

# Fundamentals of Machine Learning (ML)

What is ML? Setting: - Quantities that are easy to observe or measure  
"features"  $X$  ("covariates")

- Quantities we would like to know but cannot easily observe or measure: "response"  $Y$   
("prediction")

$\Rightarrow$  need to define a function  $\hat{y} = f(x) \approx y^*$   
algorithm's response      true value of  $Y$

• traditional science: ask a theoretician to derive a formula for  $f(x)$

• machine learning: - use a "universal" function family

$$f(x; \theta)$$

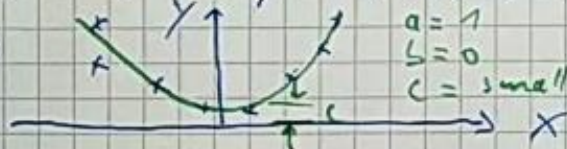
$\leftarrow$  "theta": additional parameters

By setting  $\theta$  cleverly  $f(x; \theta)$  can represent any  $f(x)$  we want  
 $f(x) = F(x; \hat{\theta})$

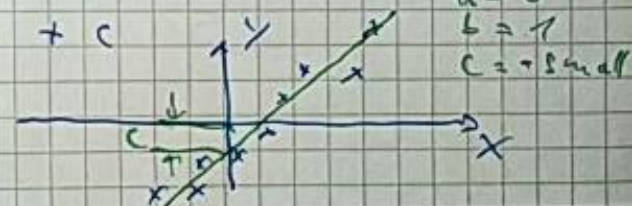
- step 2: use "training data" to find the best possible  $\hat{\theta}$   
such that  $F(x; \hat{\theta}) \approx y^*$

- Example:  $F(x; \theta = (a, b, c)) = ax^2 + bx + c$

Data 1:



Data 2:





- probabilistic solutions: sometimes, a deterministic function  $\hat{Y} = f(X)$  is not good enough: several solution  $Y$  can occur
  - the features contain only incomplete information about  $Y$ : "ambiguity"
  - "—" we worry: "aleatoric uncertainty", "variance"
  - the mapping in reality  $Y^* = f^*(X)$  may be non-deterministic  
 $\Rightarrow$  "aleatoric unc.", "variance"
  - our approximate function is not perfect, similar solutions are also plausible: "epistemic uncertainty", "bias"

$\Rightarrow$  solution: replace deterministic function with a conditional probability

$$\hat{Y} = f(X) \Rightarrow p(Y | X)$$

$\uparrow$  several solutions/outputs for  $Y$  with different probabilities

[\* we get back deterministic by  ~~$p(Y|X)$~~ ]

$$p(Y|X) = \delta(Y - f(X))$$

example:  $Y$  is a class label; kind of fruit "apple", "nectarine"

$X$ : size, color

[8 cm, red]  $\Rightarrow$

[4 cm, yellow]  $\Rightarrow$

[10 cm, green]  $\Rightarrow$

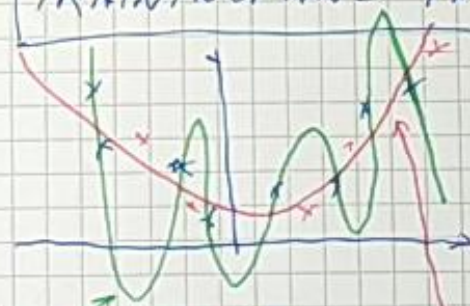
$$\begin{cases} p(Y = \text{"apple"} | 8 \text{ cm, red}) = \frac{1}{2} \\ p(Y = \text{"nectarine"} | 8 \text{ cm, red}) = \frac{1}{2} \\ p(Y = \text{"apple"} | 4 \text{ cm, yellow}) = 0.9 \end{cases}$$



• basic ML workflow:

- ① collect data - training data to determine  $\hat{\theta}$   
- test data to check if  $F(x, \hat{\theta})$   
is a good approximation for  $y^*$   
 $\Rightarrow$  detect "overfitting"

NEVER USE THE  
SAME DATA FOR  
TRAINING AND TEST



perfect interpolation of  $x$   
but bad "generalization"

no perfect  
interpolation  
but good  
generalization

- ② find  $\hat{\theta}$  such that  $\hat{y} = F(x; \hat{\theta}) \approx y^*$   
for the training data
- ③ test on the test data if  $\hat{y} = F(x; \hat{\theta}) \approx y^*$  (now  $\hat{\theta}$  is fixed)  
if not, throw away solution and start again (ex. use better  
function family  $F$ )
- ④ ~~top~~ deployment: run  $F(x, \hat{\theta})$  in a routine setting to  
predict  $\hat{y}$  for new data



Recap: predict response  $Y$  from observable features  $X$

• How good can we get? How bad can it be?

• Example: cards  
 $Y$ : color of the card (red or black)  
 $X$ : value of the card (ace, king, ...)

• Case 1: "best"  
 the feature is perfectly informative for the response  
 $\Rightarrow Y$  is always correctly predicted,

but: this almost never happens in practice

• Case 2: "worst"  
 we have no features, or features are unrelated to response  
 $\Rightarrow$  learning is impossible, but we can still use our prior knowledge:

case 2a)  $p(Y = \text{red}) = p(Y = \text{black}) = 0.5$

$\Rightarrow$  you can do no better than guessing

[observation: You will be right 50% of the time, sounds not too bad!]

case 2b)  $p(Y = \text{red}) < p(Y = \text{black})$   
 $\frac{1}{3} < \frac{2}{3}$

$\Rightarrow$  always go with the more probable answer

$\Rightarrow$  you will be right  $p(Y = \text{black})$  of the time

[if probability of majority answer is very high, 66%,  
 you can achieve impressive accuracy by doing nothing]



case 3: features give some information about the response  
 "typical"  $\Rightarrow$  model what happens with conditional probabilities:

3 kings (1 red / 2 black), 3 queens (2 red / 1 black)

$$p(Y=r | X=\text{king}) = 1/3$$

$$p(Y=b | X=\text{king}) = 2/3$$

$$p(Y=r | X=\text{queen}) = 2/3$$

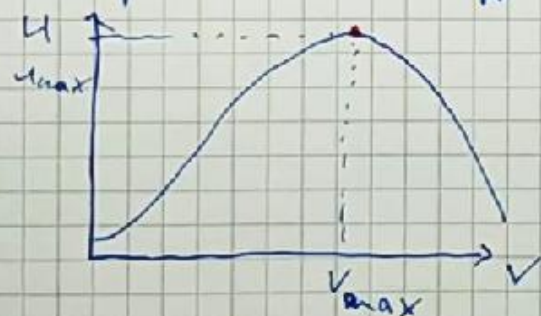
$$p(Y=b | X=\text{queen}) = 1/3$$

if you learn, that " $X=\text{king}$ ", what's the best answer

$\Rightarrow$  choose the most probable answer, given the evidence

mathematically:  $\hat{Y} = \arg \max_Y p(Y | X)$

[ Notation: arg max operator: suppose  $u = f(v)$



maximum:  $(v_{\max}, u_{\max})$

$$u_{\max} = \max_v f(v)$$

$$v_{\max} = \arg \max_v f(v)$$

if  $X=\text{king}$

$$\hat{Y} = \arg \max_{Y \in \{r, b\}} p(Y | X=\text{king})$$

$$= \arg \max \left\{ p(Y=r | X=\text{king}), p(Y=b | X=\text{king}) \right\}$$

$$= b \quad \frac{1}{3} < \frac{2}{3}$$



recap: machine learning defines generic function families, which depend on a parameter set  $\Theta$

$$Y = F(X; \Theta)$$

$\Rightarrow$  adjust  $\Theta$  such that the predicted response is as close to the truth as possible:

$$\text{choose } \hat{\Theta} \text{ such that } \hat{Y} = F(X; \hat{\Theta}) \approx Y^*$$

(measure by a "loss function": if small,  $\hat{Y}$  is close to  $Y^*$ )

- applies to deterministic responses, i.e. for each  $X$ , we return a hard decision about  $Y$  (hopefully, the best possible)
- in many situations, the uncertainty about the response should be made explicit

$\Rightarrow$  learn ~~a~~ ~~mean~~ define a generic family of conditional probabilities

$$P_{\Theta}(Y | X) \text{ or } p(Y | X; \Theta)$$

$\Rightarrow$  adjust  $\Theta$  such that the predicted probabilities for each possible response  $Y$  are as close as possible to the true probabilities ("good calibration of the probability")

$$\text{choose } \hat{\Theta} \text{ such that } \hat{p}(Y | X; \hat{\Theta}) \approx p^*(Y | X)$$

example: rain forecast:  $Y = 1$  : if rains tomorrow  
 $= 0$  : if does not rain tomorrow

if our model says  $\hat{p}(Y = \text{rain} | X) = 0.8$ ,  $\hat{p}(Y = \text{dry} | X) = 0.2$

then it should actually rain on 80% of the days with similar  $X$



- predicting a probability is strictly more general than making a deterministic prediction, because probabilities can always be reduced to a deterministic result, but not vice versa

examples: • choose the most probable answer as hard decision

"maximum a-posteriori" response (MAP)

$$\hat{Y} = \underset{Y}{\operatorname{arg\,max}} \underset{\text{weighted}}{p_{\theta}(Y|X)}$$

• choose the expected value (average prediction)

$$\hat{Y} = \mathbb{E}_{Y \sim p_{\theta}(Y|X)} [Y] = \begin{cases} \sum_Y Y \cdot p_{\theta}(Y|X) & \text{if } Y \text{ is discrete} \\ \int Y \cdot p_{\theta}(Y|X) dy & \text{if } Y \text{ is real} \end{cases}$$

• choose the median of the response

$$\hat{Y} = Y \text{ such that } p(Y < \hat{Y}|X) = p(Y > \hat{Y}|X)$$

For unimodal symmetric distributions  $p(Y|X)$ , all three methods give the same answer  $\hat{Y}$ , e.g. Gaussian

For skewed distributions, all three are different, e.g. Poisson

For multi-modal distributions (with modes that are almost equally high), reducing to a

single  $\hat{Y}$  is dangerous: miss alternative outcomes, get impossible outcomes





## Notation & terminology

18

- our observations consist of a set of instances

e.g. weather forecast: each instance may be a 6-day  
game: each round is an instance

- we assume the training set contains  $N$  instances,  
indexed by  $i \in 1, \dots, N$  (or  $i'$ ,  $i''$ )

the corresponding features and response are  $X_i, Y_i$

- we assume that  $X$  is a tuple or vector of  $D$  features  
indexed by  $j \in 1, \dots, D$  (or  $j'$ ,  $j''$ )

$X_{ij}$  is feature  $j$  of instance  $i$ .

$Y$  is a tuple or vector of  $D'$  responses indexed by  $j'$

but most of the time we consider the case  $D'=1$ , only one response

- types of variables (features and response)

• continuous  $\hat{=}$  real-valued

• discrete: - ordinal: values are ordered small < median < big  
- categorical: values not ordered apple, orange, nectarine

$\Rightarrow$  two fundamental types of machine learning problems

$Y_i$  is continuous  $\Rightarrow$  regression problem

$Y_i$  is discrete  $\Rightarrow$  classification problem

- if  $Y$  is discrete:  $C$  is the number of possible class labels / categories  
indexed by  $k \in 1, \dots, C$



- types of training sets:

- supervised learning: for each instance in TS, we know the features and the response:  $(x_i, y_i)$

[only when we apply the trained algorithm in the real world is  $y_i$  unknown]

⇒ training will adjust  $\hat{\theta}$  such that  $\hat{y}_i = F(x_i; \hat{\theta}) \approx y_i$  for all  $i \in 1, \dots, N$

- when we can observe  $y_i$  in the training set, why don't we do the same in the real application? why do we need learning at all?

⇒  $y_i$  is often only known in a "laboratory setting", not in the "wild"

- measuring  $y$  may require controlled conditions and/or very expensive equipment, not practical in the "wild"
- $y$  can be determined by a human expert (radiologist), who may not be available all the time so too expensive
- measuring  $y$  may be destructive, e.g. crash test for car
- $y$  may only become known in hindsight, e.g. weather, games

- ~~supervised~~ unsupervised learning: no  $y$  are known in TS

⇒ algorithm needs to figure out what interesting  $y$ 's might be and what their true values are, e.g. science  
existing solutions are quite limited - "data mining"

- weakly supervised learning: getting a fully "annotated" training set ( $y_i$  is known for every  $i$ ) is very expensive ⇒ looks for tricks to use less perfect TS