

## What is ML?

• setting: - quantities  $Y$  we would like to know, but cannot (easily) measure  
                  response

                  features  $X$  we can measure and are related to  $Y$   
- quantities

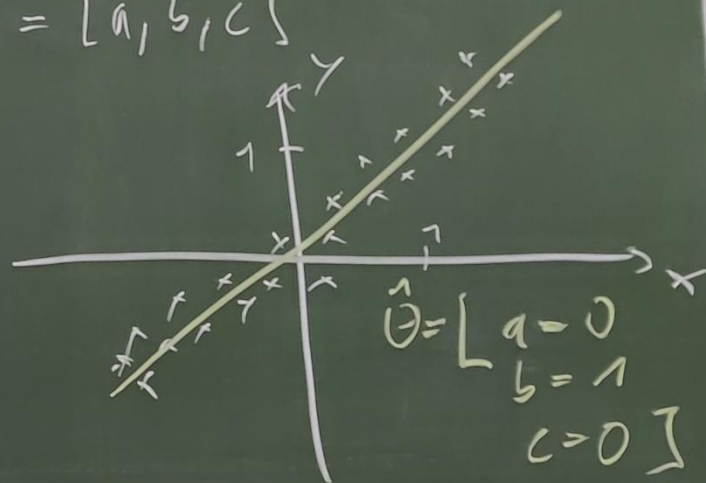
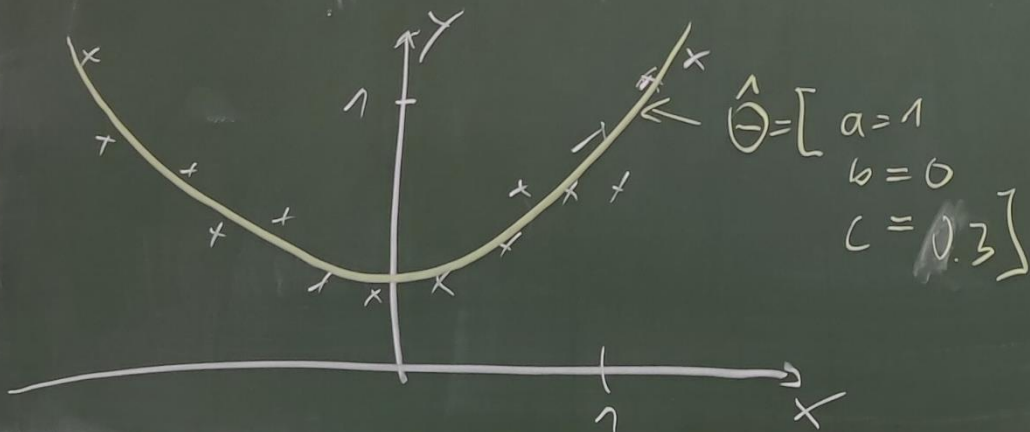
• idea: learn function  $\hat{Y} = f(X)$  such that  $\hat{Y}_i \approx Y_i^*$  truth

• traditional science: ask expert to design  $f(X)$

• machine learning: - choose a "universal function family"  $\mathcal{F}$   
                          -  $\mathcal{F}$  has parameters  $\Theta \Rightarrow$  find parameter  $\hat{\Theta}$  that fit data

$$\hat{Y} = f_{\hat{\Theta}}(X) \quad \text{s.t.} \quad \hat{Y}_i \approx Y_i^*$$

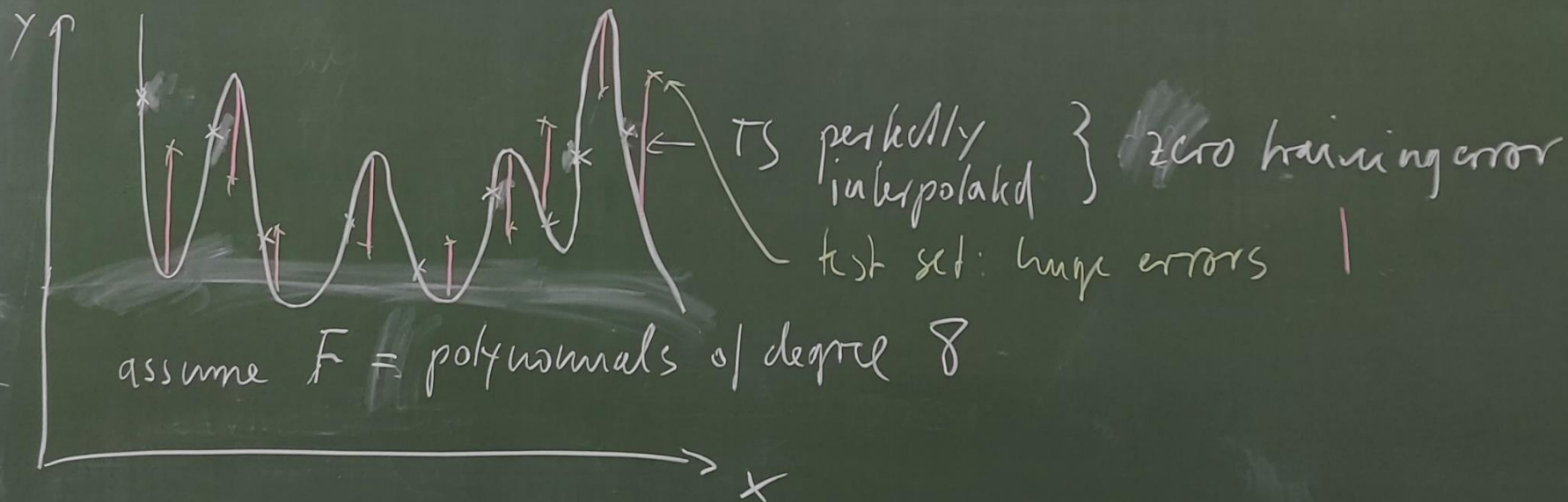
• example:  $\mathcal{F} = \{ax^2 + bx + c\}$        $\Theta = [a, b, c]$



## Basic ML workflow

- ① collect data: two datasets: - training dataset (TS) to find  $\hat{\theta}$   
- test dataset to validate quality
- ② select function family  $F$  (prior knowledge, experience, trial-and-error)  
(upstream: tools)
- ③ find best  $\hat{\theta}$  by fitting  $F$  to training set
- ④ validate quality of  $f_{\hat{\theta}}(x)$  on test set
- ⑤ deploy model in practice

separate training & test sets necessary to recognize Overfitting





## probabilistic prediction:

- reality is often more complicated: no unique true response  $y^*$   
instead: set of possible values  $\{y_1, y_2, \dots\}$   
 $\Rightarrow$  learn a conditional probability (instead of function)  $\left( \text{e.g. } y^* \in [1, 2, 3] \right)$

$$\hat{y} \sim p(y|x)$$

- define "universal probability family" with parameters  $\Theta$   
choose  $\hat{\Theta}$  such that  $y \sim p_{\hat{\Theta}}(y|x)$  matches the data

- strict generalization: deterministic case is recovered by defining

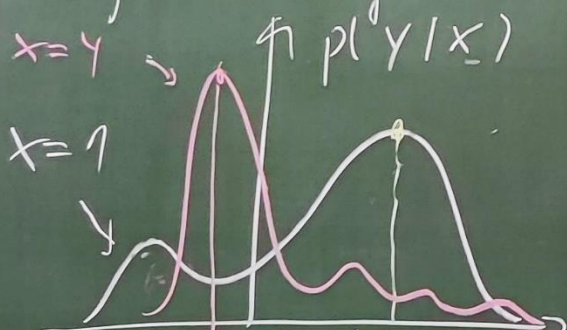
$$p_{\Theta}(y|x) = \delta(y - f_{\Theta}(x))$$

- often, we derive deterministic results from probabilistic predictions

e.g. mode of distribution  $\Rightarrow$  value with highest probability

$x=y$   $\rightarrow$   $p(y|x)$

$x=1$

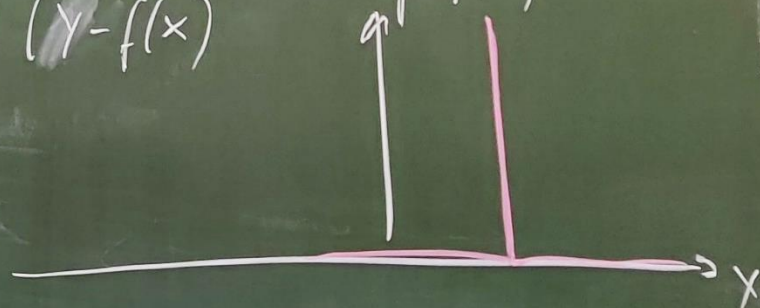


alternatives:

mean  
median  
random

$\Gamma(y - f(x))$

$p(y|x)$



## why is probabilistic prediction important?

- reality is full of uncertainty, on 4 levels
  - fundamental: nature is not fully predictable
    - quantum physics • deterministic chaos • combinatorial explosion
  - data: finite, noisy, incomplete, ambiguous
  - model: solutions  $f_{\hat{\theta}}$ ,  $p_{\hat{\theta}}$  are imperfect (function family too small)  
may not have converged perfectly, nature has changed in the meantime
  - goal: disagreement about what is relevant / desirable?

handling uncertainty is central challenge of AI and ML



Notation:

features:  $X_i$  row vector for each instance  $i$  ( $i = 1, \dots, N$ ), features  $j$  ( $j = 1, \dots, D$ )

feature matrix  $X$  rows instances, columns features,  $X \in \mathbb{R}^{N \times D}$

in Python, all float, discrete labels  $\Rightarrow$  one-hot encoding

$i \backslash j$	1 (name)	2 (height)	3 (gender)
1	Alice	1.7m	f
2	Bob	1.8m	m
3	Max	1.9m	m

$N = 3$

$i \backslash j$	1 (height)	2 (f)	3 (m)	4 (d)
1	1.7	1	0	0
2	1.8	0	1	0
3	1.9	0	1	0

elements  $X_{ij}$

response:  $Y_i$  row vector for instance  $i$ , response elements  $m$  ( $m = 1, \dots, M$ )

mostly:  $M = 1$ , scalar response (exception: one-hot encoding of discrete labels)

tasks according to type of response

(1)  $Y_i \in \mathbb{R}^M$   $Y \approx f(x)$  "regression"

(2)  $Y_i$  discrete label  $Y_i = k, k \in \{1, \dots, C\}$  "classification" (number of categories)

(2a) labels ordered "ordinal classification"  
 "tiny" < "small" < "medium" < "large" < "huge"

(2b) labels unordered "categorical classification"  
 "apples", "oranges", "bananas"

training set:  $TS = \{ (X_i, Y_i) \}_{i=1}^N$

"supervised learning" {known response for instance  $i$ }



## types of learning approaches

- ① supervised learning: true responses are known in  $TS = \{(X_i, Y_i)\}_{i=1}^N$
- ② unsupervised learning: only features are known in  $TS = \{X_i\}_{i=1}^N$   
 $\Rightarrow$  learning alg. needs to find structure in the data on its own "data mining", "research"  
few solutions that guaranteed to work: • representation learning: compute better features  
 $\tilde{X} = \varphi(X)$ , ex: dimension reduction  $\dim(\tilde{X}) < \dim(X)$   
• clustering: group similar instances into "clusters"
- ③ weakly supervised: in between
- ③a) some information about  $Y_i$  (ex: know  $Y_i$  for some instances, coarse labels: "there is a tumor in CT, but don't know where")
- ③b) hot: self-supervised: define auxiliary tasks where  $Y_i$  are easy to determine, ex: large language models drop words from text