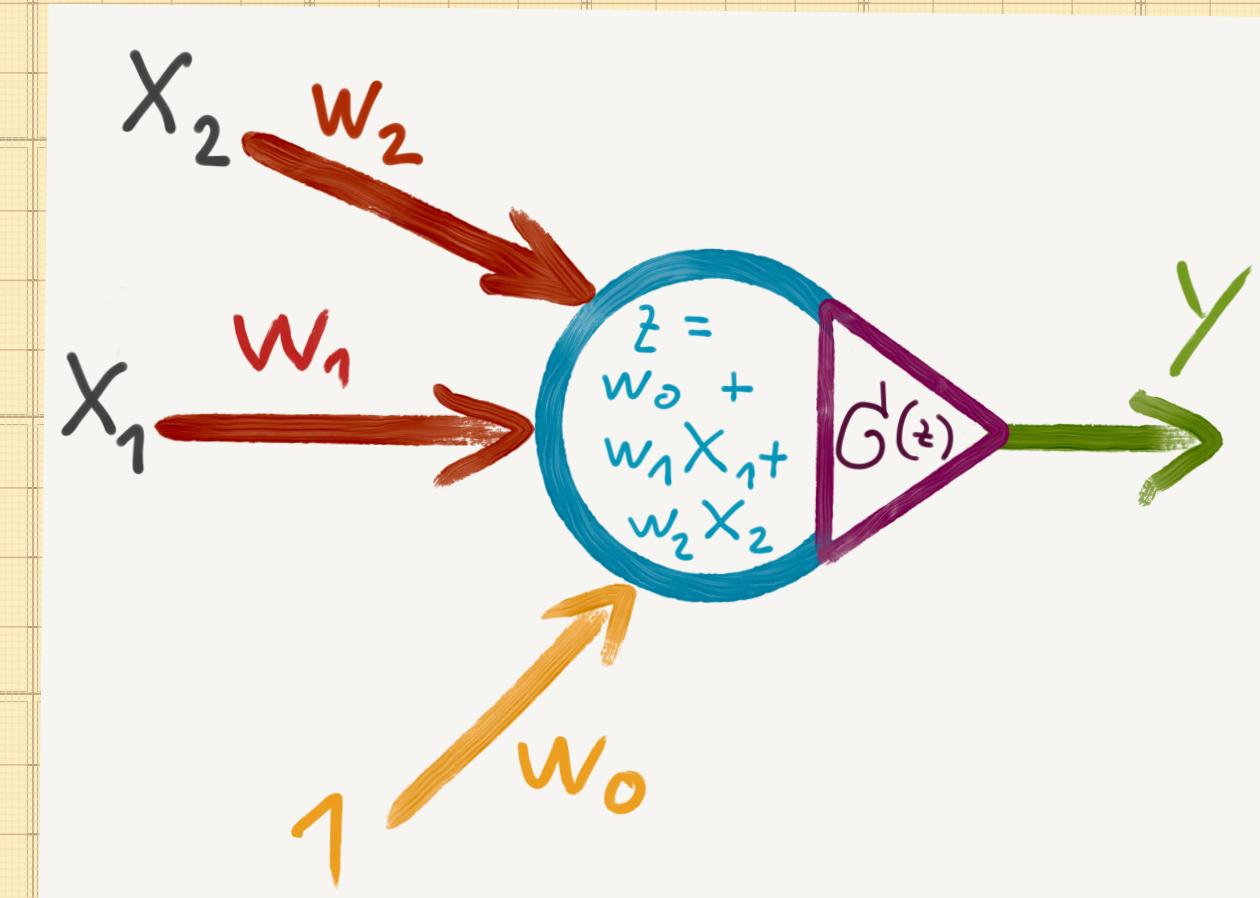
Logistic neuron :
 linear neuron
 → activation function

Sigmoid, tanh, ReLU, ...

"Perceptron": Rosenblatt ~~1956~~
 1962

→ Book recommendation:

Pattern recognition & machine learning, Ch. Bishop
 Springer



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