The background of the slide features a complex, abstract design. It consists of numerous translucent, glowing spheres of varying sizes connected by thin, dark lines, resembling a molecular network or a neural network. The colors are primarily shades of blue, teal, and orange, with some red and yellow highlights, creating a futuristic and scientific atmosphere.

# Artificial intelligence and chemistry

Tony Bonnaire

*Image generated with Midjourney « A representation of deep learning merging with chemistry »*

# Organisation of the course

Introduction to ML

Our first model

Supervised learning theory

Other basic models

Deep learning models



**Given by:** Tony Bonnaire (+ one practical session with Pablo Mas)



**Format:** 5 lectures, 1 hands-on session then data challenge (chemistry)



**Exam:** Oral presentation of your solution to the challenge



**Aim:** Introduce you to the basics of ML principles and carry out a project

Some references:

- [Pattern recognition and Machine Learning](#), Cristopher Bishop, 2006
- [Deep Learning](#), Goodfellow, Bengio and Courville, 2016
- Deep Learning with Python, Chollet, 2016
- <https://challengedata.ens.fr>: a bank of data science challenges to apply all the things we will learn in this course

# Generalities

Introduction to ML

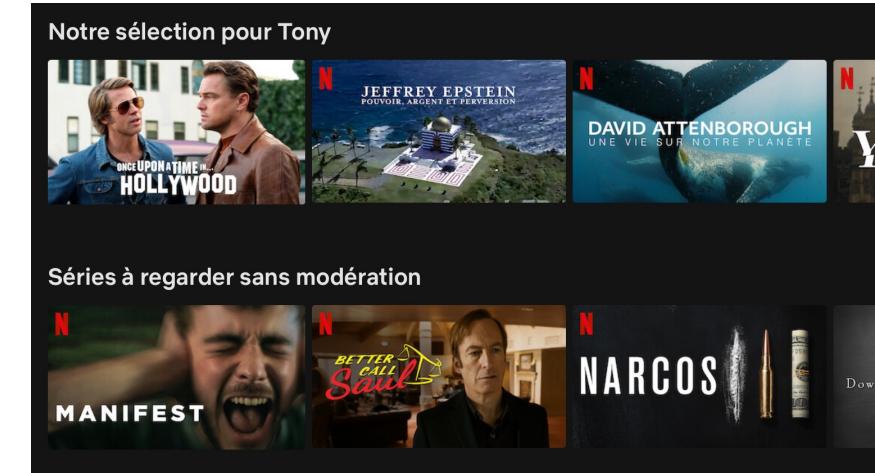
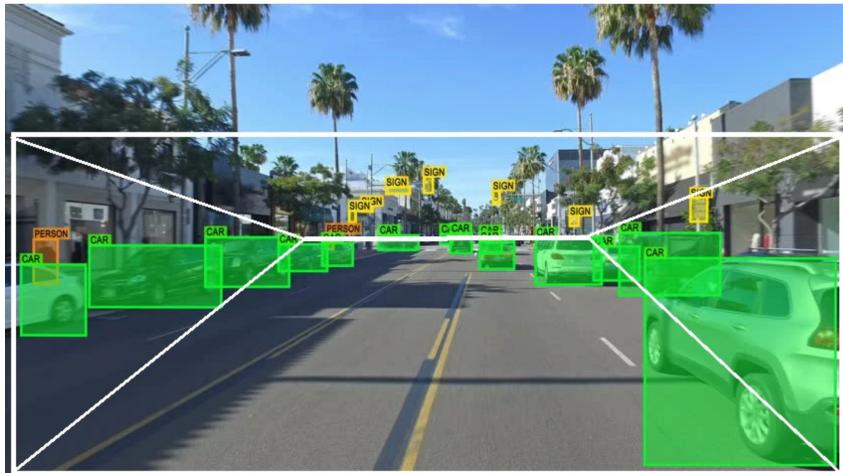
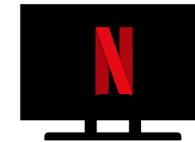
Our first model

Supervised learning theory

Other basic models

Deep learning models

AI, data science and machine learning are now embedded in our daily lives



# ML in science

Introduction to ML

Our first model

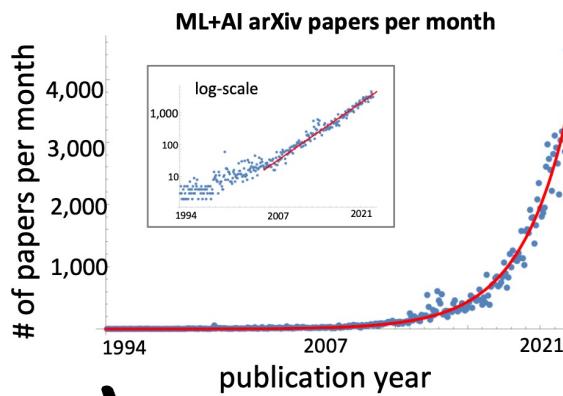
Supervised learning theory

Other basic models

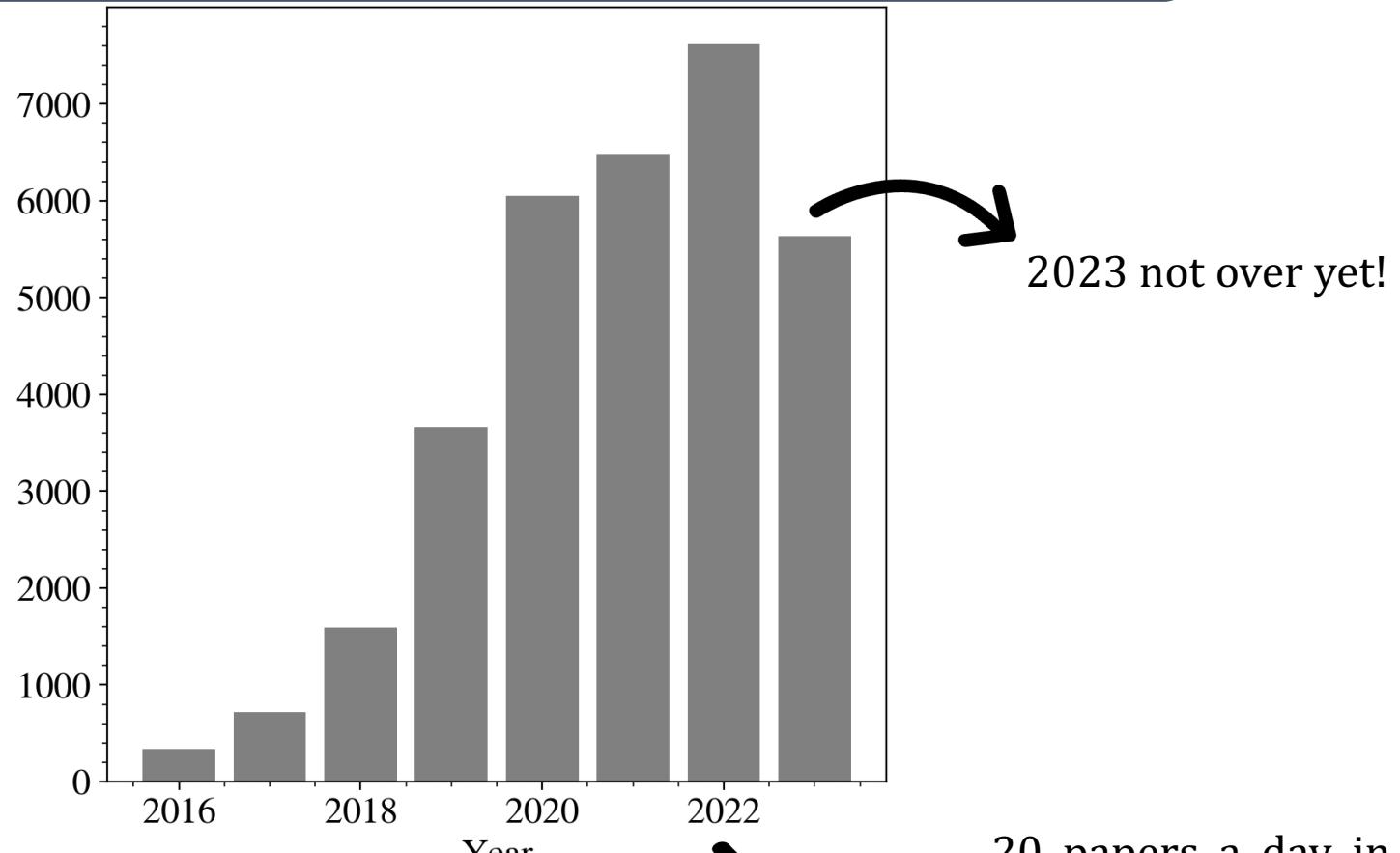
Deep learning models

But also in science

arXiv papers query with **physics, quantitative biology** and **electrical engineering** filters containing ML, AI or DL in abstract



Sidenote: ML papers from stat and cs filters grow exponentially fast! ([Kren et al., 2022](#))



2023 not over yet!

20 papers a day in average in 2022

# ML in science

Introduction to ML

Our first model

Supervised learning theory

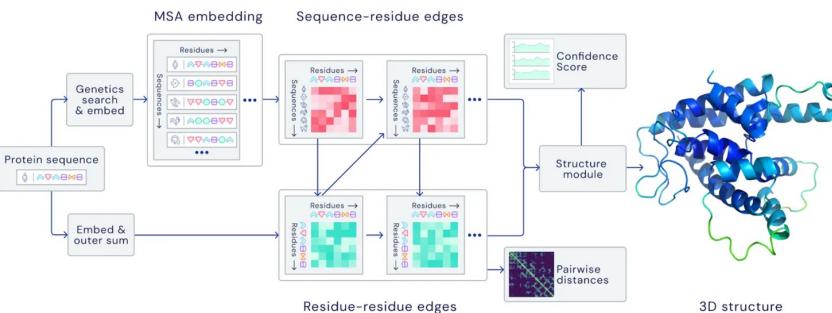
Other basic models

Deep learning models

But also in science

## Healthcare

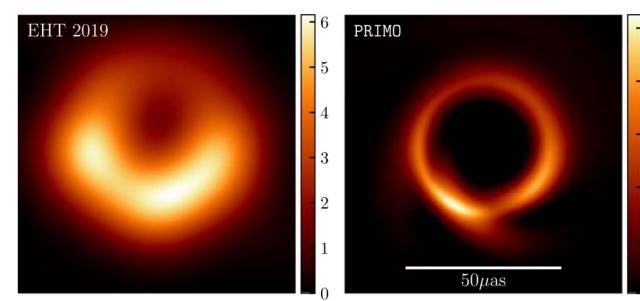
- Drug discovery
- Protein structure reconstruction



[Jumper et al., 2021](#)

## Astrophysics and cosmology

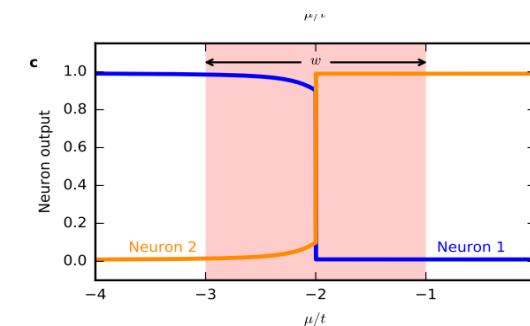
- Galaxy deblending
- Image restoration
- Source separation



[Medeiros et al., 2023](#)

## Theoretical physics

- Study phase transitions
- Discover experiments and equations



[Van Nieuwenburg et al., 2017](#)

# What is “learning”?

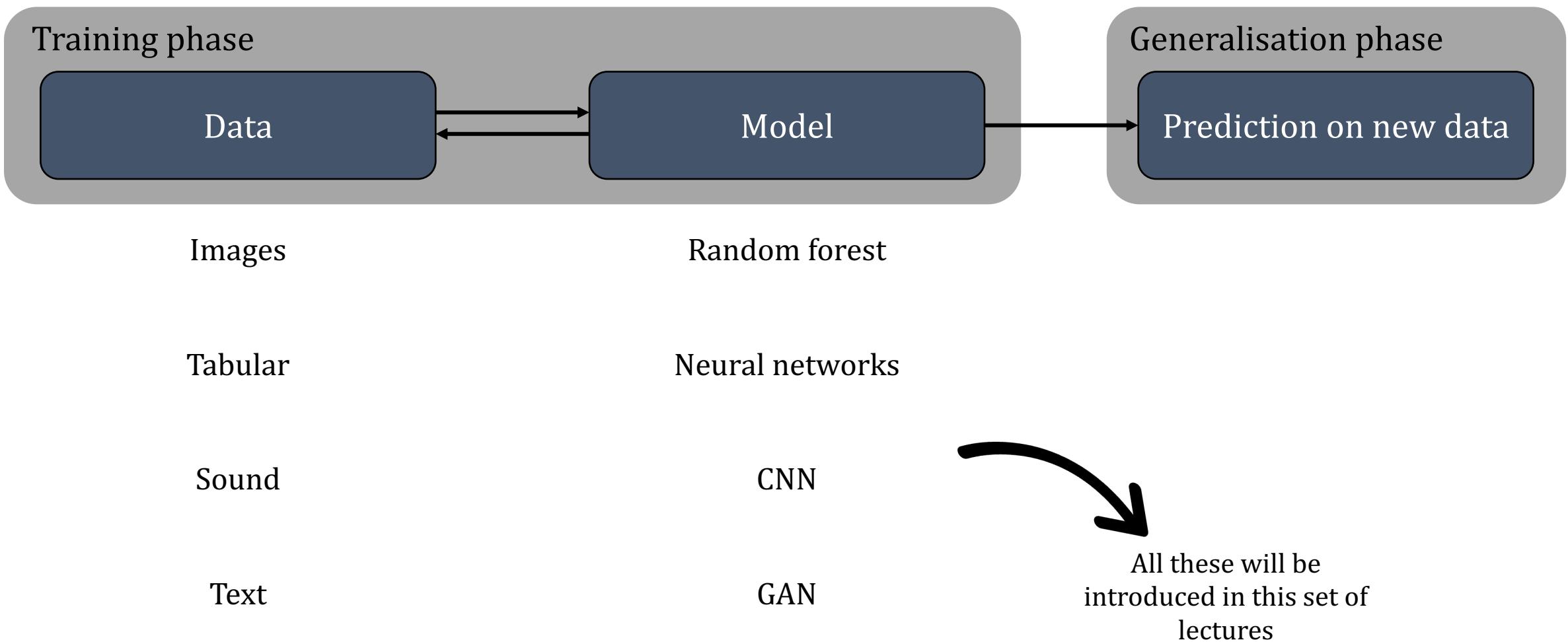
Introduction to ML

Our first model

Supervised learning theory

Other basic models

Deep learning models



...

...

# What is “learning”? An example

Introduction to ML

Our first model

Supervised learning theory

Other basic models

Deep learning models

Training phase

Data

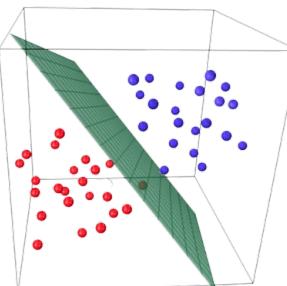
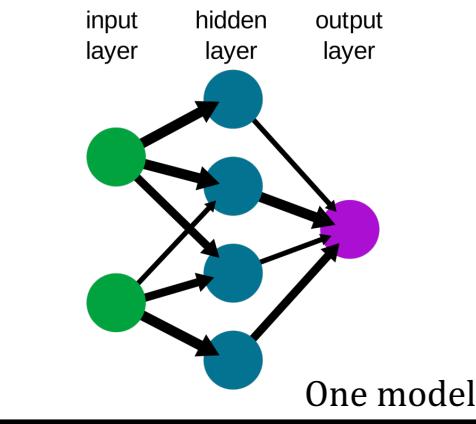
Model

Generalisation phase

Prediction on new data



Images of a “cat” or “dog”



Another model



“cat” or “dog” ?

# ML and the scientific method

Introduction to ML

Our first model

Supervised learning theory

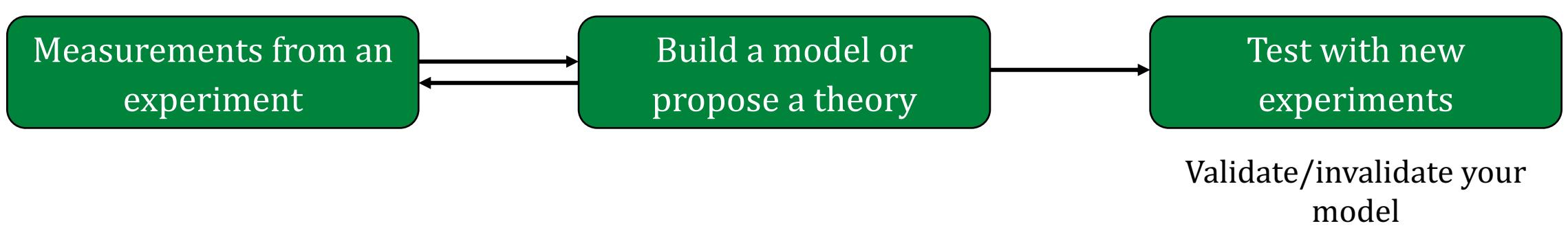
Other basic models

Deep learning models



...In fact, all this is close to what you know!

## The scientific method



# ML building blocks

Introduction to ML

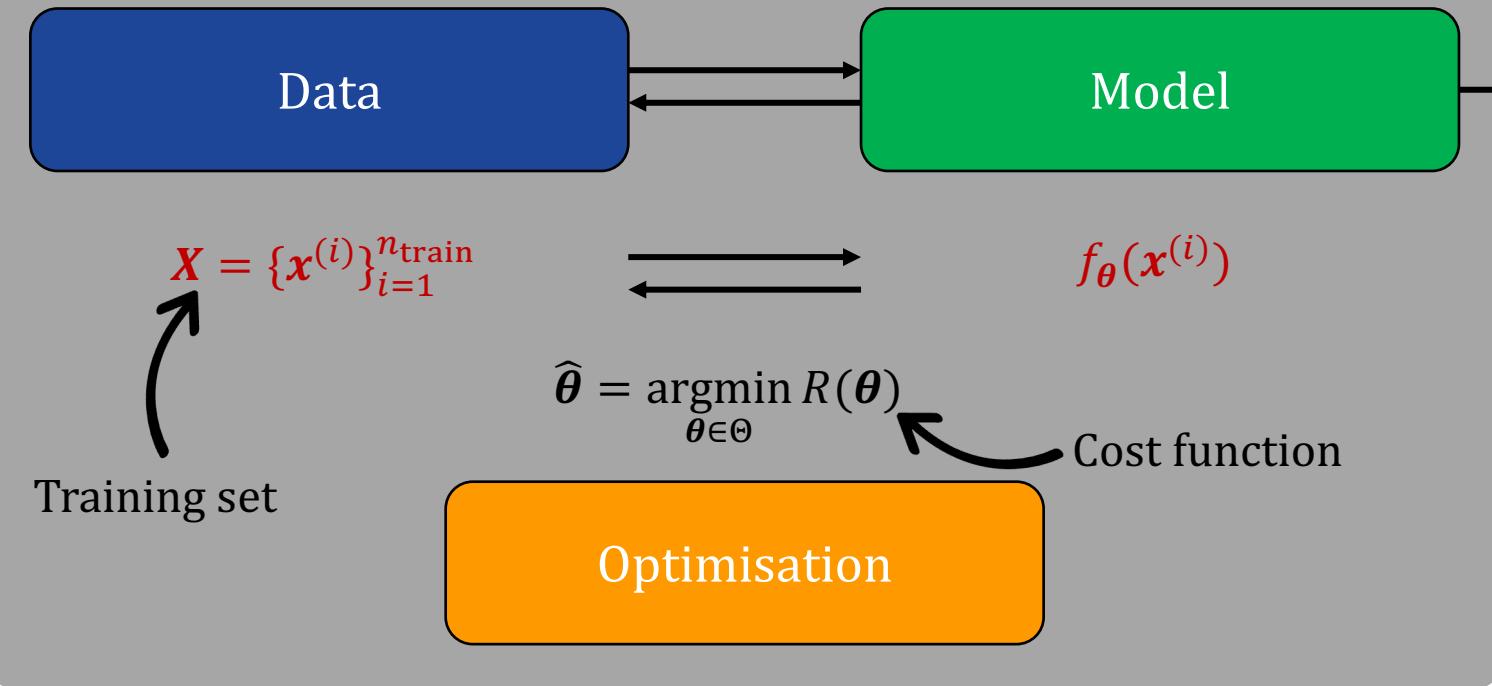
Our first model

Supervised learning theory

Other basic models

Deep learning models

## Training phase



## Generalisation phase

Prediction on new data

$$f_{\theta}(\tilde{x}^{(i)})$$

$\tilde{X}$  Test set

Some notations:

$$x^{(i)} \in \mathbb{R}^d \quad \theta \in \Theta \subset \mathbb{R}^p$$

Usually,

$$x^{(i)} \sim \rho$$



While it make take different forms depending on the problem and the model, the **cost function** measures how well (or more precisely *how bad*) your model is performing on the entire training set

# ML building blocks

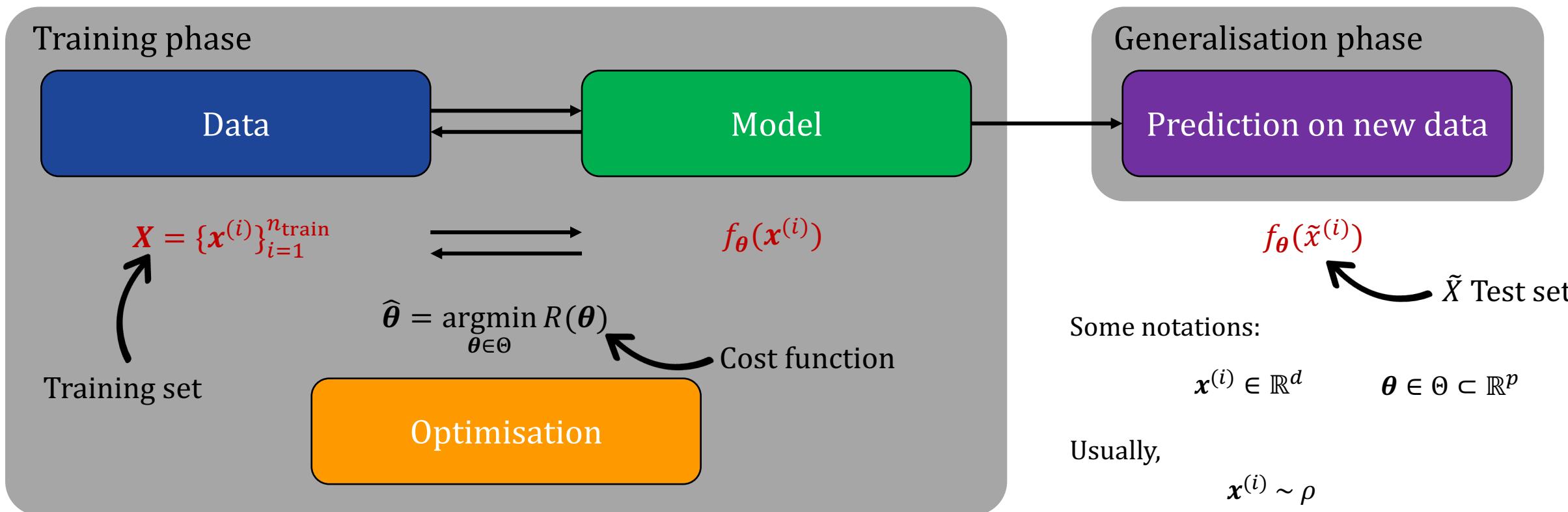
Introduction to ML

Our first model

Supervised learning theory

Other basic models

Deep learning models



Three main building blocks of *most* ML systems: some **data** feed a **parametrised model**. The best set of parameters are obtained through the **optimisation** of a cost function. The model can then be assessed on **new data**.

# Families of ML problems

Introduction to ML

Our first model

Supervised learning theory

Other basic models

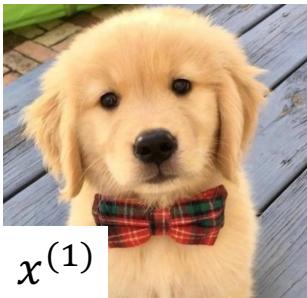
Deep learning models

## Supervised learning

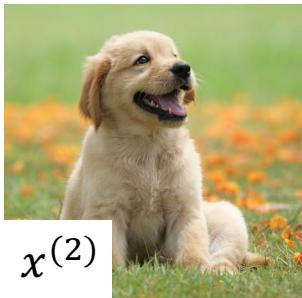
- Data are actually  $X$  and  $Y$  coming as pairs

$$\{(\mathbf{x}^{(i)}, \mathbf{y}^{(i)})\}_{i=1}^n, \quad (\mathbf{x}^{(i)}, \mathbf{y}^{(i)}) \in \mathbb{R}^d \times \mathbb{Y}, \quad \mathbb{Y} \subset \mathbb{R}^q$$

- Model  $f_\theta : X \rightarrow Y$  and usually  $f_\theta = p(y|x)$



$$X = \mathbf{x}^{(1)}$$



$$\mathbf{x}^{(2)}$$



$$\mathbf{x}^{(3)}$$



$$\mathbf{x}^{(4)}$$

Aim: Know if an image encodes a cat or a dog (called a classification task)

$$Y = \{1, 1, 0, 0\}$$

$$f_\theta(\mathbf{x}^{(i)}) = \hat{\mathbf{y}}^{(i)}$$

$$\longrightarrow$$

$$f_{\hat{\theta}}$$

$$\hat{\boldsymbol{\theta}} = \underset{\boldsymbol{\theta}}{\operatorname{argmin}} R(Y, \hat{Y}, \boldsymbol{\theta})$$

$$\text{with } R(Y, \hat{Y}, \boldsymbol{\theta}) = \frac{1}{n_{\text{train}}} \sum_{i=1}^{n_{\text{train}}} \ell(y^{(i)}, \hat{y}^{(i)})$$

**Loss function:** measures how bad your model is on a single example

Data

Model

Optimisation

# Families of ML problems

Introduction to ML

Our first model

Supervised learning theory

Other basic models

Deep learning models

## Supervised learning

- Data are actually  $X$  and  $Y$  coming as pairs

$$\{(\mathbf{x}^{(i)}, \mathbf{y}^{(i)})\}_{i=1}^n, \quad (\mathbf{x}^{(i)}, \mathbf{y}^{(i)}) \in \mathbb{R}^d \times \mathbb{Y}, \quad \mathbb{Y} \subset \mathbb{R}^q$$

- Model  $f_\theta : X \rightarrow Y$  and usually  $f_\theta = p(y|x)$

Examples of tasks

Classification

Regression

Timeseries prediction

Segmentation

Examples of models

Artificial Neural network

Random forest

Linear regression

Logistic regression

Naïve Bayes

Nearest neighbours

# Families of ML problems

Introduction to ML

Our first model

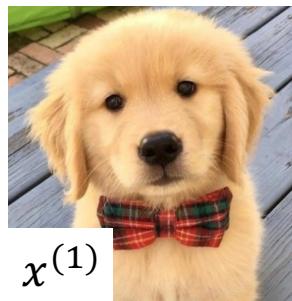
Supervised learning theory

Other basic models

Deep learning models

## Unsupervised learning

- Training data are the set of  $x^{(i)}$ 's only; no known results to predict
- In unsupervised learning, one seeks **patterns or structures** in  $X$  without prior labels



$$X = x^{(1)}$$



$$x^{(2)}$$



$$x^{(3)}$$



$$x^{(4)}$$

Aim: find two relevant classes to separate your dataset  
(called a clustering task)

$$f_{\theta}(x^{(i)}) = \hat{y}^{(i)}$$



$$f_{\hat{\theta}}$$

$$\hat{\theta} = \underset{\theta}{\operatorname{argmin}} R(\hat{Y}, \theta)$$

Data

Model

Optimisation

## Unsupervised learning

- Training data are the set of  $x^{(i)}$ 's only; no known results to predict
- In unsupervised learning, one seeks **patterns or structures** in  $X$  without prior labels

Examples of tasks

Clustering

Data augmentation

Dimensionality reduction

Sampling

Examples of models

Autoencoder

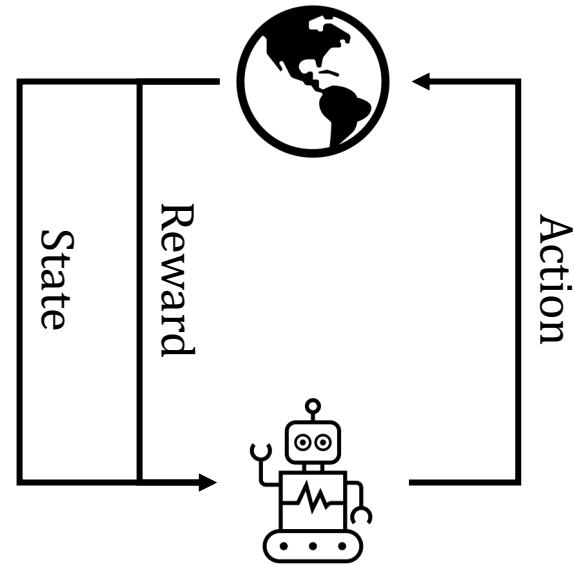
Boltzmann Machine

Gaussian mixture model

Generative Adversarial network

## Reinforcement learning

- The philosophy is different: the model does not try to “imitate” like in supervised learning nor to find patterns but “tries” things
- It is based on an **agent** interacting with an **environment**
- The agent tries to find the best possible sequence of states and actions to **maximise a reward**



Examples of tasks

Game theory

Robotics

Autonomous driving

Examples of models

Markov decision processes

Q-networks

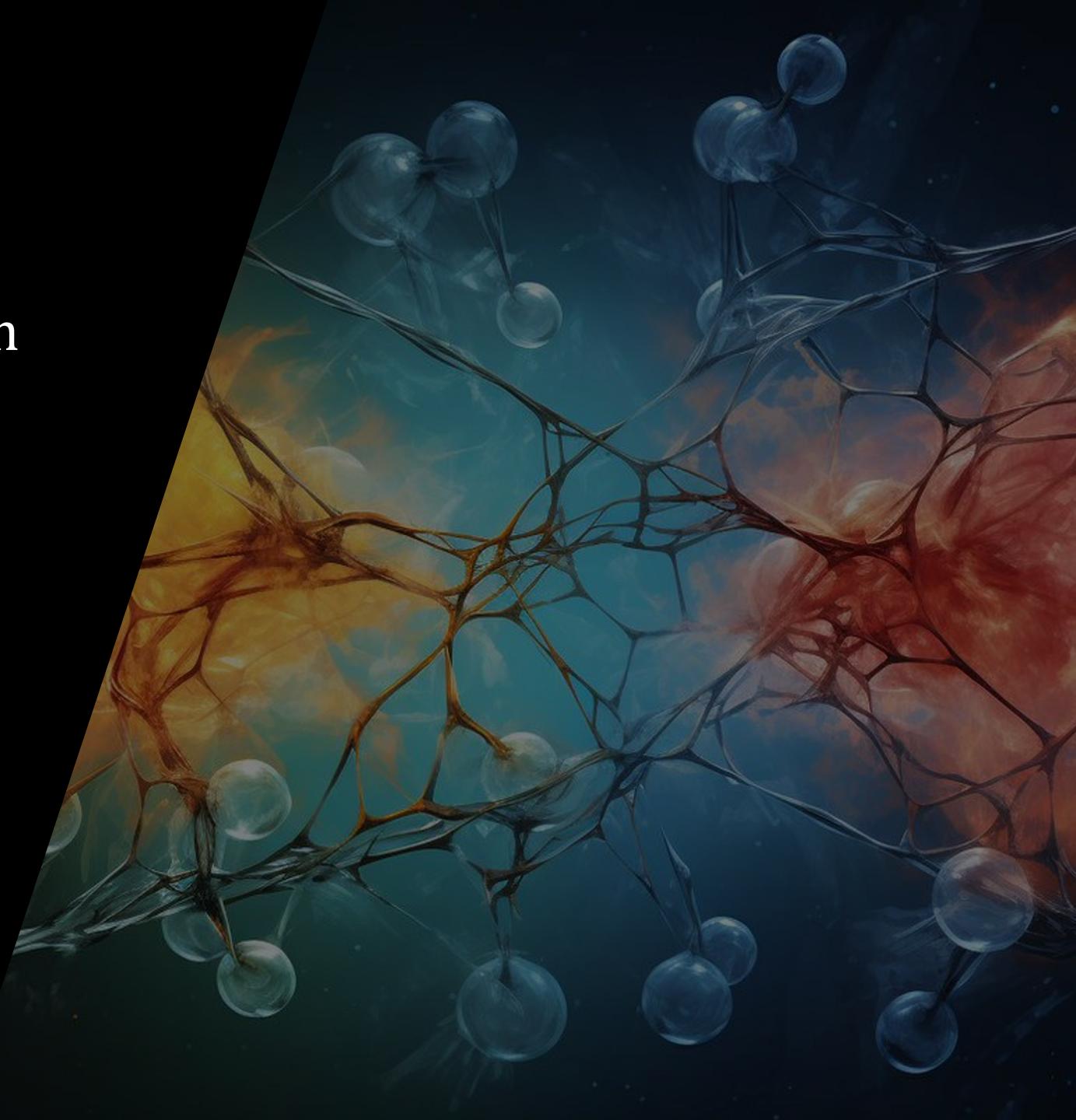
Deep policy gradient

Not discussed in this course, but a very good reference is [Reinforcement Learning – An introduction, Sutton and Barto, 2018](#)

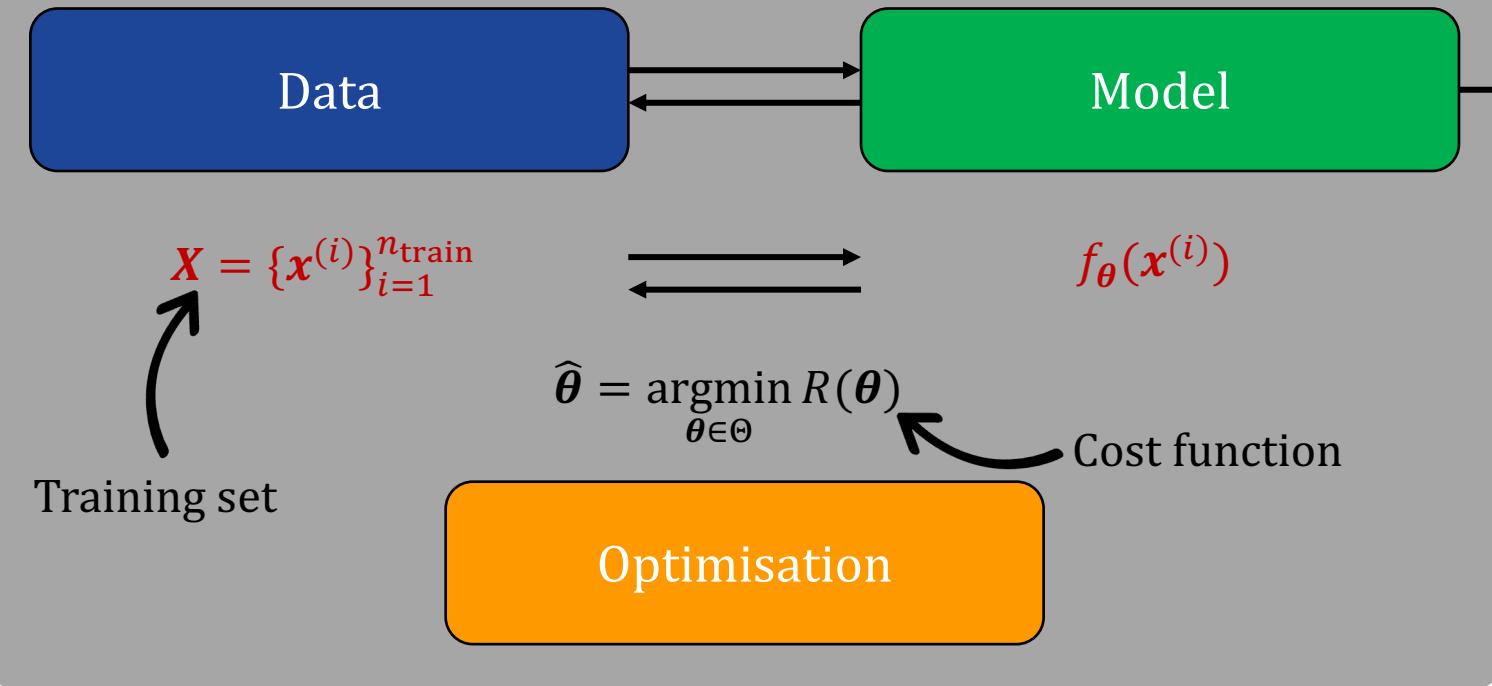
# Our first model: linear regression

Contents:

- *Linear regression model for supervised learning*
- *Optimal solution and links with maximum likelihood*
- *Gradient descent*
- *Stochastic gradient descent*



## Training phase



## Generalisation phase

Prediction on new data

$$f_{\theta}(\tilde{x}^{(i)})$$

$\tilde{X}$  Test set

Some notations:

$$x^{(i)} \in \mathbb{R}^d \quad \theta \in \Theta \subset \mathbb{R}^p$$

Usually,

$$x^{(i)} \sim \rho$$



Three main building blocks of *most* ML systems: some **data** feed a **parametrised model**. The best set of parameters are obtained through the **optimisation** of a cost function. The model can then be assessed on **new data**.

# Linear regression: the model

Introduction to ML

Our first model

Supervised learning theory

Other basic models

Deep learning models

## Linear regression

- Example: salary prediction based on the years of experience
- Data are  $n = 373$  couples  $(x^{(i)}, y^{(i)}) \Rightarrow$  **Supervised learning**
- The target variable  $y \in \mathbb{R}$  is continuous  $\Rightarrow$  **Regression**
- In linear regression, the model is

$$f_{\theta}(x_1^{(i)}) = \theta_0 + \theta_1 x_1^{(i)} = \hat{y}^{(i)},$$

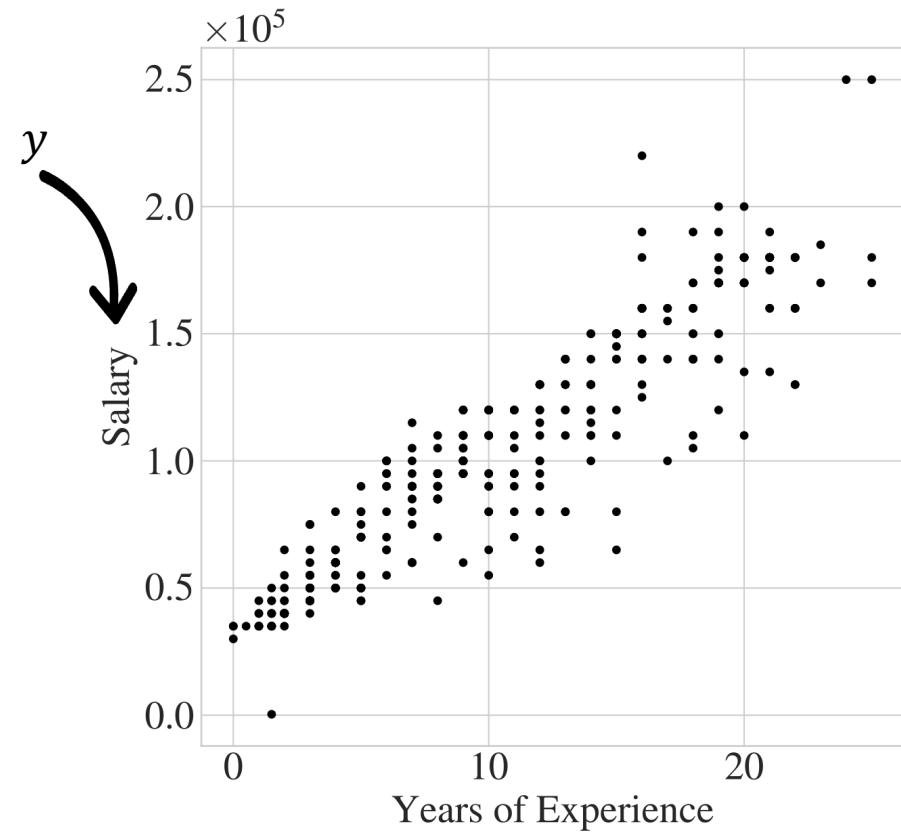
where  $x_1^{(i)}$  is the nb. of years of experience of the  $i^{\text{th}}$  training example

- It can be extended to additional features and functions of the features

$$f_{\theta}(x_1^{(i)}, x_2^{(i)}, \dots) = \theta_0 + \theta_1 x_1^{(1)} + \theta_2 x_2^{(i)} + \dots$$

$$f_{\theta}(x^{(i)}) = \boldsymbol{\theta}^T \mathbf{x}^{(i)},$$

where  $\boldsymbol{\theta} = [\theta_0, \theta_1, \theta_2, \dots]^T$ ,  $\mathbf{x}^{(i)} = [1, x_1^{(i)}, x_2^{(i)}, \dots]^T$



## Linear regression

- Now the model is fixed, how to find  $\hat{\theta}$ , the best possible parameters for our model and data?
- A natural way is assuming a **squared loss function**

$$\hat{\theta} = \underset{\theta}{\operatorname{argmin}} R(\theta) = \underset{\theta}{\operatorname{argmin}} \frac{1}{n} \sum_{i=1}^n (\hat{y}^{(i)} - y^{(i)})^2$$

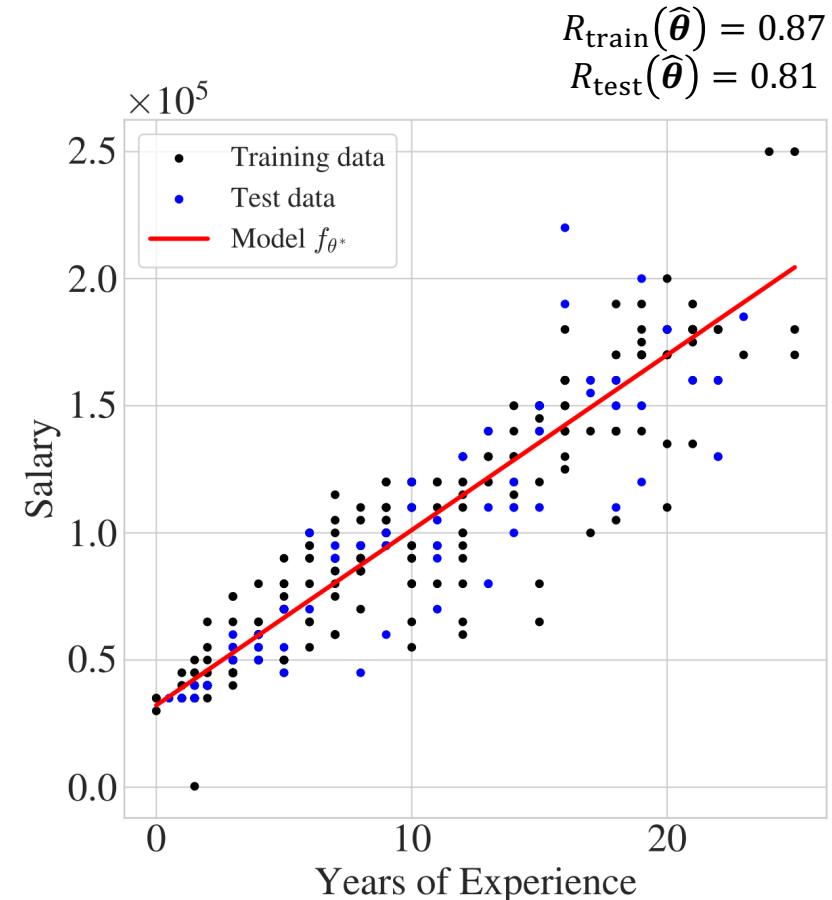
- Here**, the optimisation problem can be solved analytically in closed-form giving (usually not the case):  $\hat{\theta} = (X^T X)^{-1} X^T y$
- Note that this is the maximum likelihood estimator, assuming Gaussian error between  $\hat{y}^{(i)} = f_{\theta}(x^{(i)})$  and  $y$



Exactly solvable, low variance



Cannot represent local relationships, may be biased



I separated the dataset into **training and test sets**,  $n_{\text{train}} = 0.8n$  and  $n_{\text{test}} = 0.2n$  chosen randomly

# How to perform ERM in general?

Introduction to ML

Our first model

Supervised learning theory

Other basic models

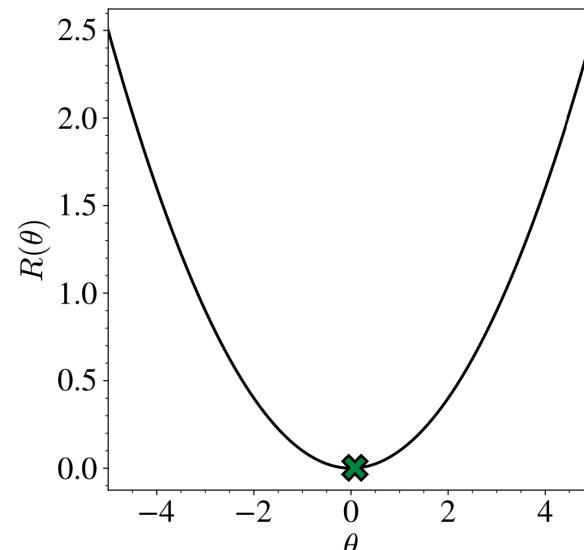
Deep learning models

- In linear regression, we could directly optimise the empirical risk in closed-form, but **in general it is not possible**
- How to perform the **Empirical Risk Minimisation** (ERM) in general? The aim is to compute

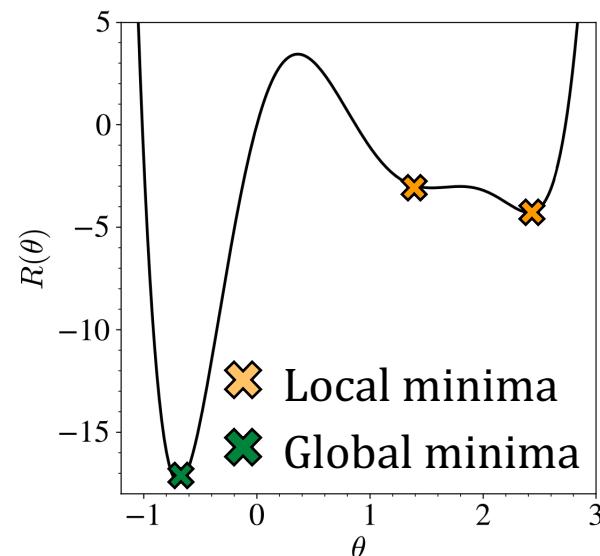
$$\hat{\boldsymbol{\theta}} = \operatorname{argmin}_{\boldsymbol{\theta}} R(\boldsymbol{\theta}) = \operatorname{argmin}_{\boldsymbol{\theta}} \frac{1}{n_{\text{train}}} \sum_{i=1}^{n_{\text{train}}} \ell(f_{\boldsymbol{\theta}}(x^{(i)}), y^{(i)})$$

- Our goal is hence to minimize a function of (possibly many) parameters under a loss function encoding how bad the model approximates the target variable

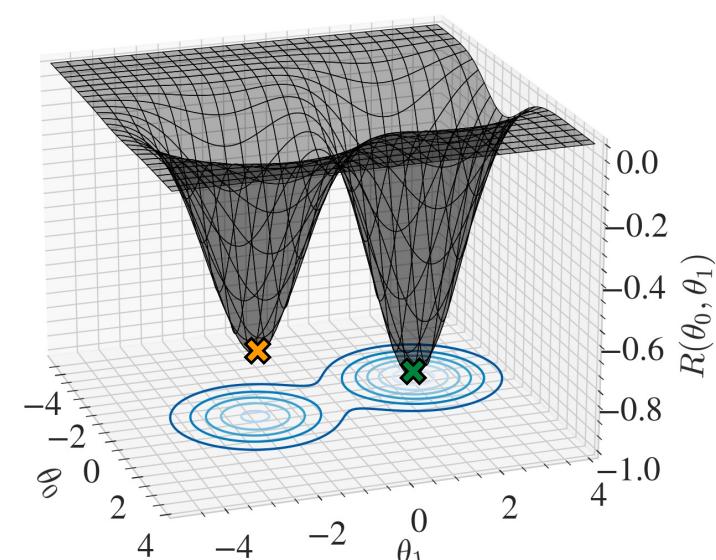
1D convex



1D non-convex



2D non-convex



# Naïve view and curse of dimensionality

20

Introduction to ML

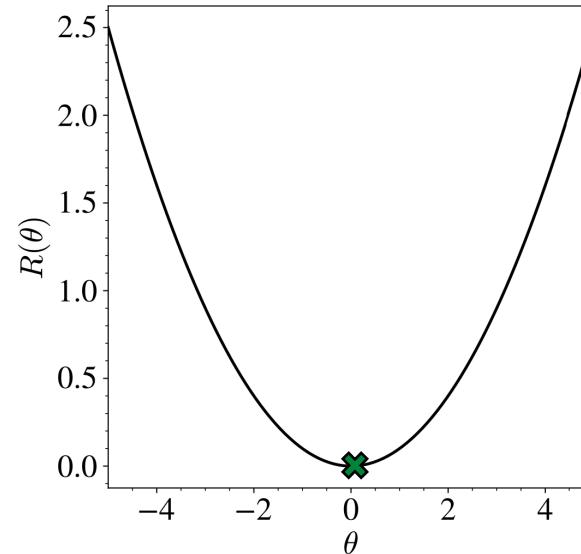
Our first model

Supervised learning theory

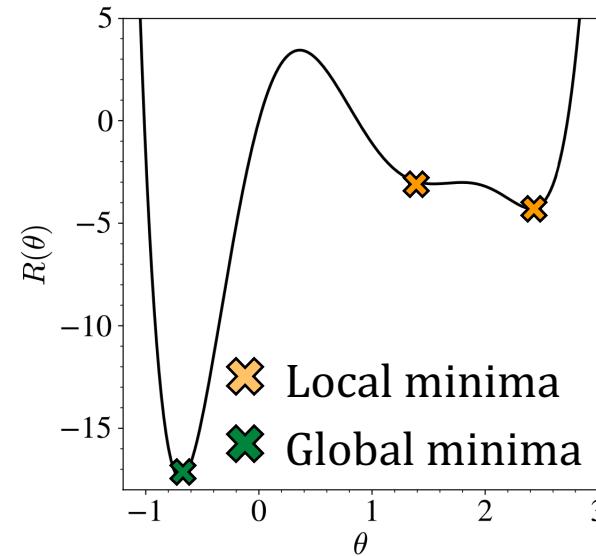
Other basic models

Deep learning models

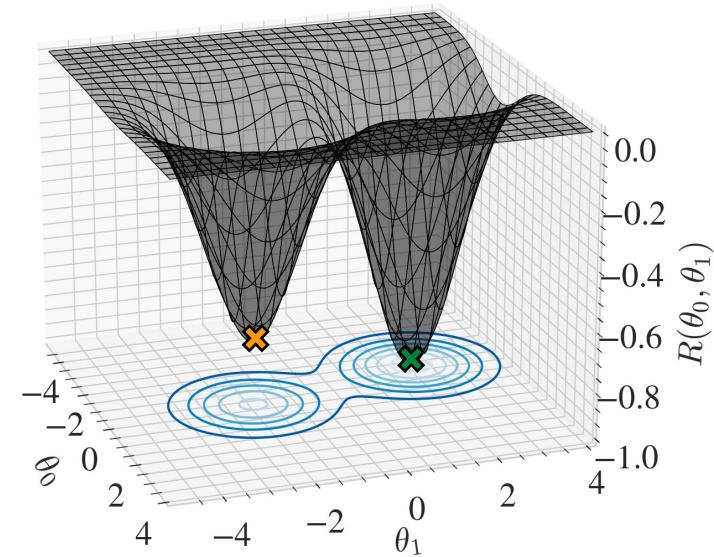
1D convex



1D non-convex



2D non-convex



- A naïve way of minimizing such functions could be to uniformly pave the parameter space and choose the one with the smallest value as being the minimum: this is **global optimization (grid search optimization)**
- This is where the **curse of dimensionality** kicks in. In general, we optimize models over many parameters  $d \gg 1$  (imagine the pixels of an image) and all the points are far away from each other in high dimensions
- Sampling uniformly  $[0, 1]^{10}$  with a step of 0.01 requires  $10^{20}$  evaluations (think of GPT-3 and its 175 billion parameters!)

# Gradient descent algorithm

Introduction to ML

Our first model

Supervised learning theory

Other basic models

Deep learning models

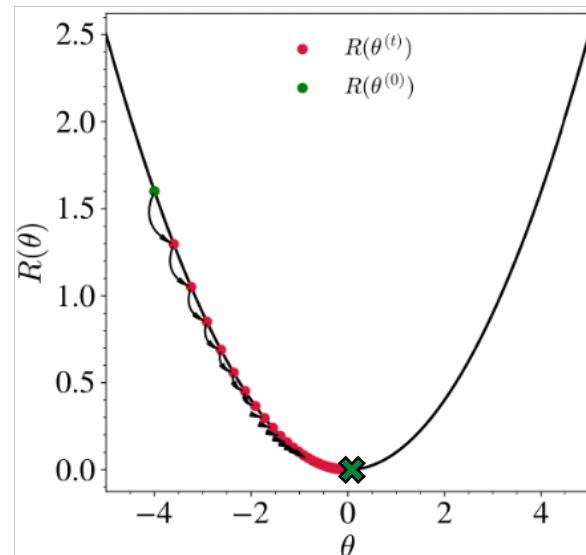
- A solution: local, directed search to navigate through the landscape
  - Numerical optimisation by **gradient descent**

Note: the superscript does not have to do with the training example here but with the time step ( $\theta$  is a parameter).

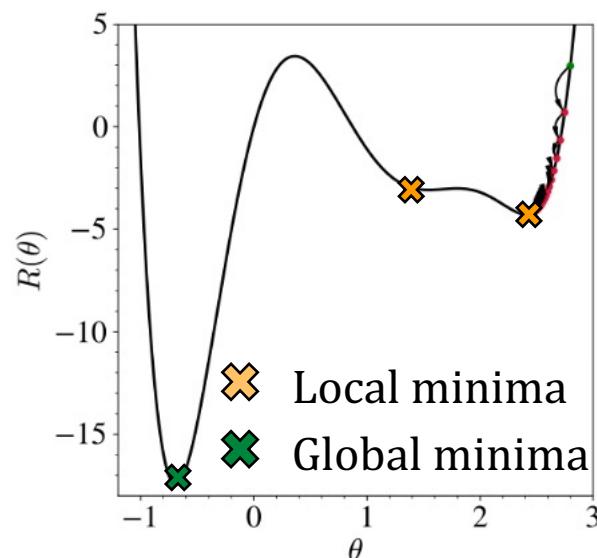
## Algorithm: Gradient descent

- Initialise  $\theta_0$  randomly
- Compute  $\theta^{(t+1)} = \theta^{(t)} - \eta \nabla_{\theta} R(\theta^{(t)})$
- Repeat 2 until  $\|\theta^{(t+1)} - \theta^{(t)}\|_2^2 \leq \epsilon$

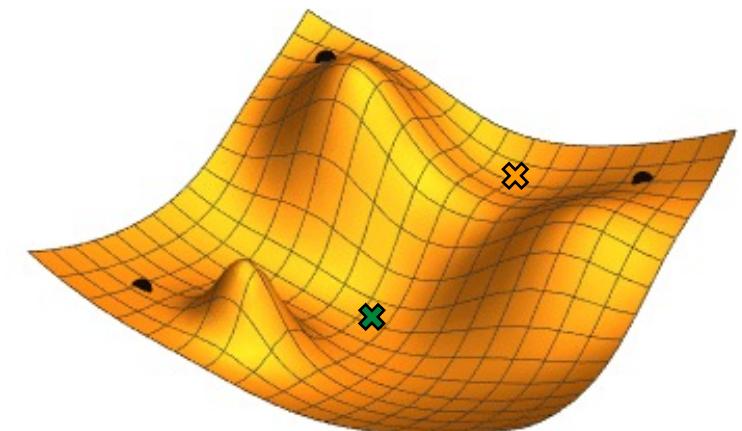
1D convex



1D non-convex



2D non-convex



# A word about hyperparameters

Introduction to ML

Our first model

Supervised learning theory

Other basic models

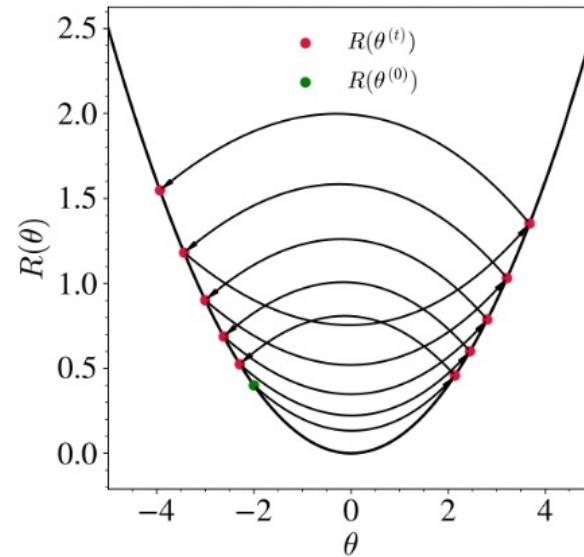
Deep learning models

- One **hyperparameter**: the learning rate  $\eta$
- A value that is too large can lead to divergence while, when too small, the computational cost explodes (+ stuck in small asperities)

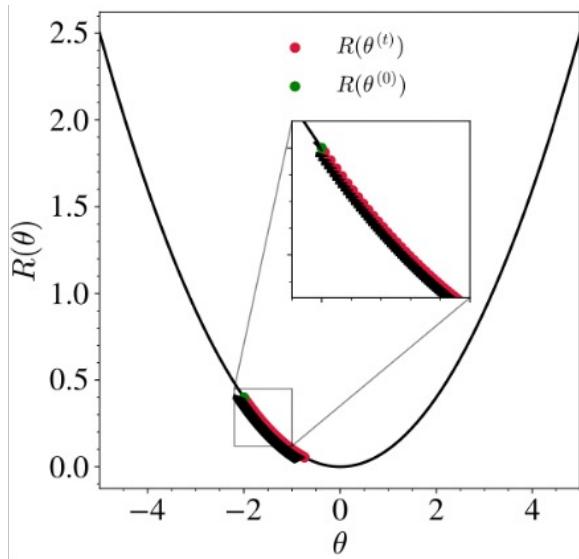


**Hyperparameter:** parameter that is **not learned** during the optimisation.  
Ex: depth of tree in DTs, # of trees in RFs, learning rates, etc.

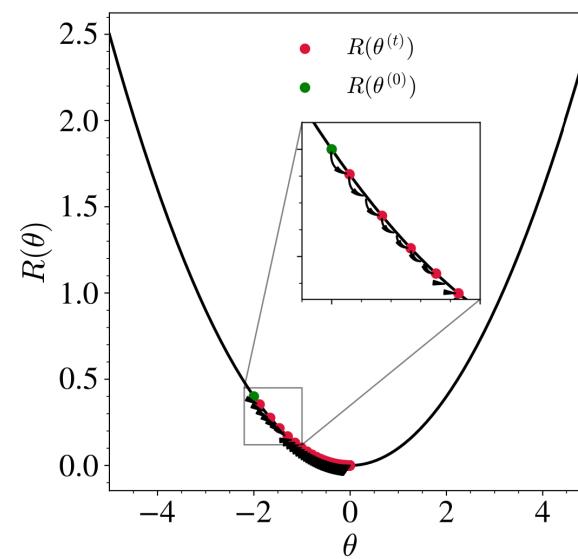
$\eta$  too large



$\eta$  too small



appropriate  $\eta$



- Usually, we use **grid search** to find the hyperparameter performing best on a third dataset: the **validation set (or use cross-validation)**

# Example: our linear regression

23

Introduction to ML

Our first model

Supervised learning theory

Other basic models

Deep learning models

- In the linear regression, the risk is

$$R(\boldsymbol{\theta}) = \frac{1}{n} \sum_{i=1}^n (\theta_0 + \theta_1 x_1^{(i)} - y^{(i)})^2 \quad \boldsymbol{\theta} = [\theta_0, \theta_1]^T$$

- Let's apply the gradient descent algorithm starting from  $\boldsymbol{\theta}^{(0)}$  random
- Then, we need to compute the gradient  $\nabla_{\boldsymbol{\theta}} R(\boldsymbol{\theta})$

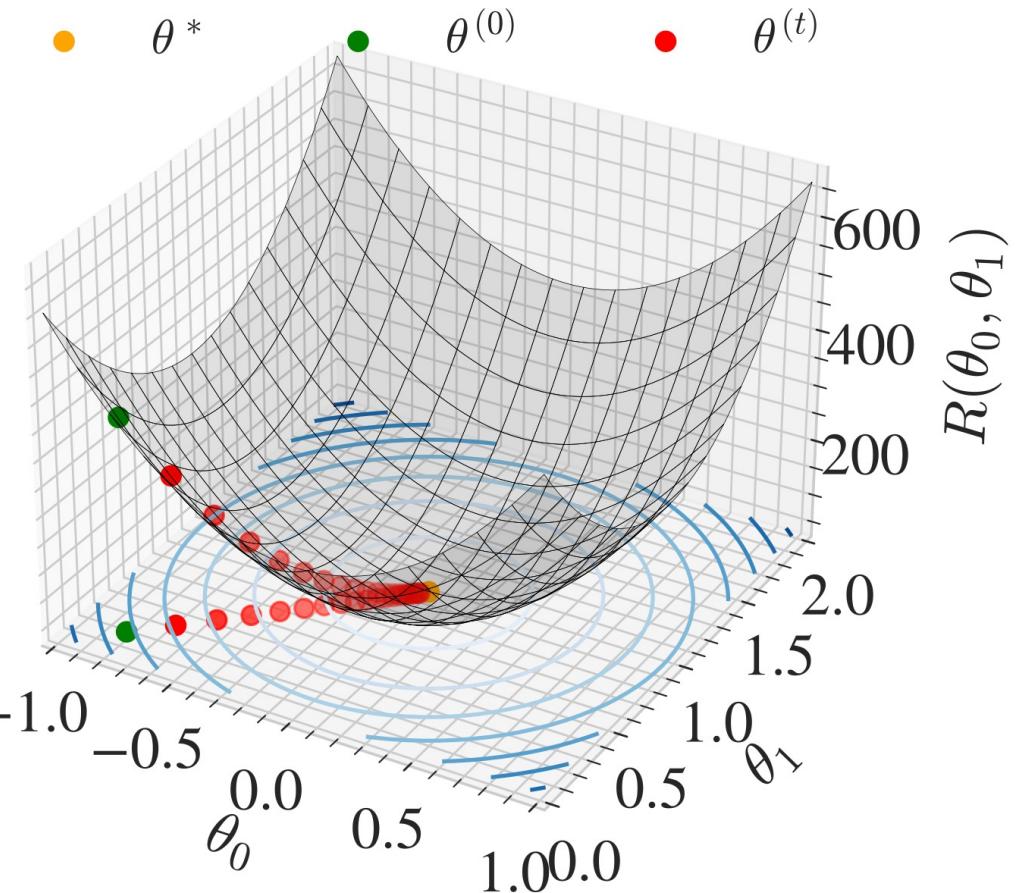
$$\frac{\partial R(\boldsymbol{\theta})}{\partial \theta_0} = \frac{2}{n} \sum_{i=1}^n (\theta_0 + \theta_1 x_1^{(i)} - y^{(i)}),$$

$$\frac{\partial R(\boldsymbol{\theta})}{\partial \theta_1} = \frac{2}{n} \sum_{i=1}^n x_1^{(i)} (\theta_0 + \theta_1 x_1^{(i)} - y^{(i)})$$

- Therefore, the update is

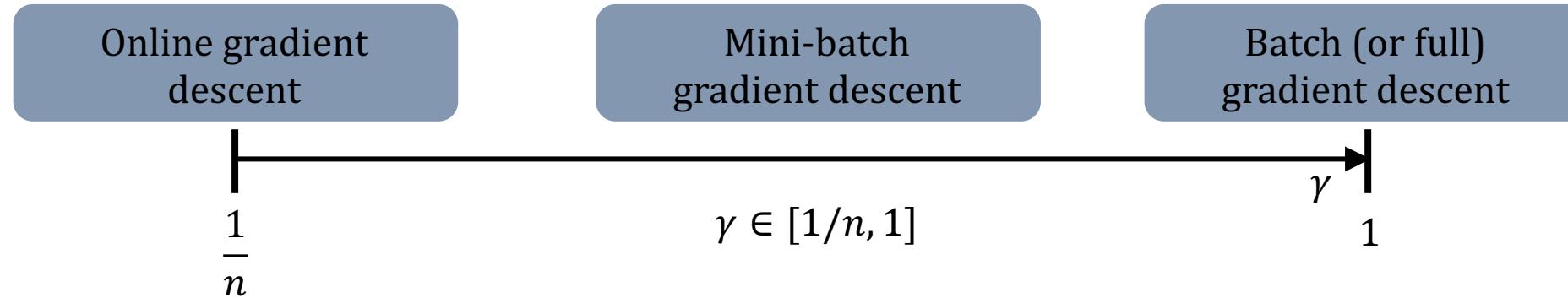
$$\theta_j^{(t+1)} = \theta_j^{(t)} - \frac{\eta}{n} \sum_{i=1}^n x_j^{(i)} (\theta_0 + \theta_1 x_1^{(i)} - y^{(i)})$$

where  $\forall i \in [1, \dots, n]$   $x_0^{(i)} = 1$



Note that I **standardized** the features. This is sometimes required when optimising some models. In GD, it allows faster convergence.

- Problem of gradient descent: we **need the entire dataset** to compute  $\nabla_{\theta} R(\theta^{(t)})$
- Solution: what about using only a fraction  $\gamma$  of the dataset chosen randomly?



## Algorithm: Stochastic gradient descent

1. Initialise  $\theta_0$  randomly
2. For  $e \in [1, \dots, E]$ 
  - 2.1. Shuffle the dataset
  - 2.2. For  $i \in [1, \dots, |\gamma n|]$ 
    - 2.2.1. Compute  $\theta^{(i+1)} = \theta^{(i)} - \eta \widehat{\nabla}_{\theta} R(\theta^{(i)})$

$$\text{where } \widehat{\nabla} R(\theta) = \sum_{j=i\gamma n}^{i\gamma n + \gamma n} \nabla R_j(\theta)$$

- $e$  is called an **epoch** and a set of  $\gamma n$  training examples is called a **mini-batch**
- Usually, **SGD often converges faster than full-batch GD**
- It may however oscillate around the true minimum
- Under some technical assumptions, **SGD provides an almost sure convergence** to a local (resp. global) in non-convex (resp. convex) landscapes

# SGD landscapes on our linear regression

25

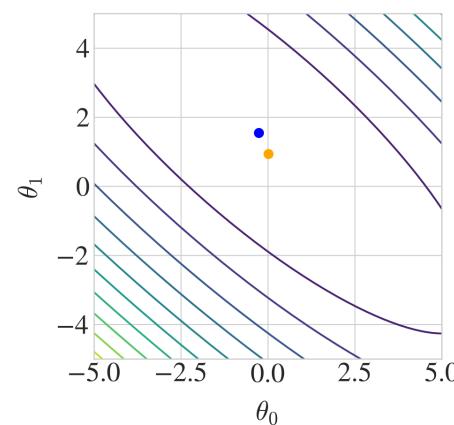
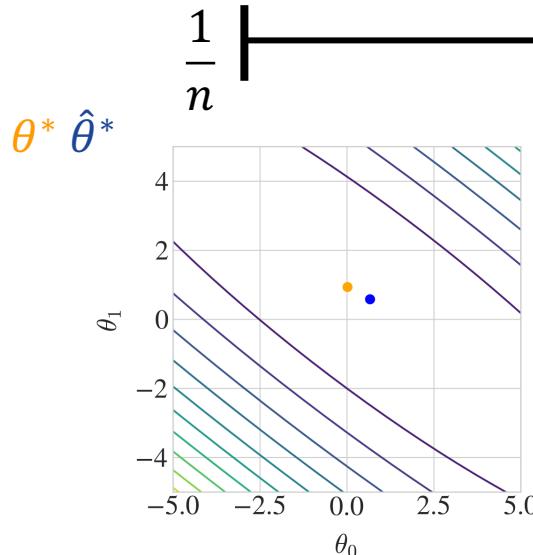
Introduction to ML

Our first model

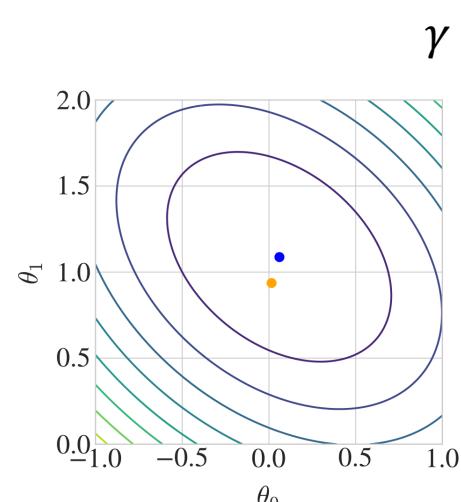
Supervised learning theory

Other basic models

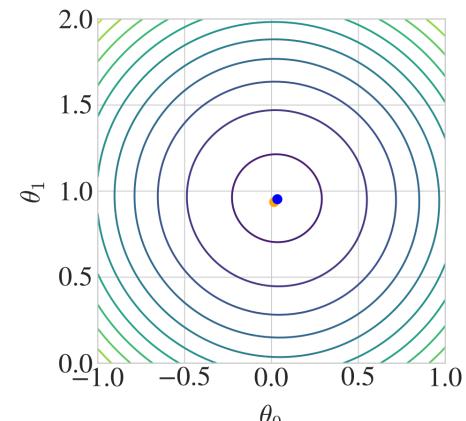
Deep learning models

Online gradient  
descentMini-batch  
gradient descentBatch (or full)  
gradient descent

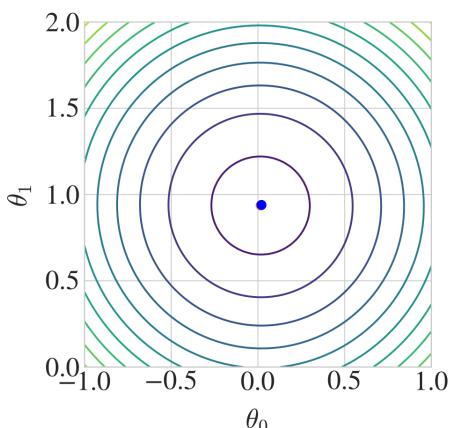
$$\gamma = 2/n$$



$$\gamma = 0.1$$



$$\gamma = 0.5$$



$$\gamma = 1$$

$$\gamma \in [1/n, 1]$$

# SGD on our linear regression

26

Introduction to ML

Our first model

Supervised learning theory

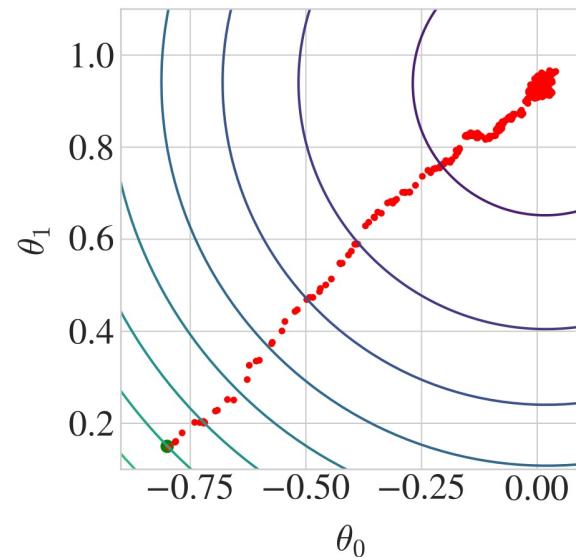
Other basic models

Deep learning models

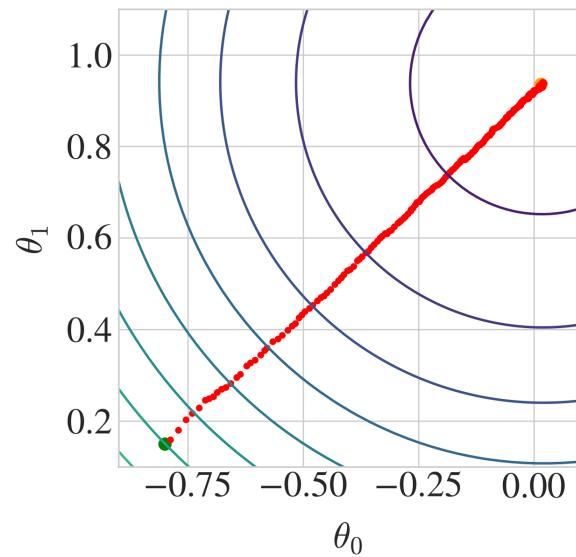
Online gradient  
descentMini-batch  
gradient descentBatch (or full)  
gradient descent

$$\frac{1}{n} \rightarrow 1$$

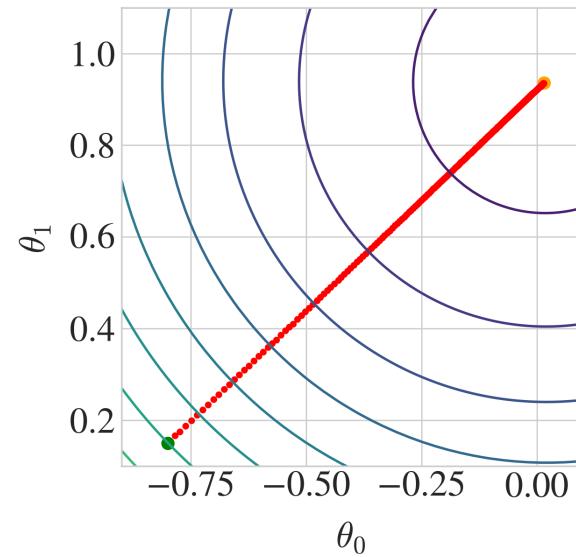
$$\gamma \in [1/n, 1]$$



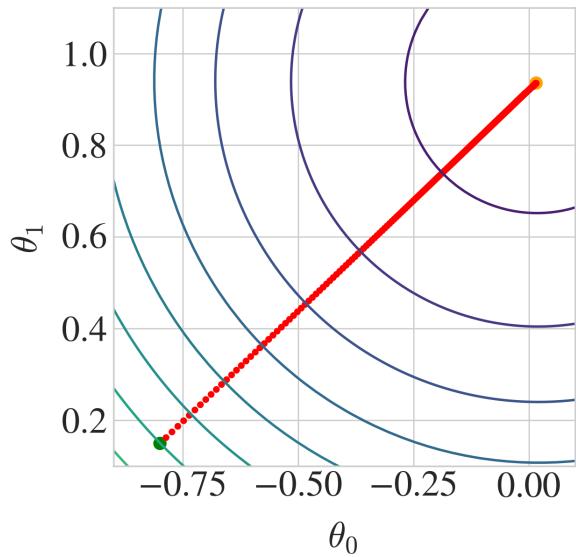
$$\gamma = 2/n$$



$$\gamma = 0.1$$



$$\gamma = 0.5$$

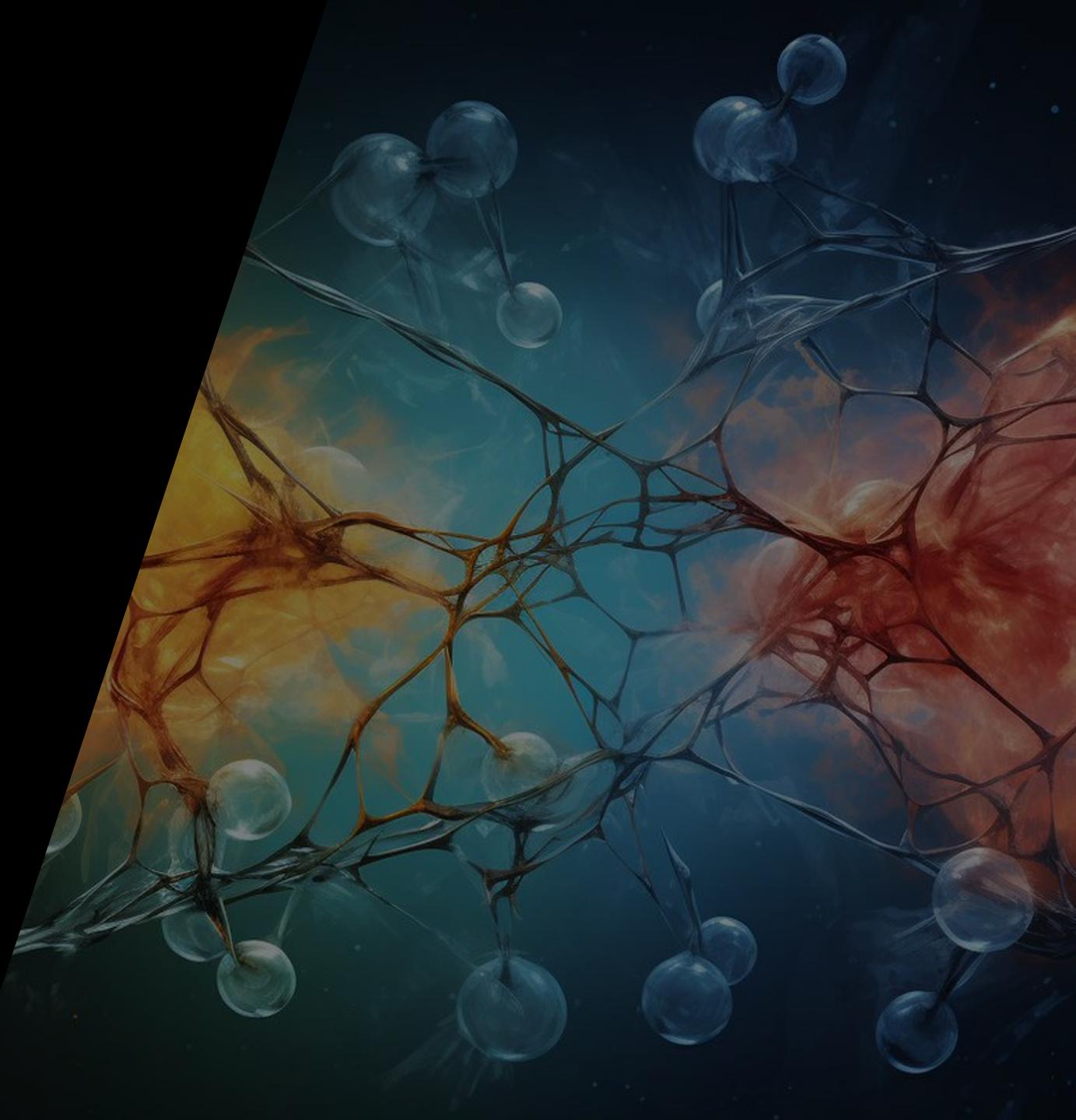


$$\gamma = 1$$

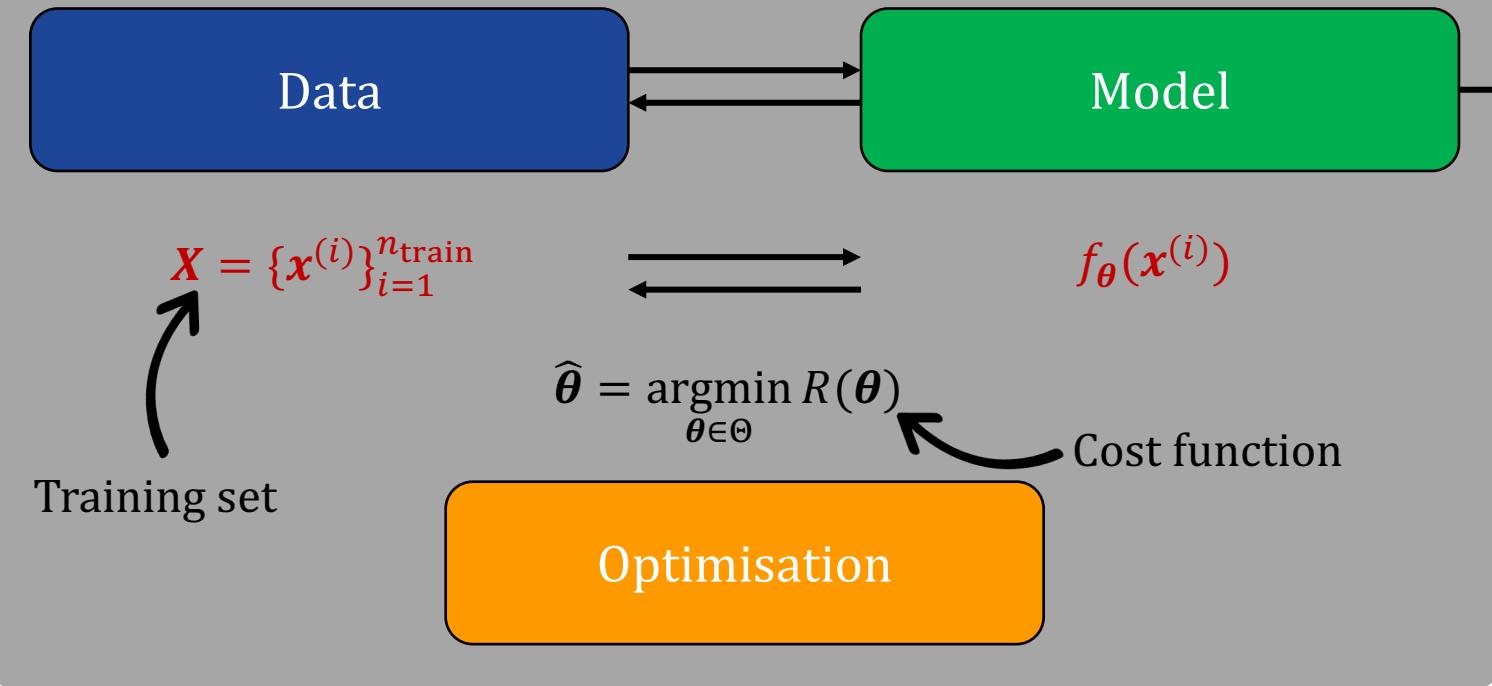
# A theory of supervised learning

*Contents :*

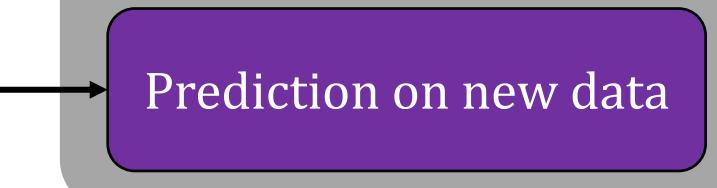
- *Bias-variance trade-off*
- *Generalisation: overfitting and underfitting*
- *Introduction to regularisation*
- *Application to linear regression*



## Training phase



## Generalisation phase



Some notations:

$$x^{(i)} \in \mathbb{R}^d \quad \theta \in \Theta \subset \mathbb{R}^p$$

Usually,

$$x^{(i)} \sim \rho$$



Three main building blocks of *most* ML systems: some **data** feed a **parametrised model**. The best set of parameters are obtained through the **optimisation** of a cost function. The model can then be assessed on **new data**.

# The bias-variance trade-off

Introduction to ML

Our first model

Supervised learning theory

Other basic models

Deep learning models

- The whole aim of training supervised models is to generalise well on unseen data. In practice, we would like to minimise  $\mathbb{E}_{x,y}(\ell(f_{\theta}(x), y))$  that we approximate by  $\frac{1}{n_{\text{train}}} \sum_{i=1}^{n_{\text{train}}} \ell(f_{\theta}(x^{(i)}), y^{(i)})$
- In a regression context, suppose there exists  $f$  such that  $y^{(i)} = f(x^{(i)}) + \epsilon^{(i)}$ , with  $\mathbb{E}[\epsilon^{(i)}] = 0, \mathbb{E}[\epsilon^{(i)2}] = \sigma_{\epsilon}^2$
- We build a model  $f_{\theta}$  of  $f$  minimising the squared error  $\ell(f_{\theta}(x^{(i)}), y^{(i)}) = (y^{(i)} - f_{\theta}(x^{(i)}))^2$
- We can show that the expected error on a test example  $\tilde{x}$  decomposes as

$$\mathbb{E}[(y - f_{\theta}(\tilde{x}))^2] = \mathbb{E}[f_{\theta}(\tilde{x}) - f(\tilde{x})]^2 + \mathbb{E}[(f(\tilde{x}) - \mathbb{E}[f_{\theta}(\tilde{x})])^2] + \sigma_{\epsilon}^2$$

$$\mathbb{E}[(y - f_{\theta}(\tilde{x}))^2] = \text{Bias}[f_{\theta}(\tilde{x})]^2 + \text{Var}[f_{\theta}(\tilde{x})] + \sigma_{\epsilon}^2$$

Note: Every terms are >0

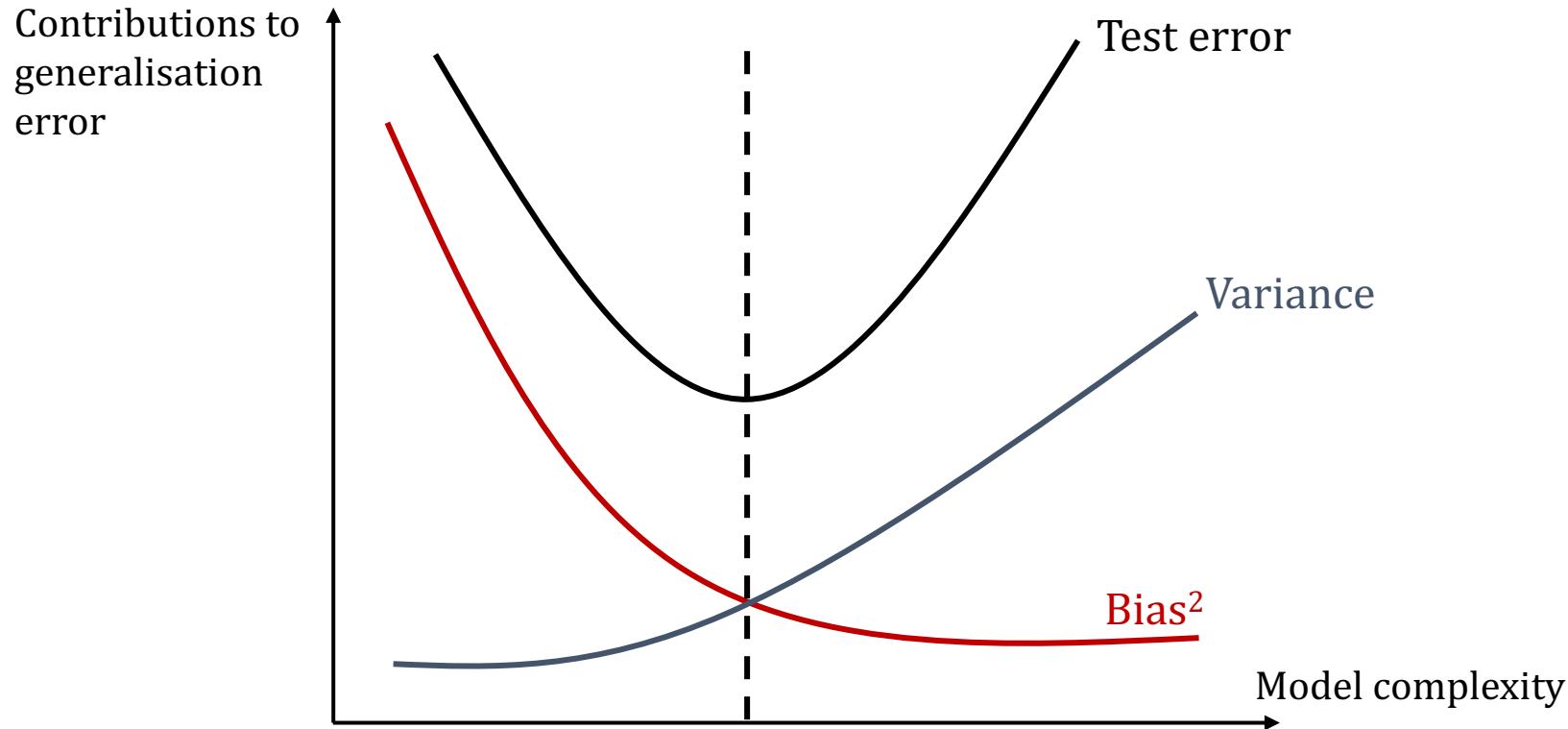
$\downarrow$   
**Modelling  
error**

$\downarrow$   
 Model  
variability

$\downarrow$   
 Irreducible  
error

- Note that a similar expression holds for classification

- Usually, “simple” models have large bias and low variance while “complex” ones have large variance and small bias\*



# Example on polynomial regression

Introduction to ML

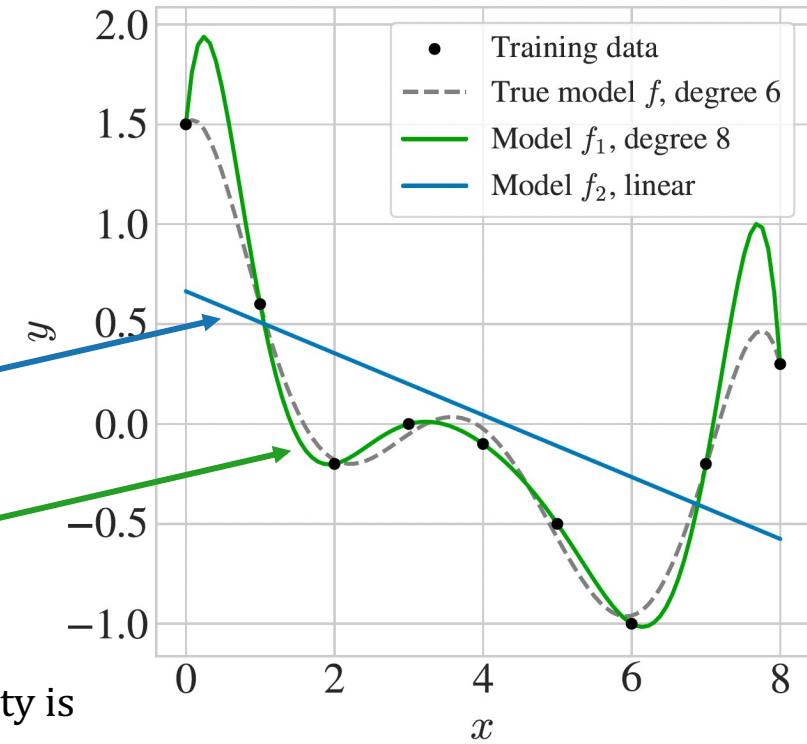
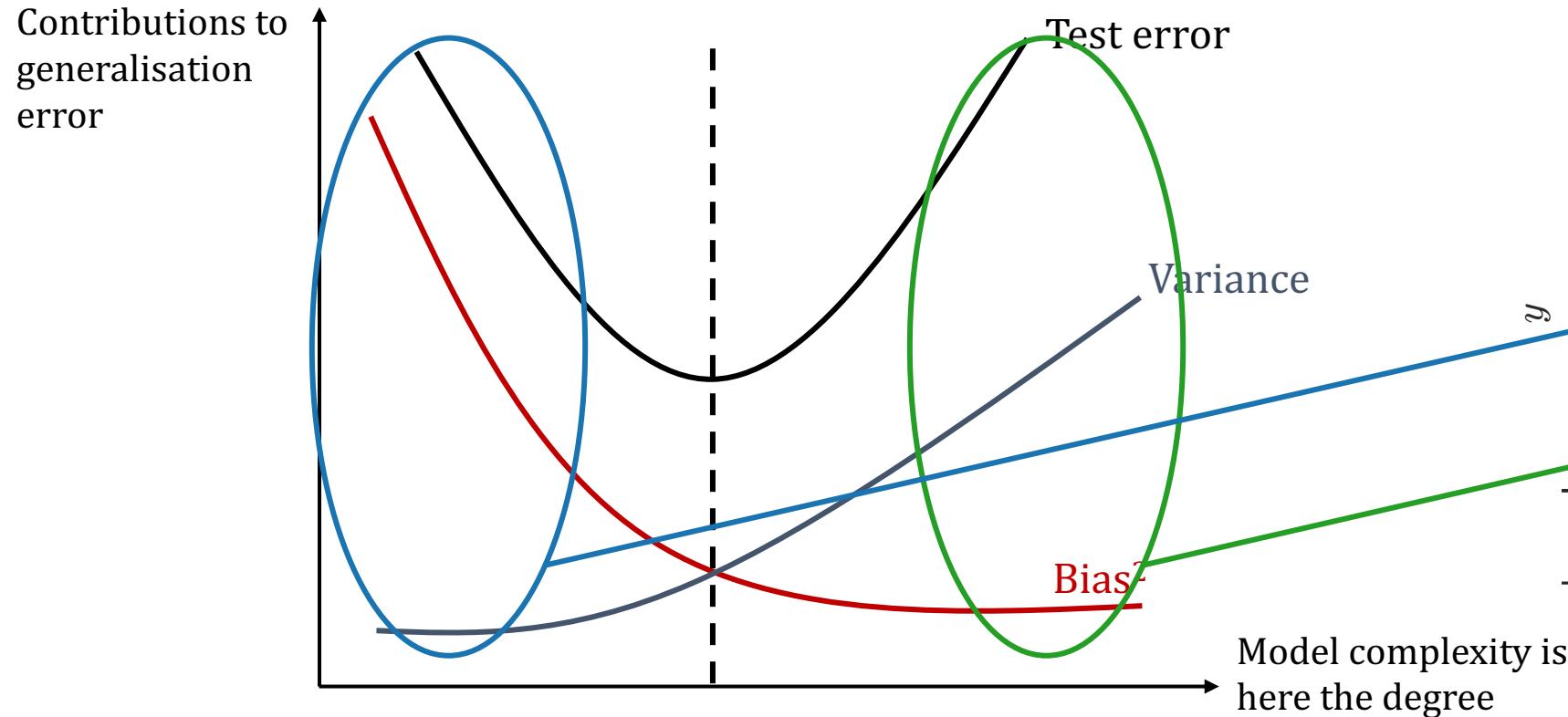
Our first model

Supervised learning theory

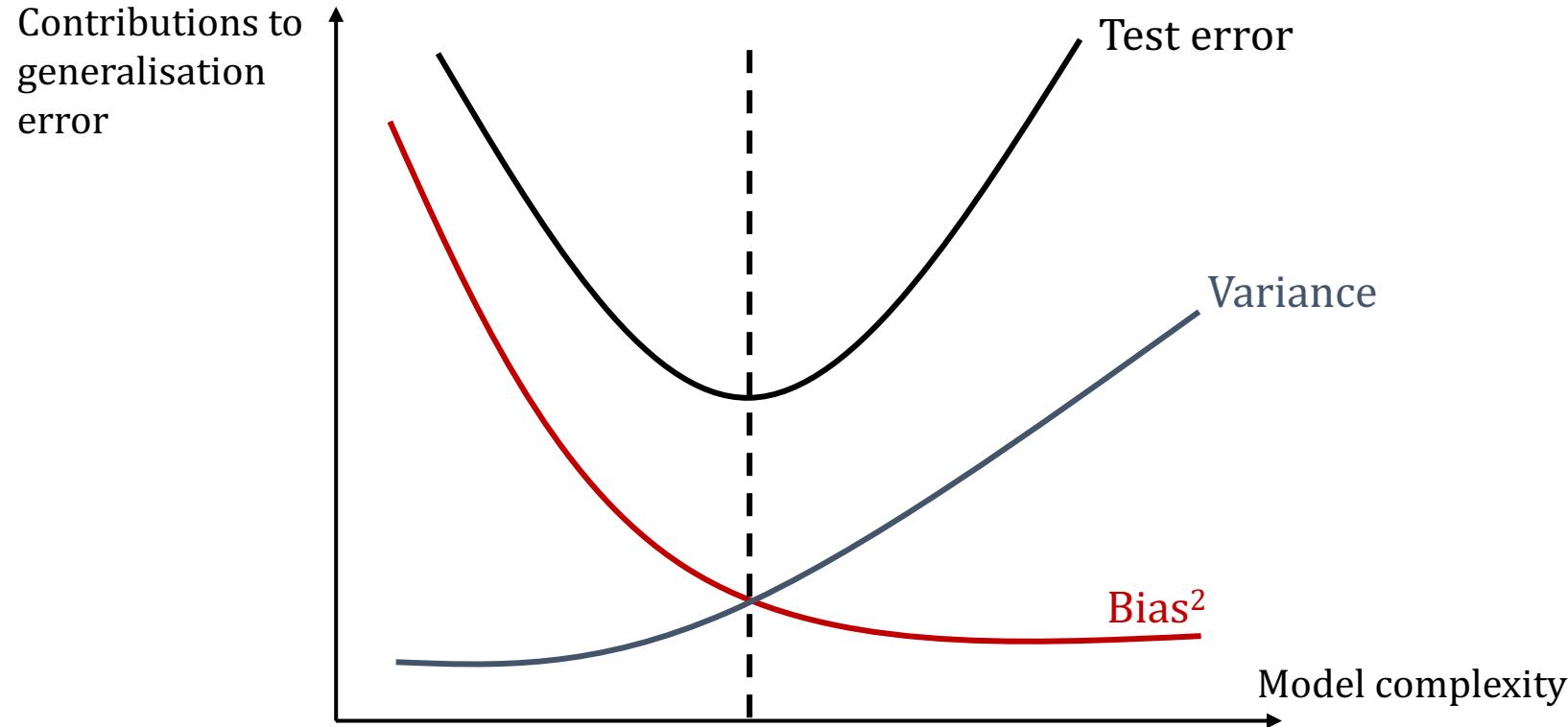
Other basic models

Deep learning models

- Usually, “simple” models have large bias and low variance while “complex” ones have large variance and small bias



- Usually, “simple” models have large bias and low variance while “complex” ones have large variance and small bias\*



Models need to be built such that they are not too flexible to fit the noise in the data but also not too restrictive to avoid bias

\* Heavily overparametrized models in fact go beyond this view and the test error decreases again in a **double descent** phenomenon (see [Belkin+18](#) and [Nakkiran+19](#))

# Model complexity and generalisation error

33

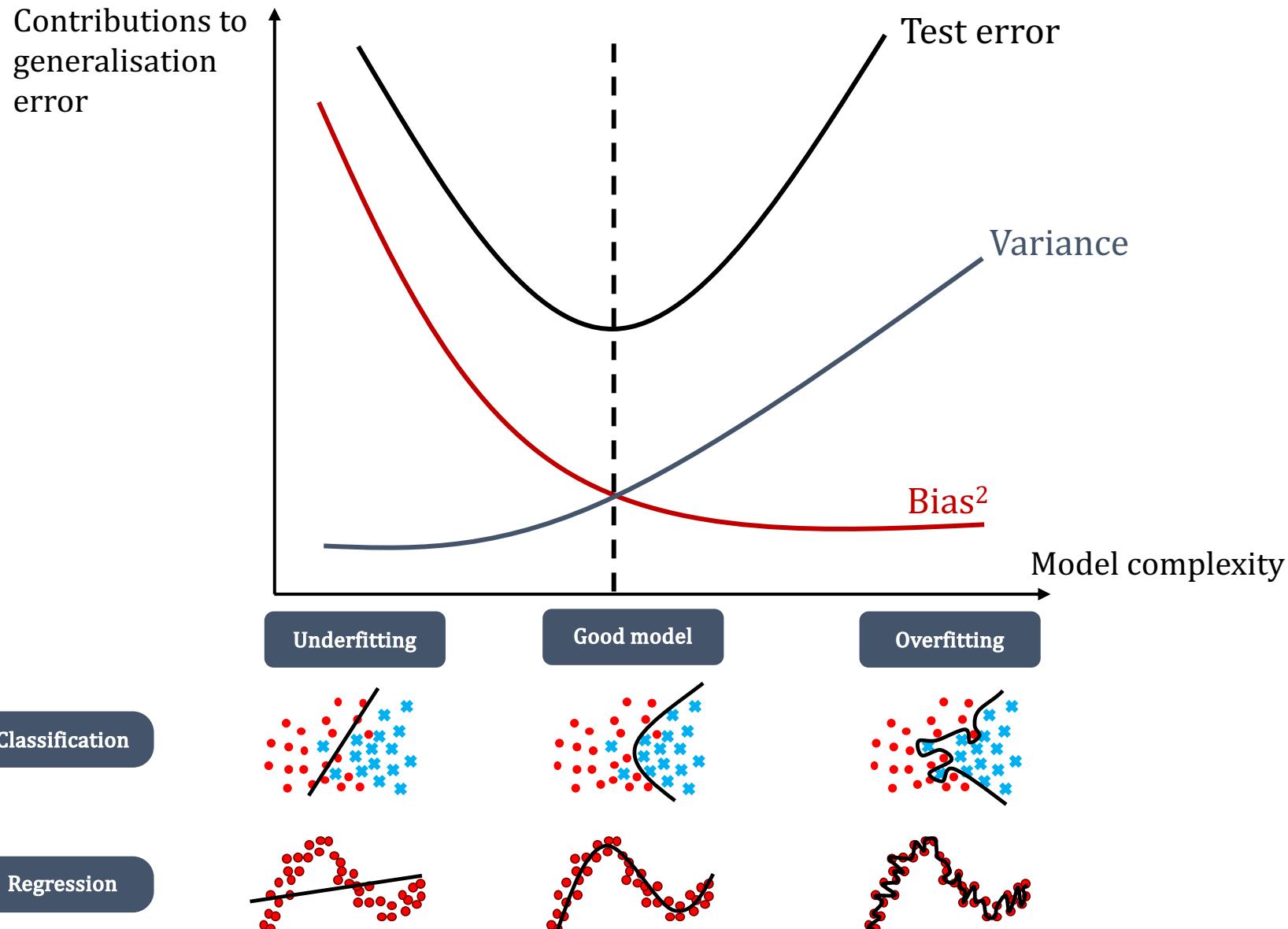
Introduction to ML

Our first model

Supervised learning theory

Other basic models

Deep learning models



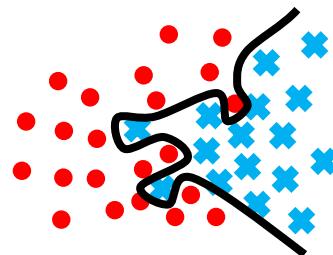
Classification

Regression

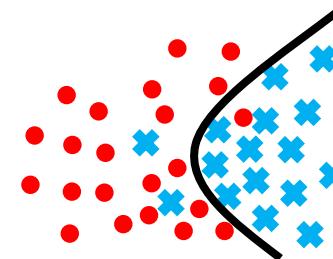
- To measure if we really learnt something useful in supervised learning in practice, we use a **test dataset** that the model has never seen but for which we know the labels and check that  $R_{\text{train}}(\hat{\theta})$  and  $R_{\text{test}}(\hat{\theta})$  are of the same order
- Based on the previous view, there are two regimes where things could go wrong: high variance and low bias models vs high bias and low variance models, respectively defining **overfitting** and **underfitting**

Classification

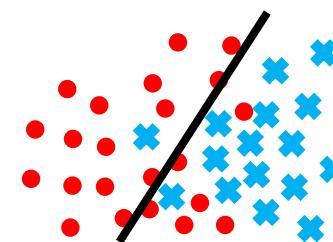
Overfitting



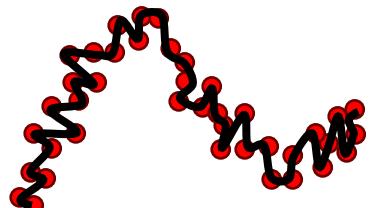
Good model



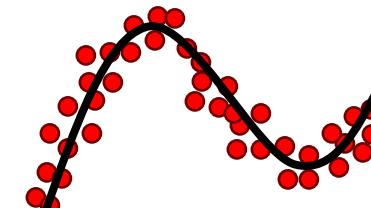
Underfitting



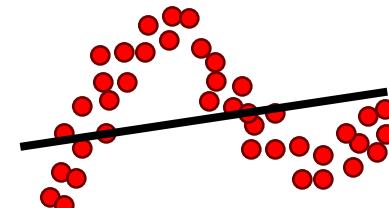
Regression



$$R_{\text{train}}(\hat{\theta}) \ll R_{\text{test}}(\hat{\theta})$$



$$R_{\text{train}}(\hat{\theta}) \approx R_{\text{test}}(\hat{\theta})$$



$$R_{\text{train}}(\hat{\theta}) \approx R_{\text{test}}(\hat{\theta}) \text{ but large}$$

# Some fixes to under and over fitting

35

Introduction to ML

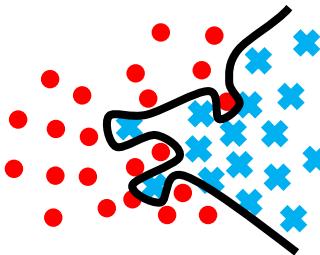
Our first model

Supervised learning theory

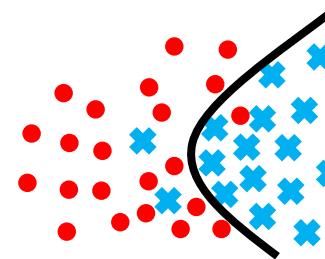
Other basic models

Deep learning models

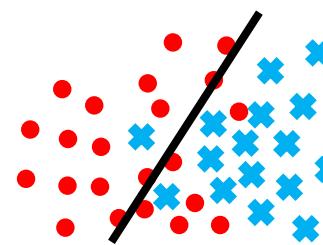
## Overfitting



## Good model

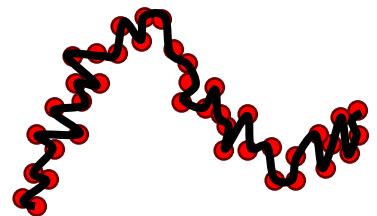


## Underfitting

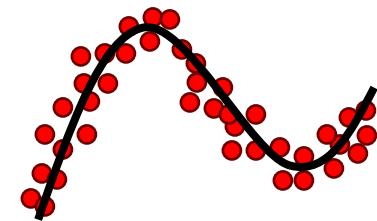


## Classification

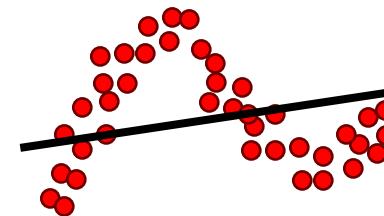
## Regression



$$R_{\text{train}}(\hat{\theta}) \ll R_{\text{test}}(\hat{\theta})$$



$$R_{\text{train}}(\hat{\theta}) \approx R_{\text{test}}(\hat{\theta})$$



$$R_{\text{train}}(\hat{\theta}) \approx R_{\text{test}}(\hat{\theta}) \text{ but large}$$

## Possible fixes

- Add more data
- Remove features
- Stop the training earlier
- Add regularisation

- You did a good job!

- Try a more complex model
- Train your model longer

# A gentle introduction to regularisation

36

Introduction to ML

Our first model

Supervised learning theory

Other basic models

Deep learning models

- One generic way to deal with **overfitting** is by incorporating some idea of **regularity** about the parameters: this is called explicit **regularization** and appears in the optimization problem as a constraint on parameter space

$$\hat{\theta} = \operatorname{argmin}_{\theta \in \Theta} R(\theta)$$

Unregularized  
optimization

$$\hat{\theta} = \operatorname{argmin}_{\theta \in \Theta} R(\theta) \text{ s. t. } P(\theta) \leq \epsilon$$

Regularized  
optimization

- The function  $P(\theta)$  is called **penalty function** and the regularized problem can in fact be written equivalently as (by Lagrange duality)

$$\hat{\theta} = \operatorname{argmin}_{\theta \in \Theta} R(\theta) + \lambda P(\theta)$$

1.  $\lambda$  is a **regularization** (hyper)**parameter**
2.  $P(\theta)$  can take different forms depending on the penalty we chose to impose on the set of parameters

- Common penalty functions are the  $L_p$ -norms with  $p = 0, 1$  or  $2$

$$P(\theta) = \|\theta\|_p$$

# Regularisation as a smoothness constraint

37

Introduction to ML

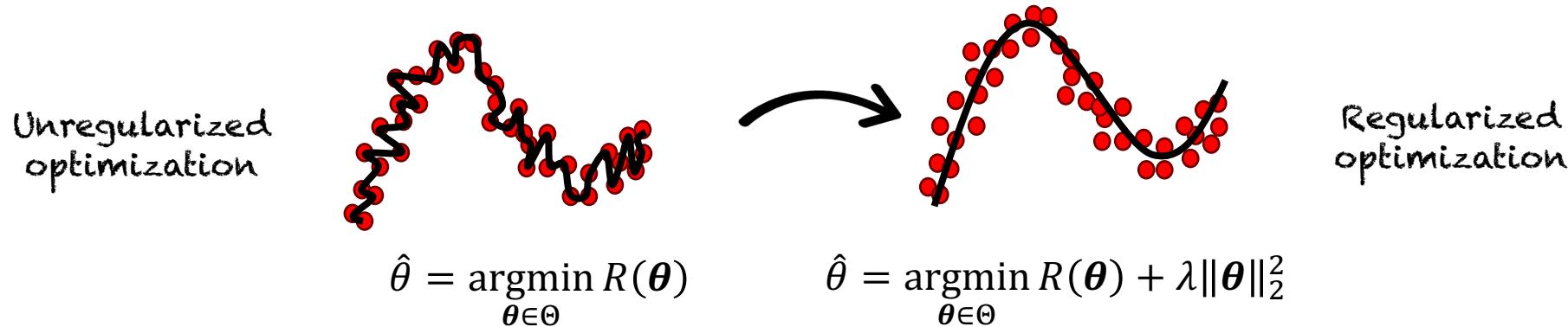
Our first model

Supervised learning theory

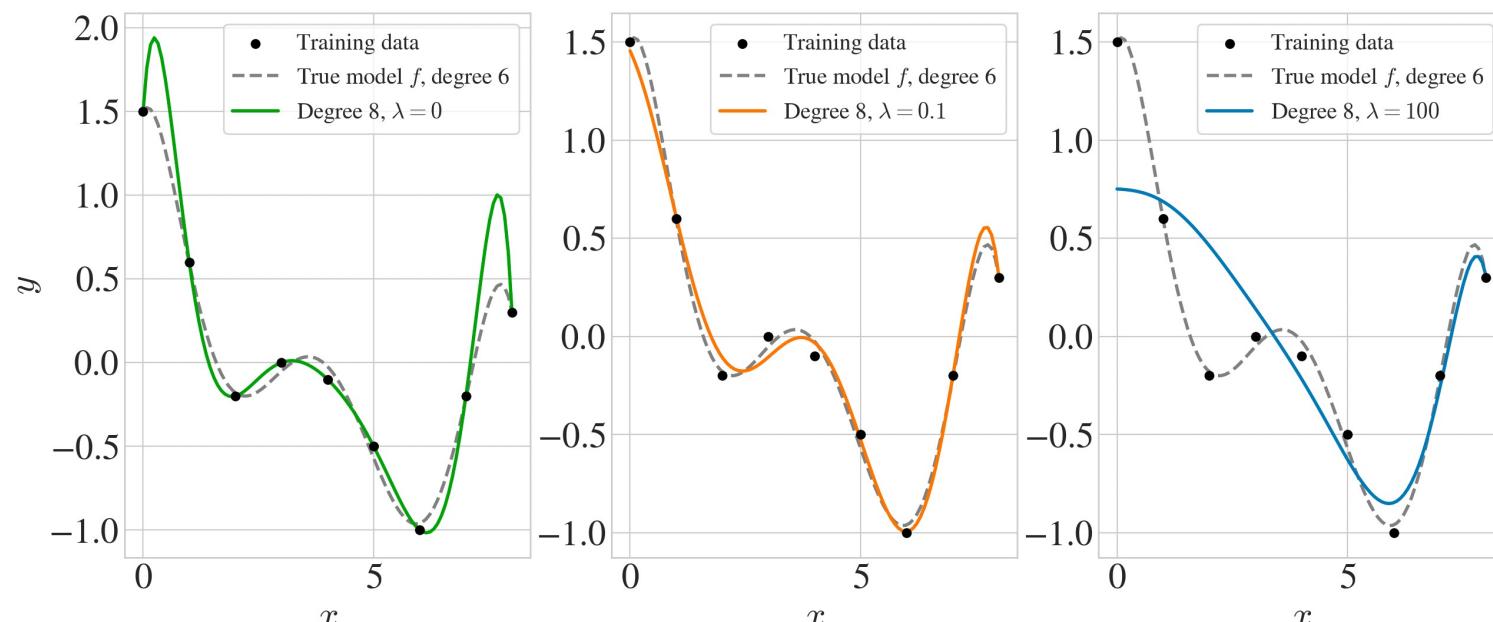
Other basic models

Deep learning models

- By solving the regularized optimization with **an appropriate value of  $\lambda$** , we can reduce the overfitting



- The unregularized models lacks smoothness** which is introduced by the  $L_2$  penalty
- The additional term here penalizes large weight values and **reduces the variance** of the estimator



# A probabilistic view of optimisation

38

Introduction to ML

Our first model

Supervised learning theory

Other basic models

Deep learning models

- Let us consider that  $y^{(i)} = f_{\theta}(\mathbf{x}^{(i)}) + \epsilon^{(i)}$ , meaning our model is correct up to an error with  $\epsilon^{(i)}$  that we assume i.i.d.
- In this, case we can write the **likelihood** of the data as

$$p(\mathbf{Y}|\mathbf{X}, \boldsymbol{\theta}) = \prod_i p(y^{(i)}|\mathbf{x}^{(i)}, \boldsymbol{\theta})$$

- Using the **Bayes theorem**, we can express the posterior probability distribution of the parameters given the dataset as

$$p(\boldsymbol{\theta}|\mathbf{X}) \propto p(\mathbf{Y}|\mathbf{X}, \boldsymbol{\theta})p(\boldsymbol{\theta})$$

- Assuming a Gaussian likelihood  $\epsilon^{(i)} \sim \mathcal{N}(0, \sigma_{\epsilon}^2)$ ,  $\theta_i \sim \mathcal{N}(0, \sigma_{\theta}^2)$ , and taking the log of this expression yields

$$\log p(\boldsymbol{\theta}|\mathbf{X}) = -\frac{1}{2\sigma_{\epsilon}^2} \sum_i (y^{(i)} - \hat{y}^{(i)})^2 - \frac{1}{2\sigma_{\theta}^2} \sum_i \theta_i^2 + \text{cst.}$$

- Finally,

$$\max_{\boldsymbol{\theta}} \log p(\boldsymbol{\theta}|\mathbf{X}) \Leftrightarrow \min_{\boldsymbol{\theta}} \sum_i (y^{(i)} - \hat{y}^{(i)})^2 + \lambda \|\boldsymbol{\theta}\|^2, \quad \text{with } \lambda = \sigma_{\epsilon}^2 / \sigma_{\theta}^2$$

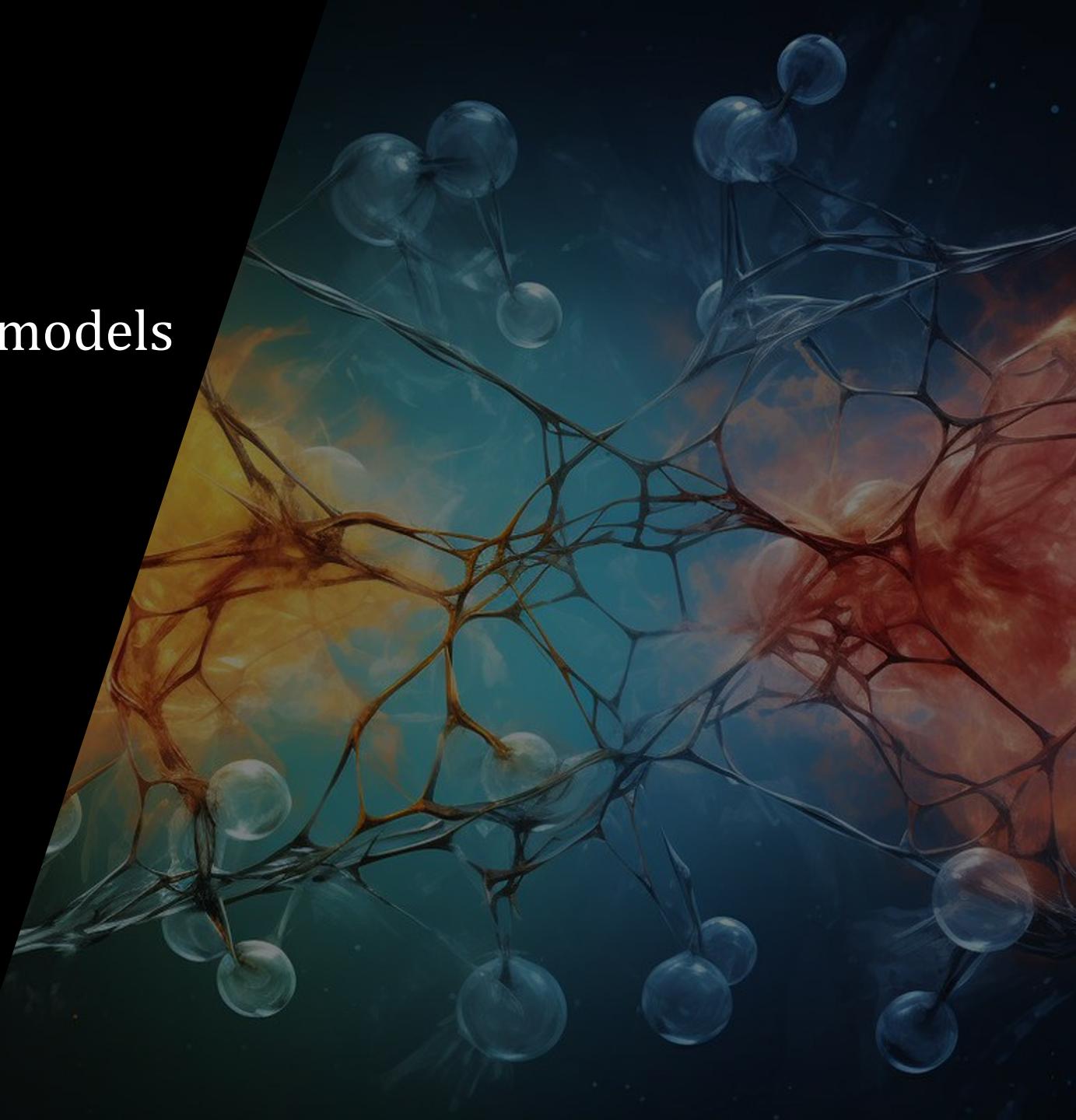


The maximum a posteriori estimator under Gaussian likelihood and Gaussian prior is equivalent to minimizing the square loss with an  $L_2$  penalty. The **prior plays the role of the regularization**, and the **square-loss comes from the Gaussian error** assumption.

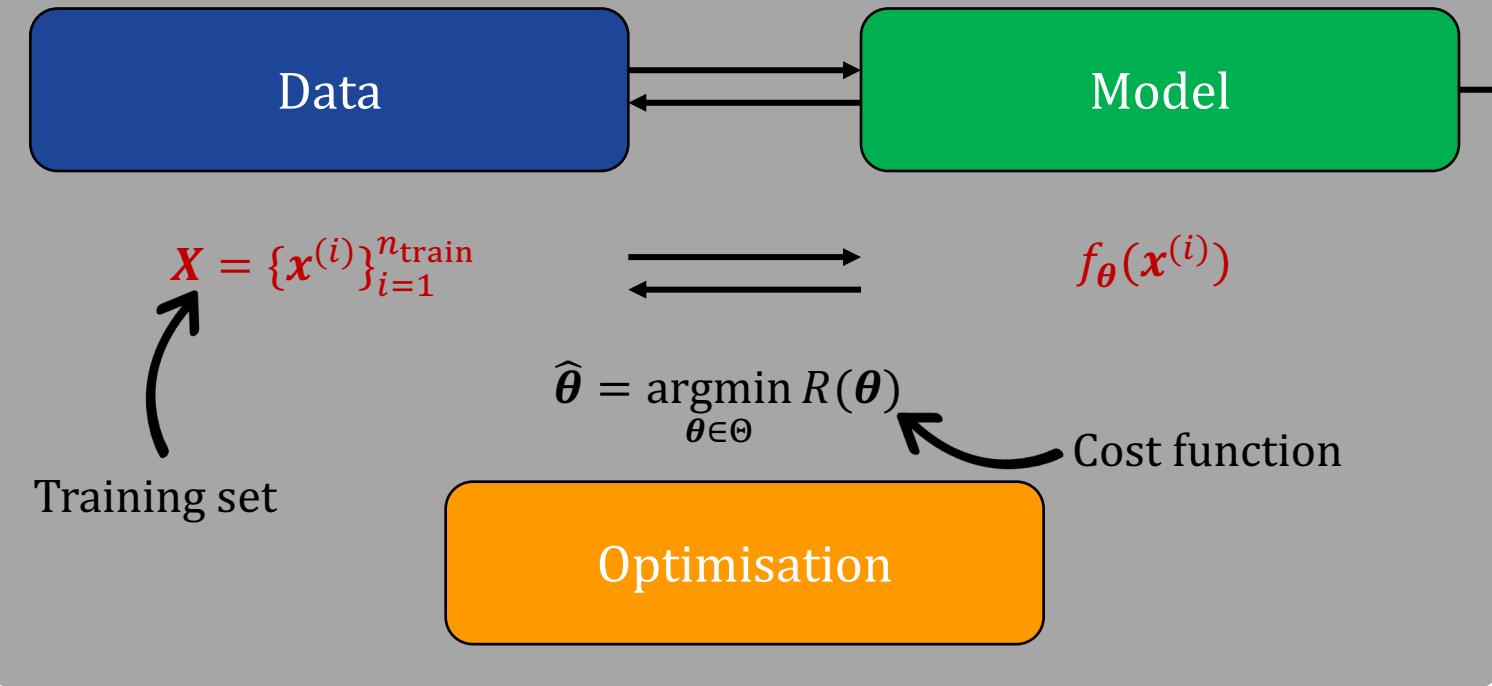
# Other basic supervised learning models

Contents :

- *A first non-linear model: decision trees*
- *Ensembling: bagging and boosting*
- *Random forest and boosted trees*
- *Feed-forward neural networks: definition*
- *Backpropagation of errors*



## Training phase



## Generalisation phase

Prediction on new data

$$f_{\theta}(\tilde{x}^{(i)})$$

$\tilde{X}$  Test set

Some notations:

$$x^{(i)} \in \mathbb{R}^d \quad \theta \in \Theta \subset \mathbb{R}^p$$

Usually,

$$x^{(i)} \sim \rho$$



Three main building blocks of *most* ML systems: some **data** feed a **parametrised model**. The best set of parameters are obtained through the **optimisation** of a cost function. The model can then be assessed on **new data**.

# A classification task: the XOR dataset

Introduction to ML

Our first model

Supervised learning theory

Other basic models

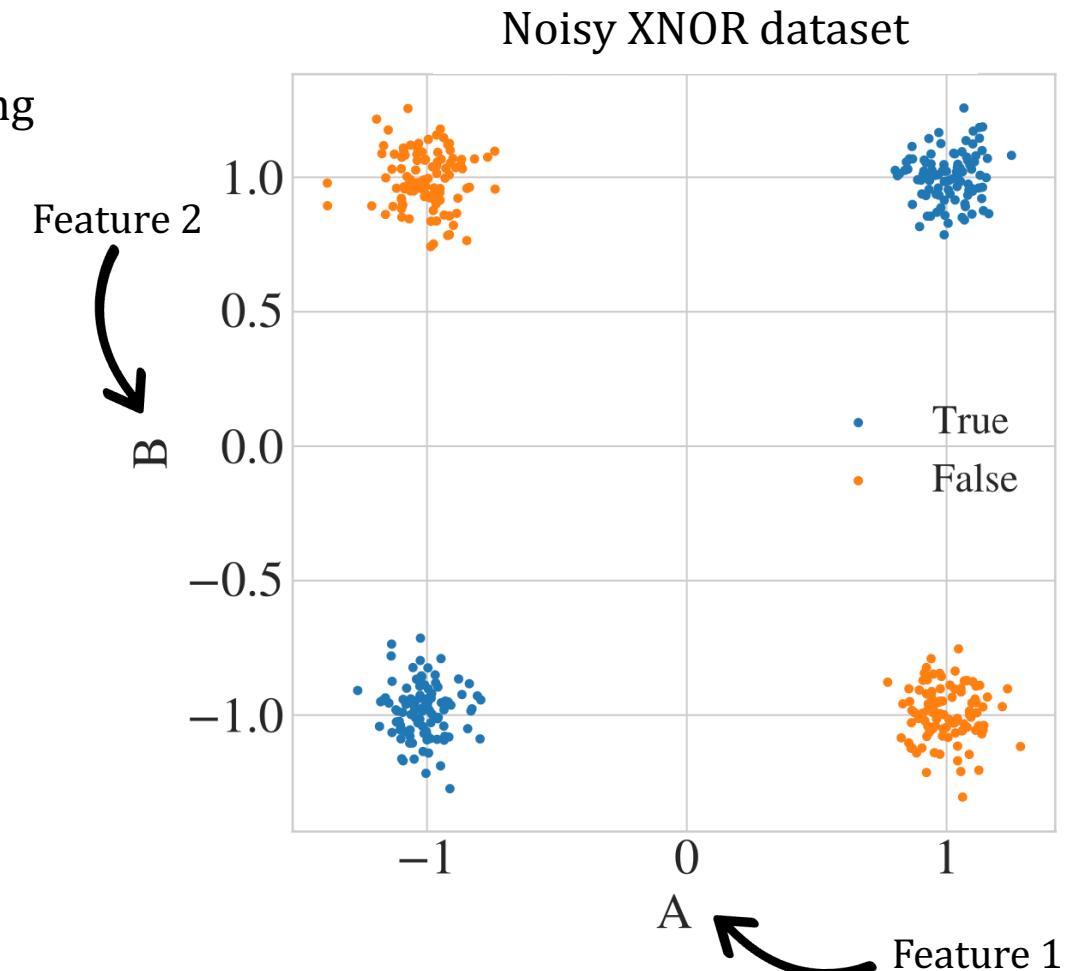
Deep learning models

## Decision trees

- Consider a **classification task** on an artificial dataset replicating the XNOR function

A	B	XNOR
True	True	True
True	False	False
False	True	False
False	False	True

- Data are  $n = 500$  couples  $(x^{(i)}, y^{(i)}) \Rightarrow$  **Supervised learning**
- The target variable  $y \in \{-1, 1\}$  is discrete  $\Rightarrow$  **Classification**
- A linear classification would not be able to learn such a function\*
- Decision trees to the rescue!



\*In fact, an alternative (not discussed here) would be to use **kernels** to find a higher dimensional space in which the data separates linearly and then use a linear classifier.

# Intuition behind decision trees

Introduction to ML

Our first model

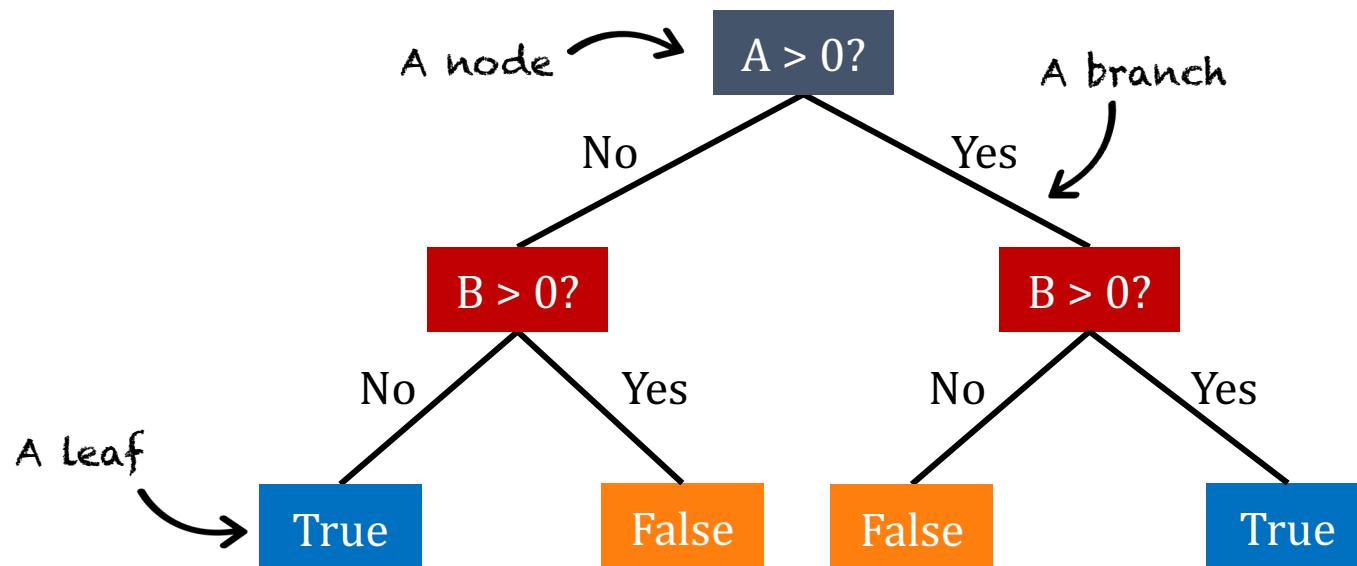
Supervised learning theory

Other basic models

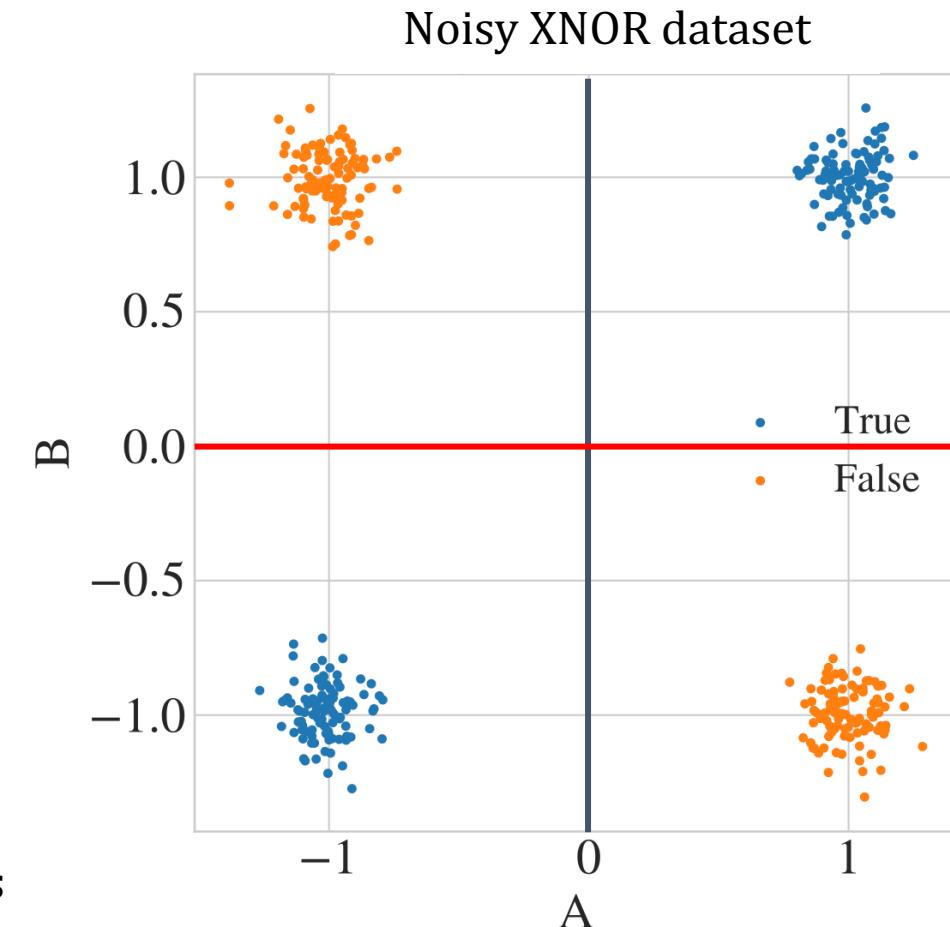
Deep learning models

## Decision trees

- Decision trees incrementally ask questions about the features



- All root and inner nodes question the value of a feature, and branches split the dataset into **different regions to which a datapoint can belong uniquely**



# A more formal view of decision trees

Introduction to ML

Our first model

Supervised learning theory

Other basic models

Deep learning models

## Decision trees

- More formally, at a given node in parent region  $R_t$  asking a question about the  $j^{\text{th}}$  feature, we create two regions:

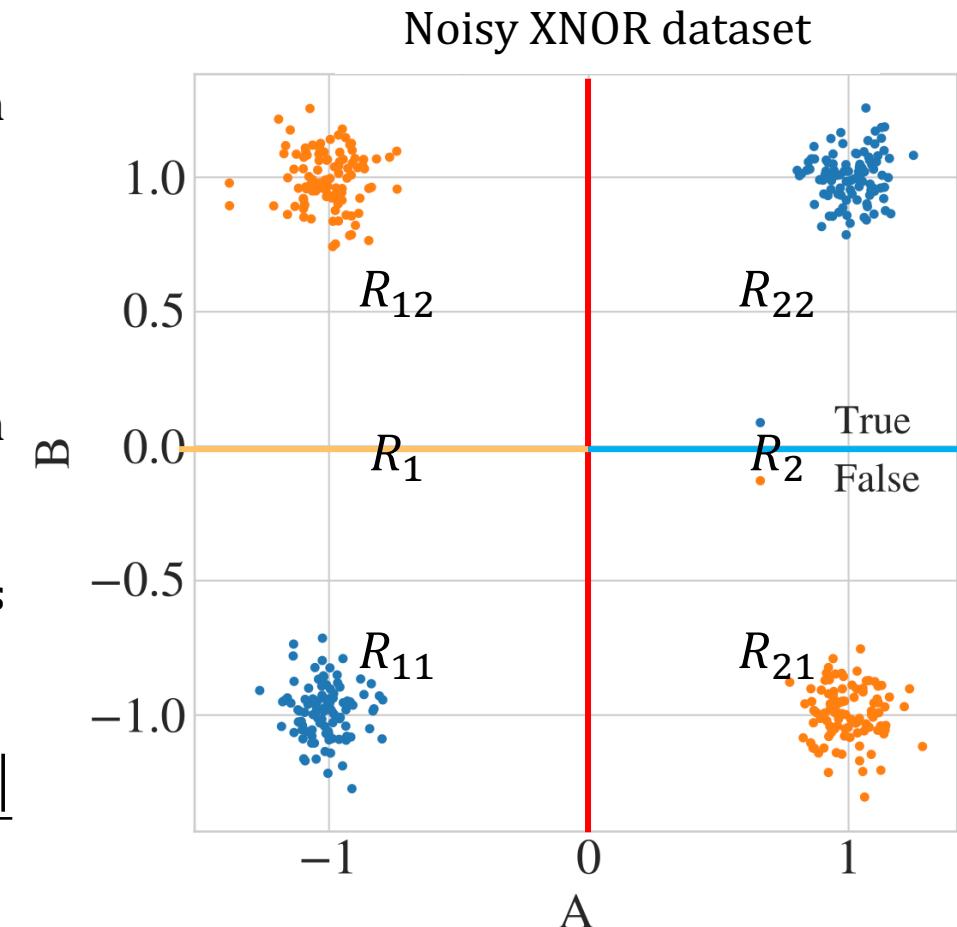
$$R_1 = \{x \mid x_j < \alpha_{t1}^j, x \in R_t\}$$

$$R_2 = \{x \mid x_j \geq \alpha_{t2}^j, x \in R_t\}$$

- The parameters  $\theta$  of decision trees are the threshold values at each nodes (the sequence of  $\alpha$ )
- To find the optimal values, decision tree algorithms minimise a loss function in each region, and **a good choice is the cross-entropy**

$$L(R_i) = - \sum_{k=1}^q \rho_k^i \log_2 \rho_k^i,$$

$$\rho_k^i = \frac{|\{x^{(j)} \mid x^{(j)} \in R_i, y^{(j)} = k\}|}{|\{x^{(j)} \mid x^{(j)} \in R_i\}|}$$



can handle categorical values, easy to interpret, fast to compute, both regression and classification



Shallow trees: weak learners (underfit) and high bias, deep trees: high variance estimators (overfit)

# A more formal view of decision trees

44

Introduction to ML

Our first model

Supervised learning theory

Other basic models

Deep learning models

## Decision trees

- More formally, at a given node in parent region  $R_t$  asking a question about the  $j^{\text{th}}$  feature, we create two regions:

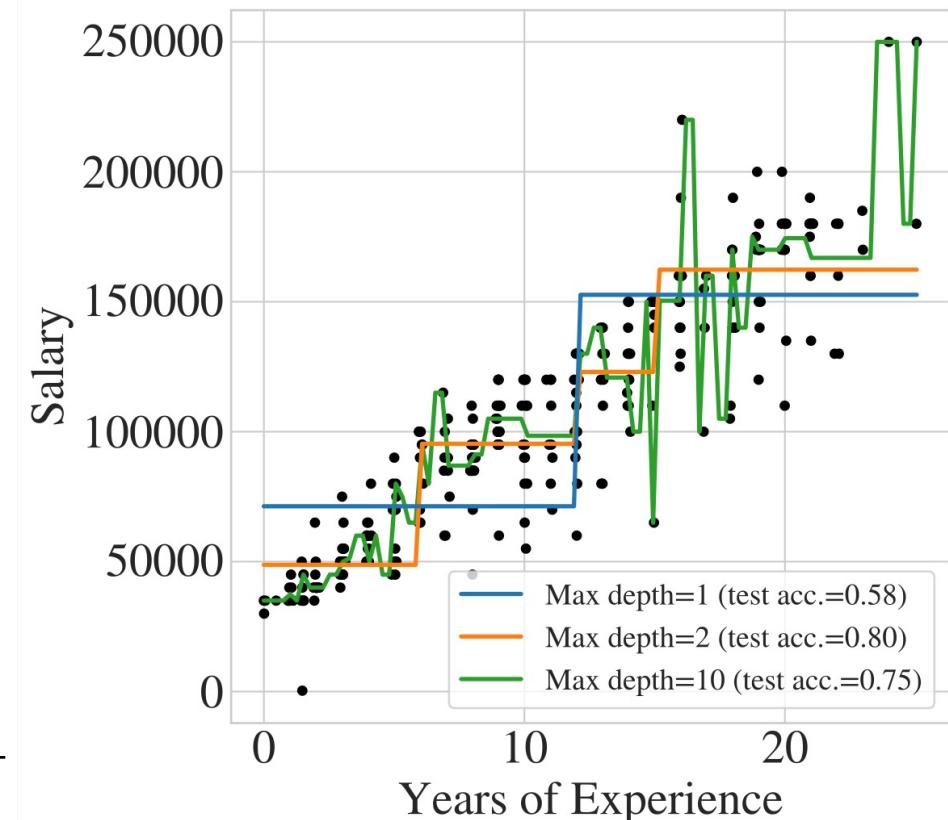
$$R_1 = \{x \mid x_j < \alpha_{t1}^j, x \in R_t\}$$

$$R_2 = \{x \mid x_j \geq \alpha_{t2}^j, x \in R_t\}$$

- The parameters  $\theta$  of decision trees are the threshold values at each nodes (the sequence of  $\alpha$ )
- To find the optimal values, decision tree algorithms minimise a loss function in each region, and **a good choice is the cross-entropy**

$$L(R_i) = - \sum_{k=1}^q \rho_k^i \log_2 \rho_k^i,$$

$$\rho_k^i = \frac{|\{x^{(j)} \mid x^{(j)} \in R_i, y^{(j)} = k\}|}{|\{x^{(j)} \mid x^{(j)} \in R_i\}|}$$



can handle categorical values, easy to interpret, fast to compute, both **regression** and classification



Shallow trees: weak learners (underfit) and high bias, deep trees: high variance estimators (overfit)

# A more formal view of decision trees

45

Introduction to ML

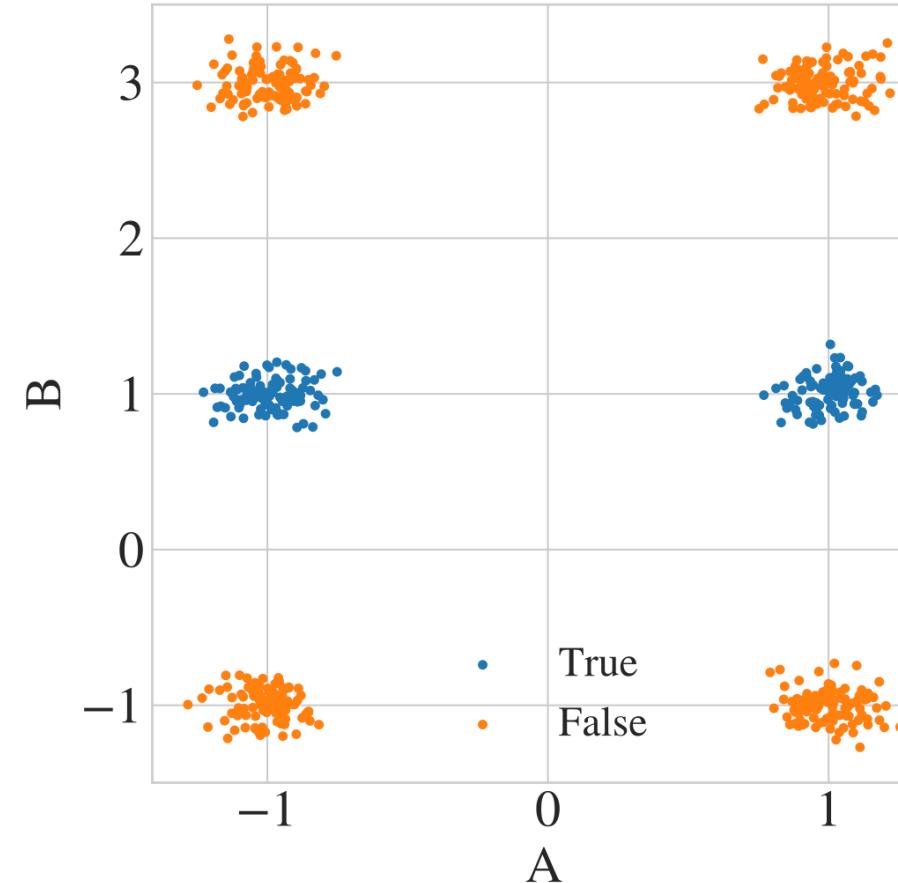
Our first model

Supervised learning theory

Other basic models

Deep learning models

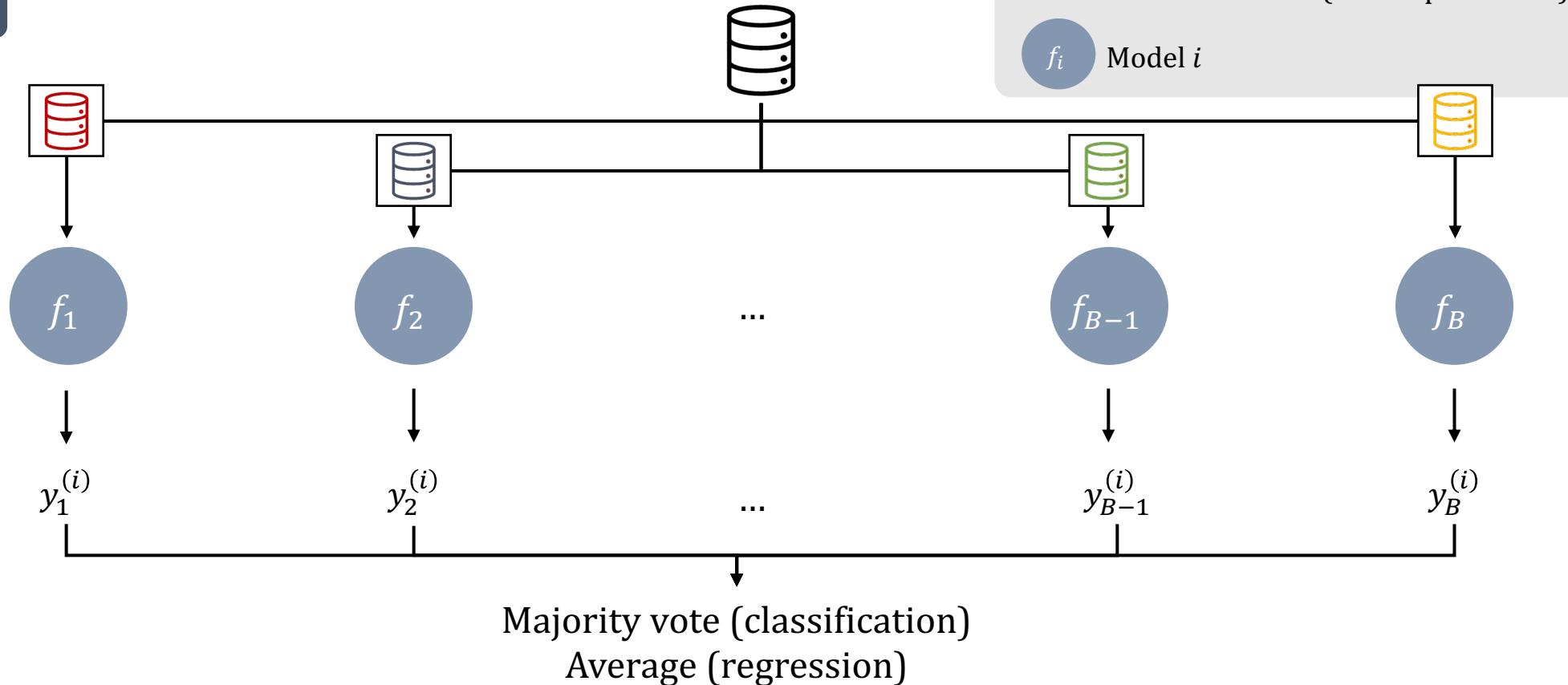
## Decision trees



Practice: Can you build a decision tree solving the binary classification problem?

- To circumvent the problems that can have weak learners like decision trees, **ensembling methods** were proposed

## Bagging



- To reduce the variance, models need to be **uncorrelated**: this is achieved by using **random sampling of the dataset**

# Ensembling methods: random forest

47

Introduction to ML

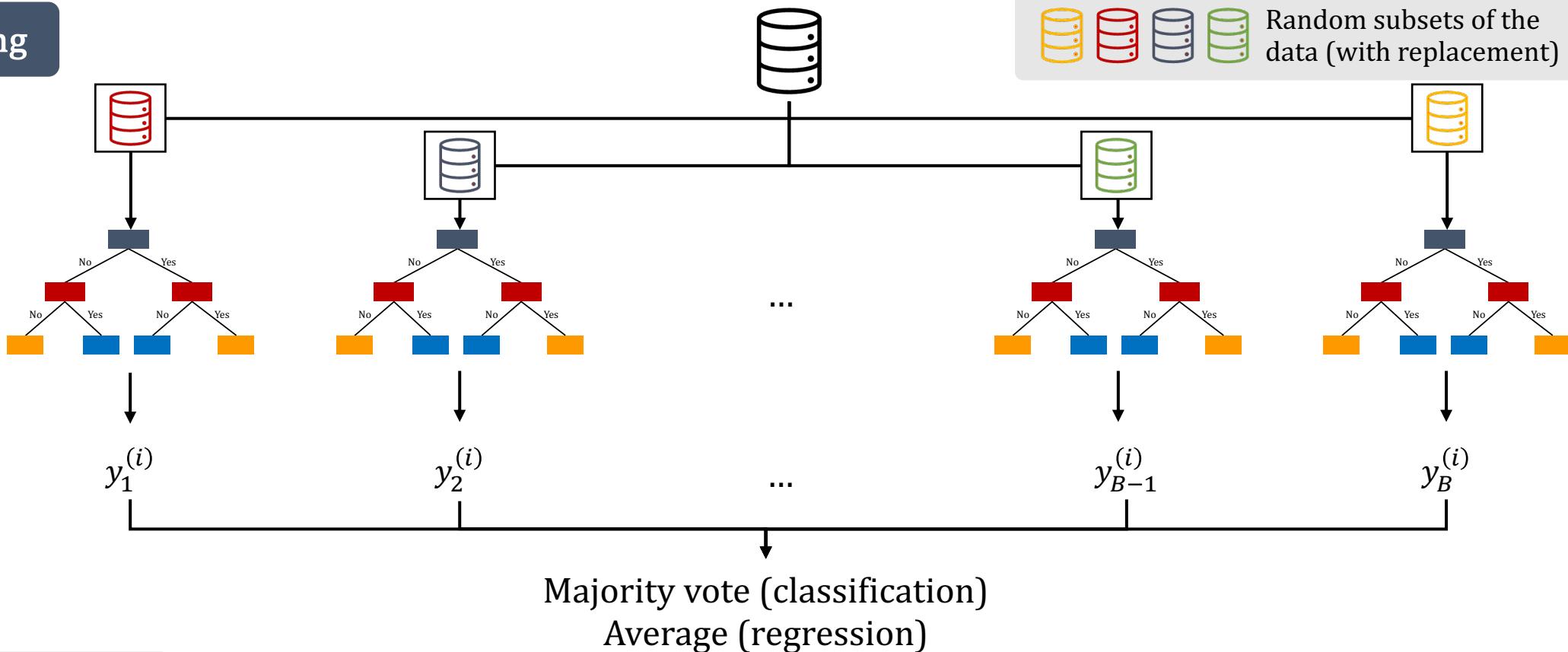
Our first model

Supervised learning theory

Other basic models

Deep learning models

## Bagging



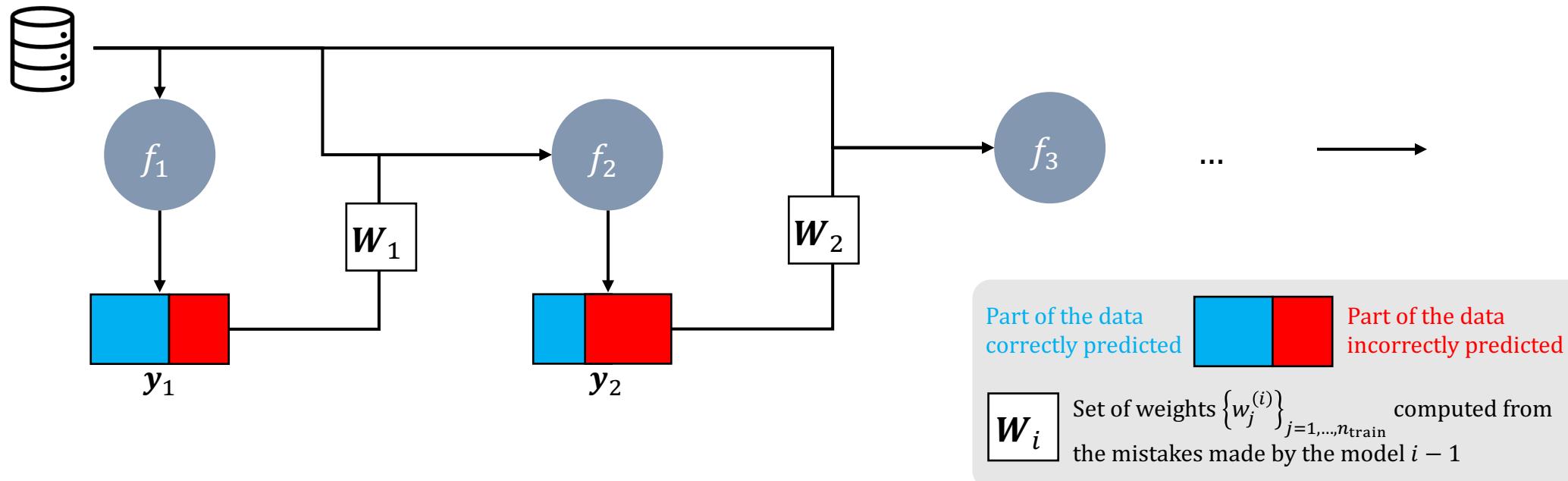
## Random Forest

= Bagging of decision trees + feature bagging

- + Can still handle categorical values, both regression and classification, **reduced variance**
- More expensive to compute (need to train  $B$  trees instead of one), harder to interpret

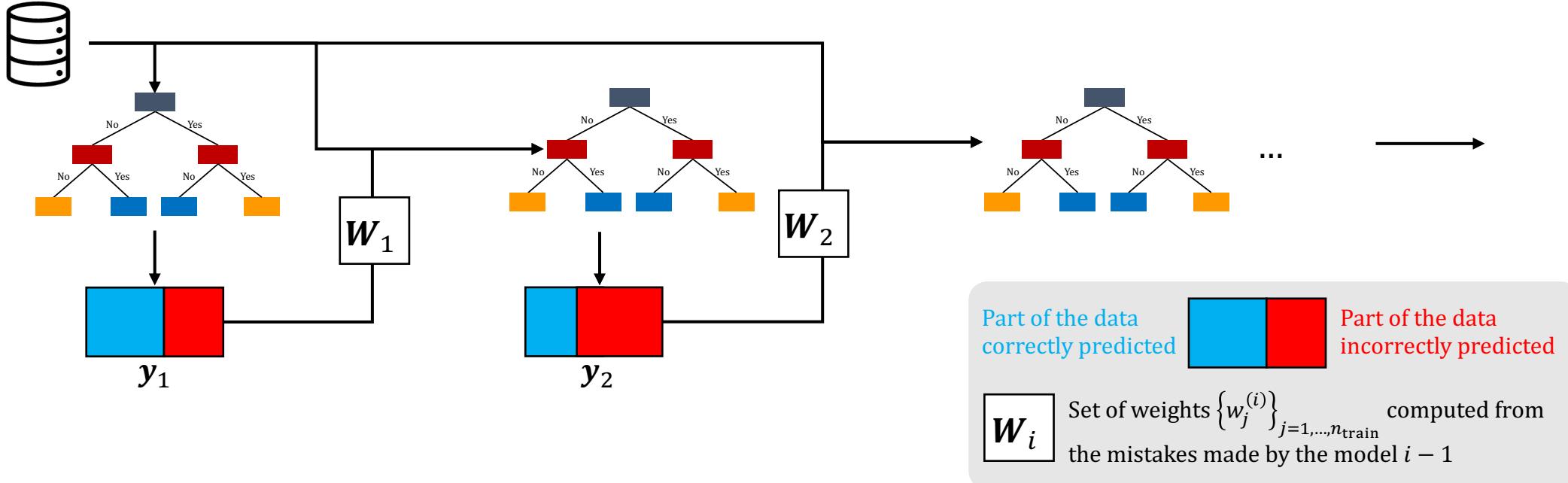
- While bagging trains high-variance models in parallel to reduce the variance of the combined estimate, **boosting trains high-bias models sequentially to reduce the overall bias**

## Boosting



- Successive learners  $f_i$  are fed by data  $X_i$ , a weighted version of the initial dataset  $X$ , **giving more weights to the errors committed by the previous model  $f_{i-1}$**
- The output is, as in bagging, **a linear combination of all the learners**
- The choice of weighting and training depends on the algorithm and context (see [Adaboost](#) or [gradient boosting](#))

## Boosting



**Boosted trees**

= Boosting of decision trees

- + Both regression and classification (residuals or weighted classification error), **reduced bias**, good performances
- More expensive to compute, increased variance, subject to overfitting

# Comparison between DT, RF and GB

50

Introduction to ML

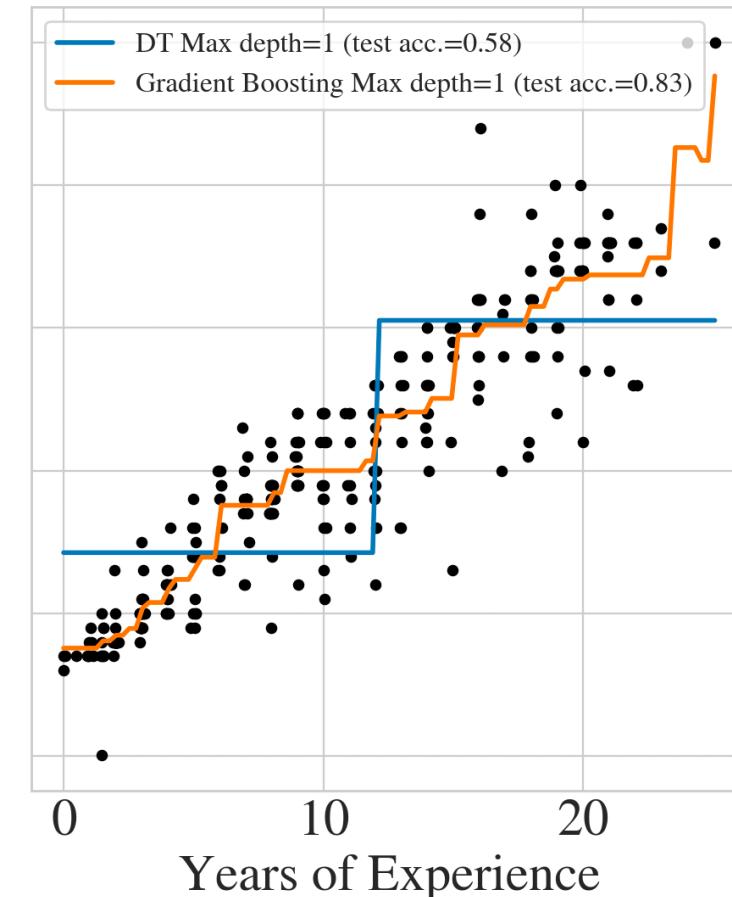
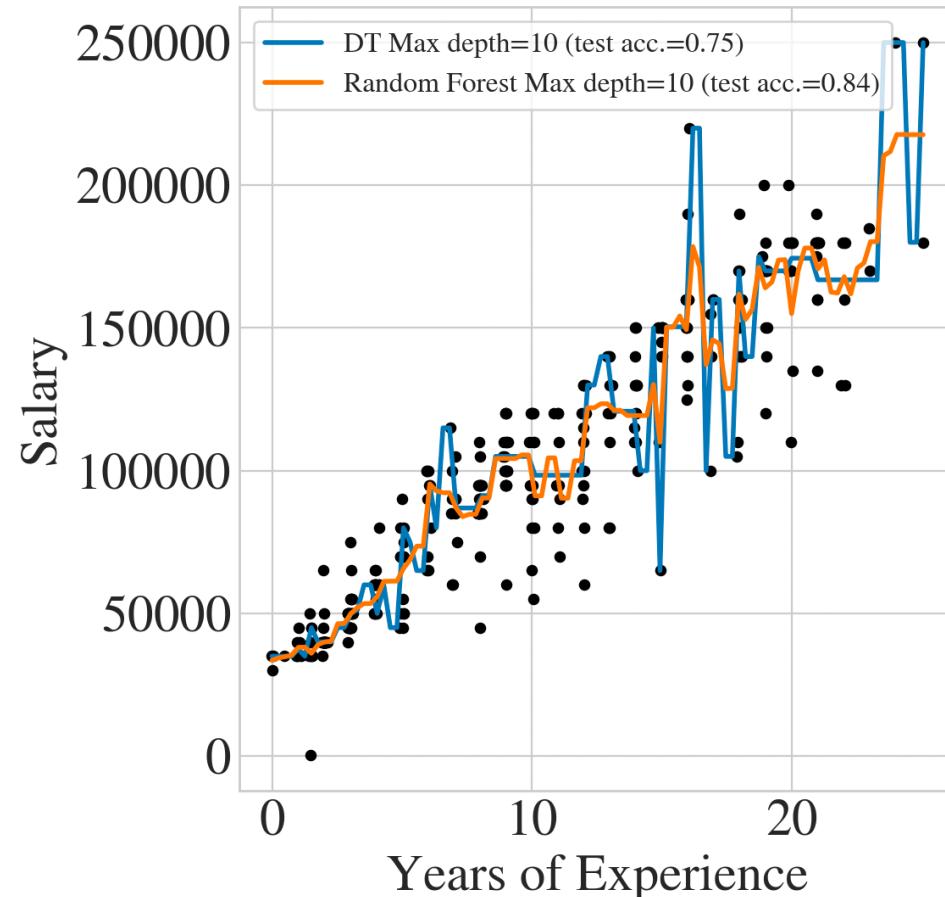
Our first model

Supervised learning theory

Other basic models

Deep learning models

Comparison on our regression problem



# Artificial neural networks: neurons

Introduction to ML

Our first model

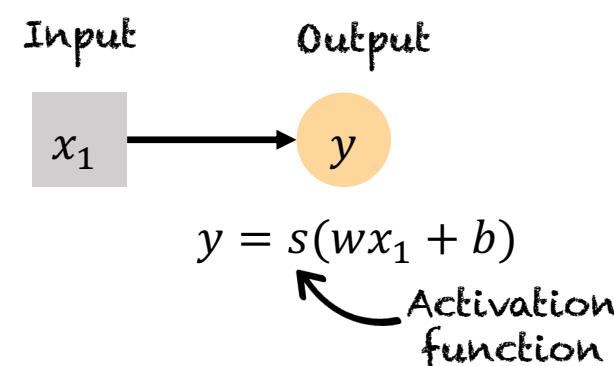
Supervised learning theory

Other basic models

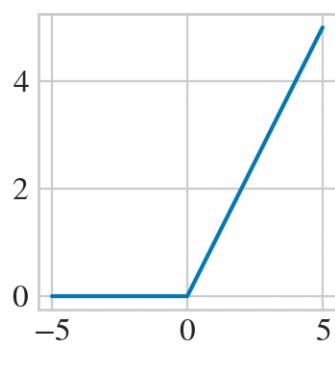
Deep learning models

## Artificial neural networks

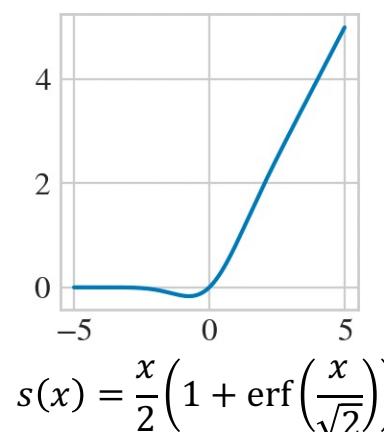
- Building block of neural networks: the **neuron**
- Made of three operations: it first multiplies the input by a **weight**, then adds a **bias**, and finally applies an **activation function**  $s$  to the result
- If  $s$  is the identity, you recognize the linear regression model with  $\theta_0 = b$  and  $\theta_1 = w$ . To represent non-linear functions,  $s$  must be non-linear



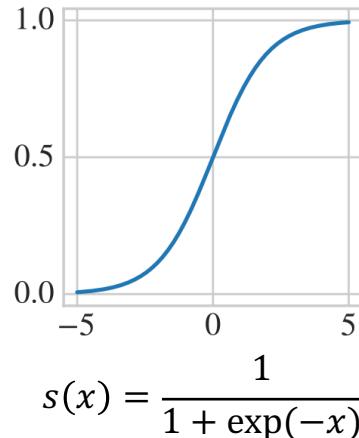
ReLU



GeLU



Sigmoid



- This model is motivated by biological neurons and the activation function mimics the activation or inhibition through **non-linear (and differentiable) functions**
- Can take different forms but some examples include the **ReLU or sigmoid functions**

# Artificial neural networks: more features

52

Introduction to ML

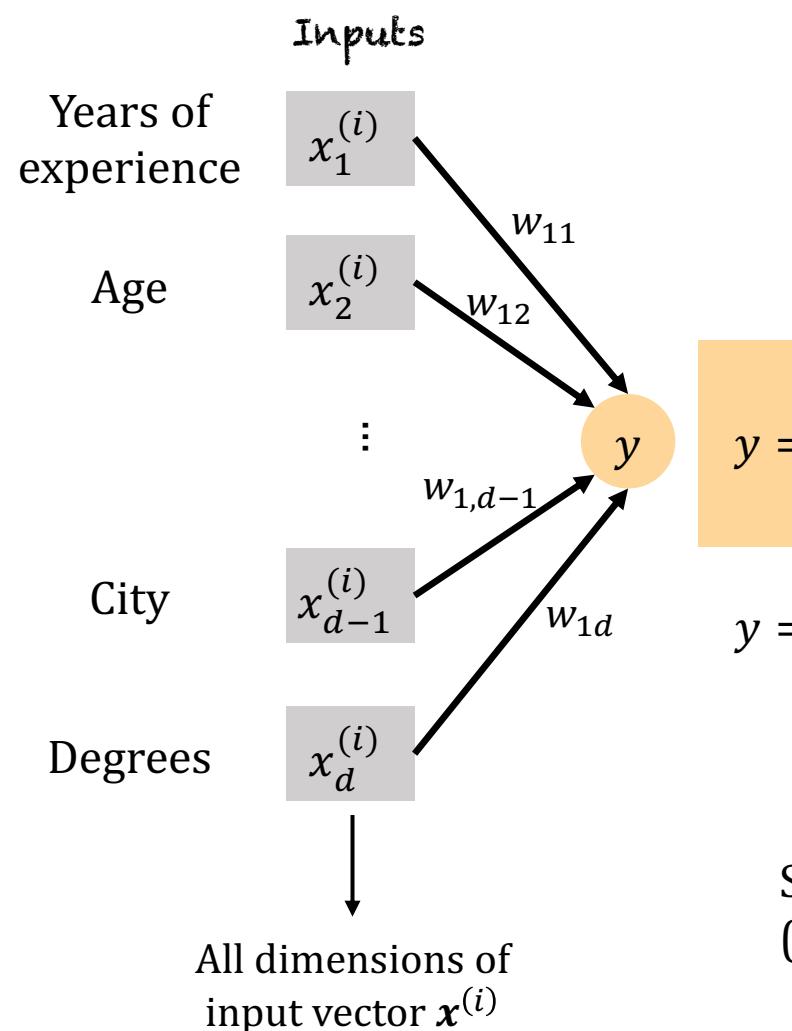
Our first model

Supervised learning theory

Other basic models

Deep learning models

## Artificial neural networks



Notation:  $w_{ij}$  is the weight linking the unit  $j$  of the first layer to the unit  $i$  of the next layer

### Output

$$y = s \left( \sum_{j=1}^d w_{1j} x_j^{(i)} + b \right)$$

$$y = s(\mathbf{w}_1 \cdot \mathbf{x}^{(i)} + b)$$

Such a network has  $p = d + 1$  parameters  
(as in multi-feature linear regression)

# More layers: the art of learning features

53

Introduction to ML

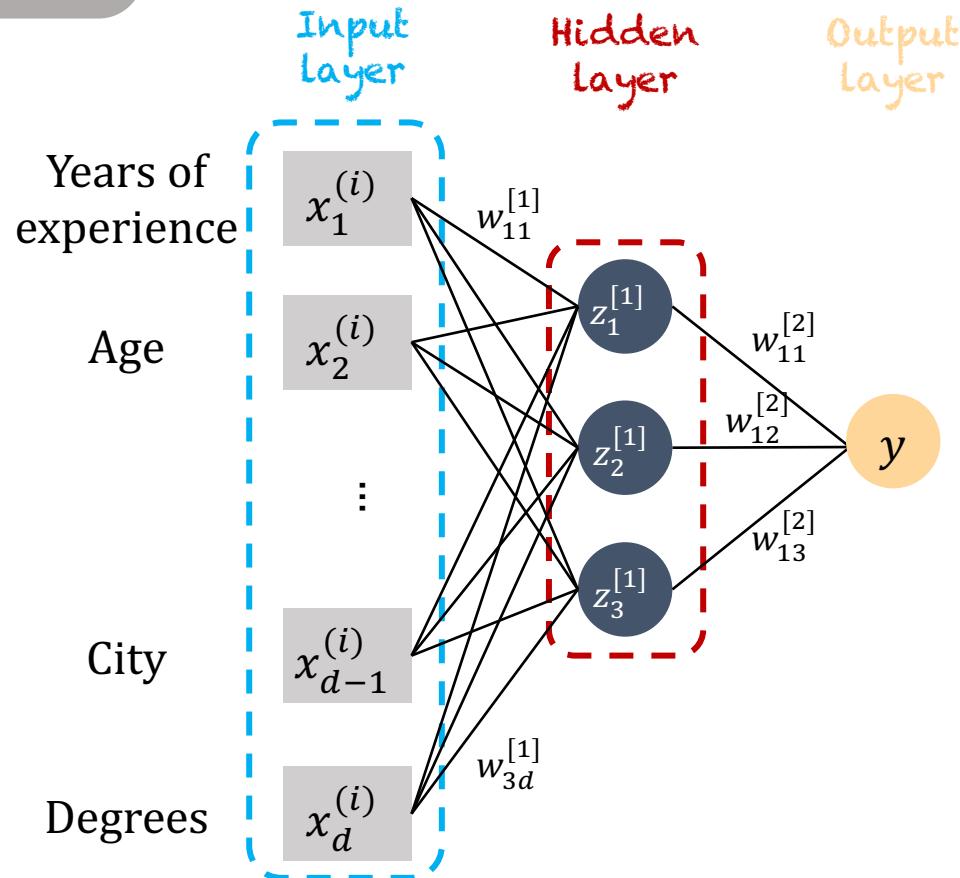
Our first model

Supervised learning theory

Other basic models

Deep learning models

## Artificial neural networks



- Units in a same layer do not interact
- Data go from input to output in a **feed-forward way**
- The width of the output layer corresponds to the number of classes/values that you want to predict
- The **activation** of the unit  $j$  in layer one is given by

Layer 1

Weights of Layer 1 for unit  $j$  (vectorized)

Bias of the unit  $j$  in Layer 1

$$z_j^{[1]} = s(w_j^{[1]} \cdot x^{(i)} + b_j^{[1]})$$

- We also define the **pre-activation**

$$u_j^{[1]} = w_j^{[1]} \cdot x + b_j^{[1]}$$

# More layers: the art of learning features

54

Introduction to ML

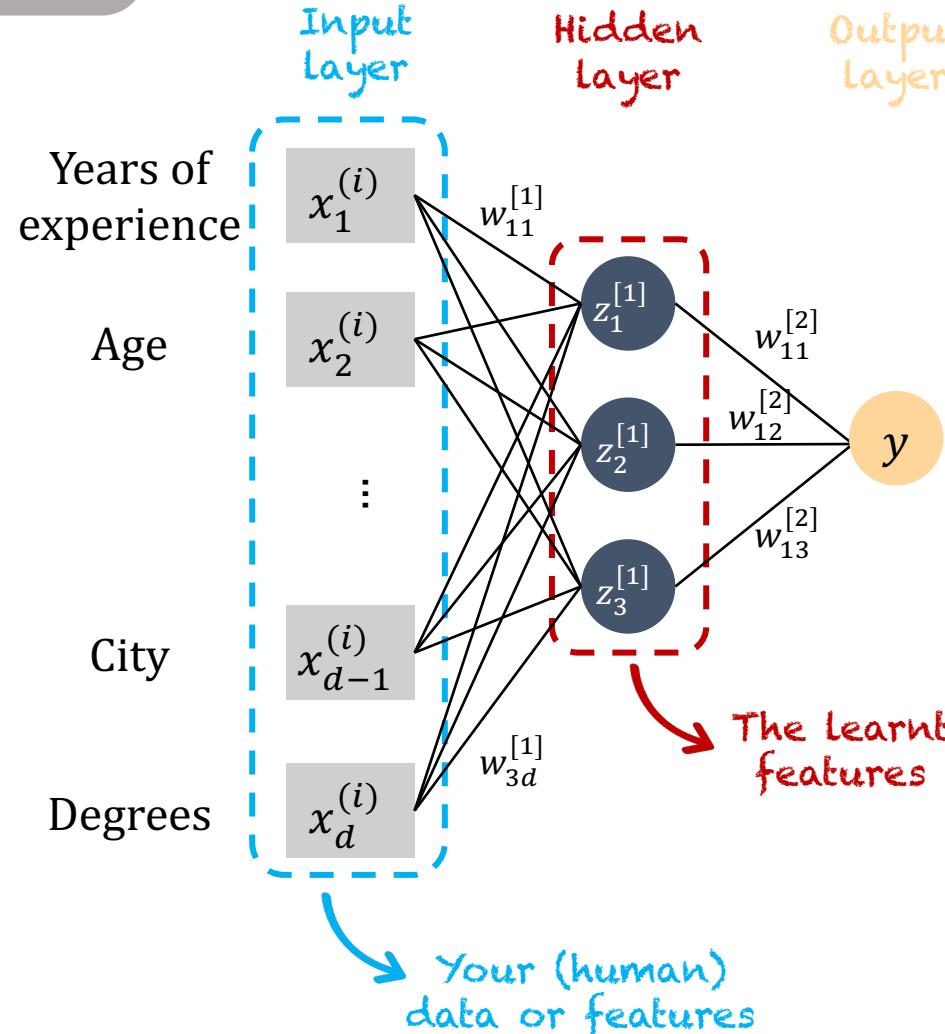
Our first model

Supervised learning theory

Other basic models

Deep learning models

## Artificial neural networks



Let's have a look at the output of this network

$$y = s \left( \sum_{j=1}^d w_{1j}^{[2]} z_j + b^{[2]} \right)$$

$\{z_j\}_{j=1,\dots,3}$  act as **new features** to predict  $y$

## Artificial neural networks

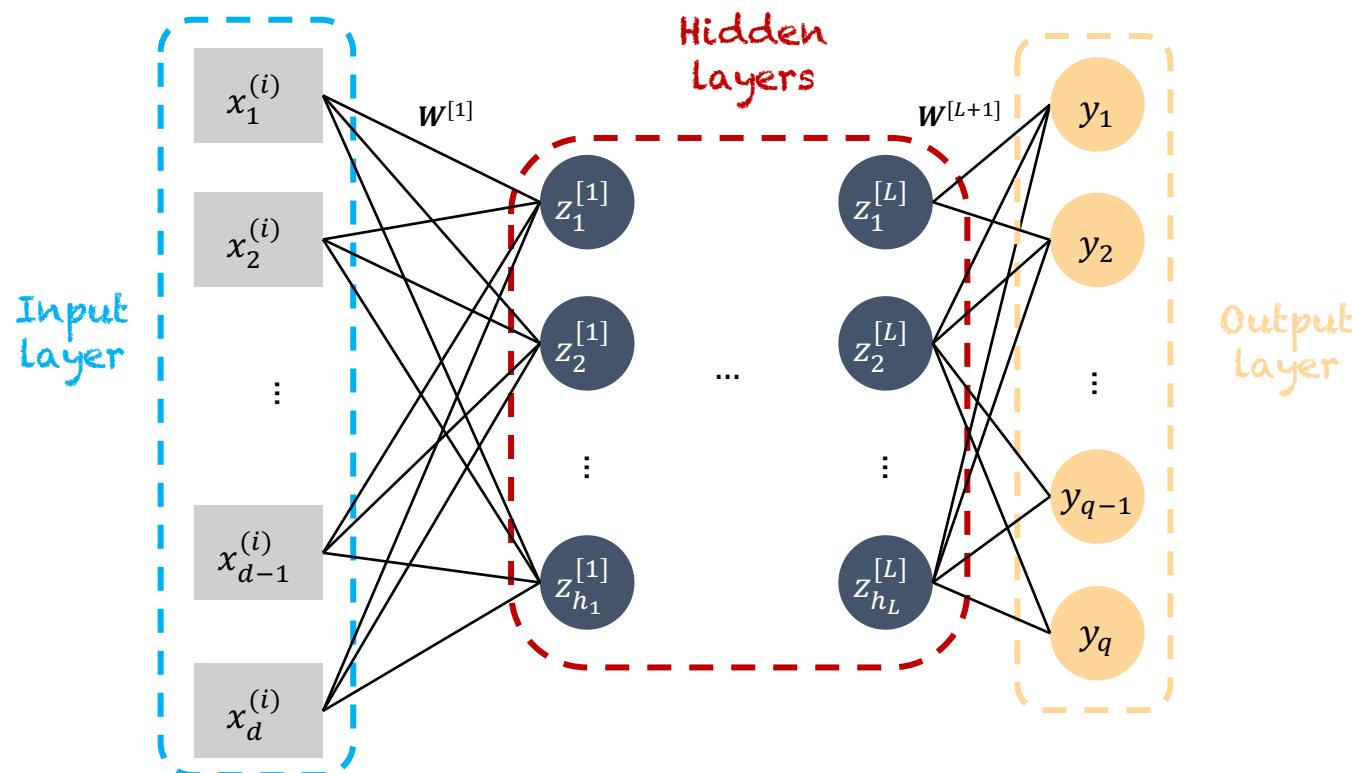
- A fully-connected neural network with  $d$  inputs,  $L$  hidden layers of width  $h_1, h_2, \dots, h_L$  and an output layer of size  $q$
- The  $j^{\text{th}}$  output is computed as

$$y_j = s^{[L+1]}(\mathbf{W}^{[L+1]}\mathbf{z}^{[L]}) \quad \text{with } \mathbf{z}^{[L]} = [1, z_1^{[L]}, z_2^{[L]}, \dots, z_{h_L}^{[L]}]^T$$

$$y_j = s^{[L+1]} \left( \mathbf{W}^{[L+1]} s^{[L]} \left( \mathbf{W}^{[L]} \dots s^{[1]}(\mathbf{W}^{[1]} \mathbf{x}) \right) \right)$$

$$\mathbf{w}_j^{[l]} = [b_j^{[l]}, w_{1j}^{[l]}, w_{2j}^{[l]}, \dots, w_{h_L j}^{[l]}]$$

$$\mathbf{W}^{[l]} = [\mathbf{w}_1^{[l]}, \mathbf{w}_2^{[l]}, \dots, \mathbf{w}_{h_l}^{[l]}]^T \in \mathbb{R}^{h_l \times (h_{l-1} + 1)}$$



- In the end, a **neural network** is a function  $f_{\theta}: X \rightarrow Y$  of some parameters ( $\theta = \text{weights and biases}$ )
- $f_{\theta}$  is a composition of non-linear function when  $s$  is non-linear; allowing to build non-linear estimators
- The cascade of layers learn successive features to predict the outputs

# General view of deep neural networks

56

Introduction to ML

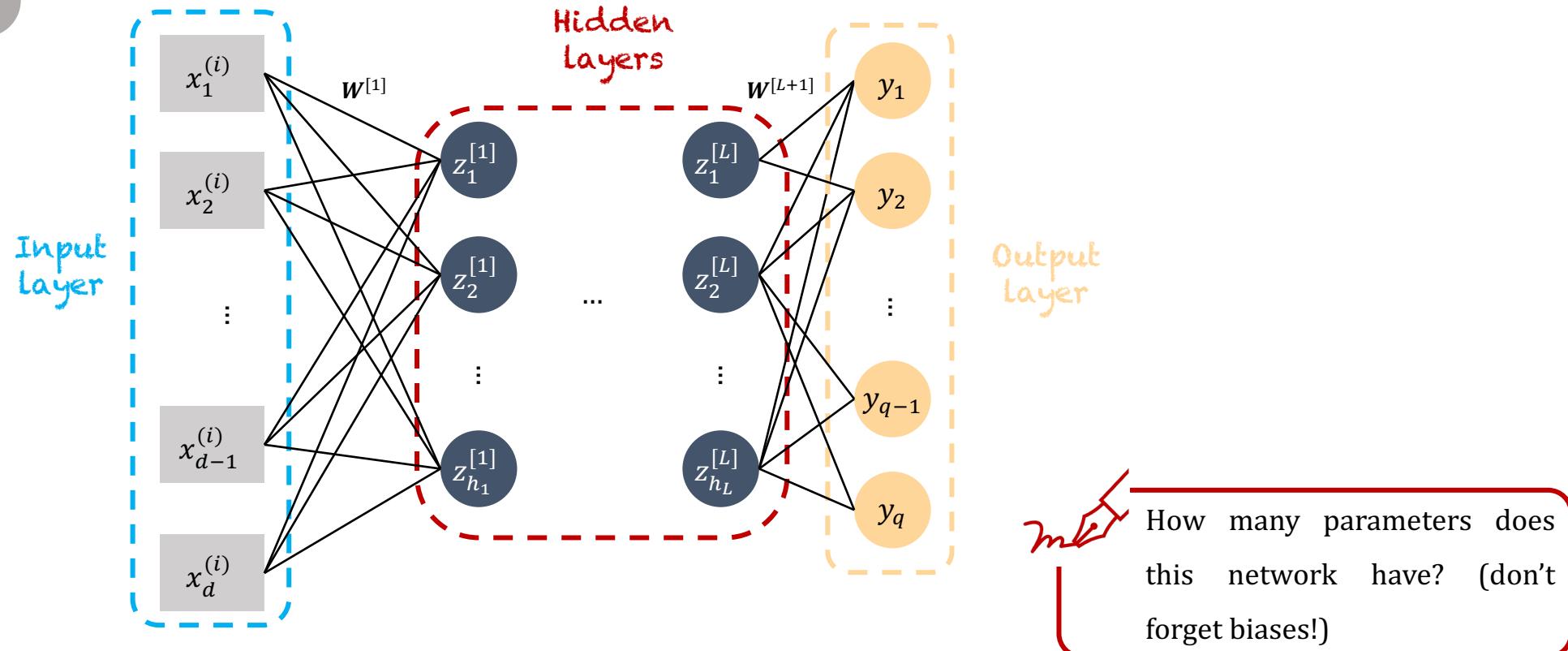
Our first model

Supervised learning theory

Other basic models

Deep learning models

## Artificial neural networks



+ Both regression and classification, learns features, good performances when enough data

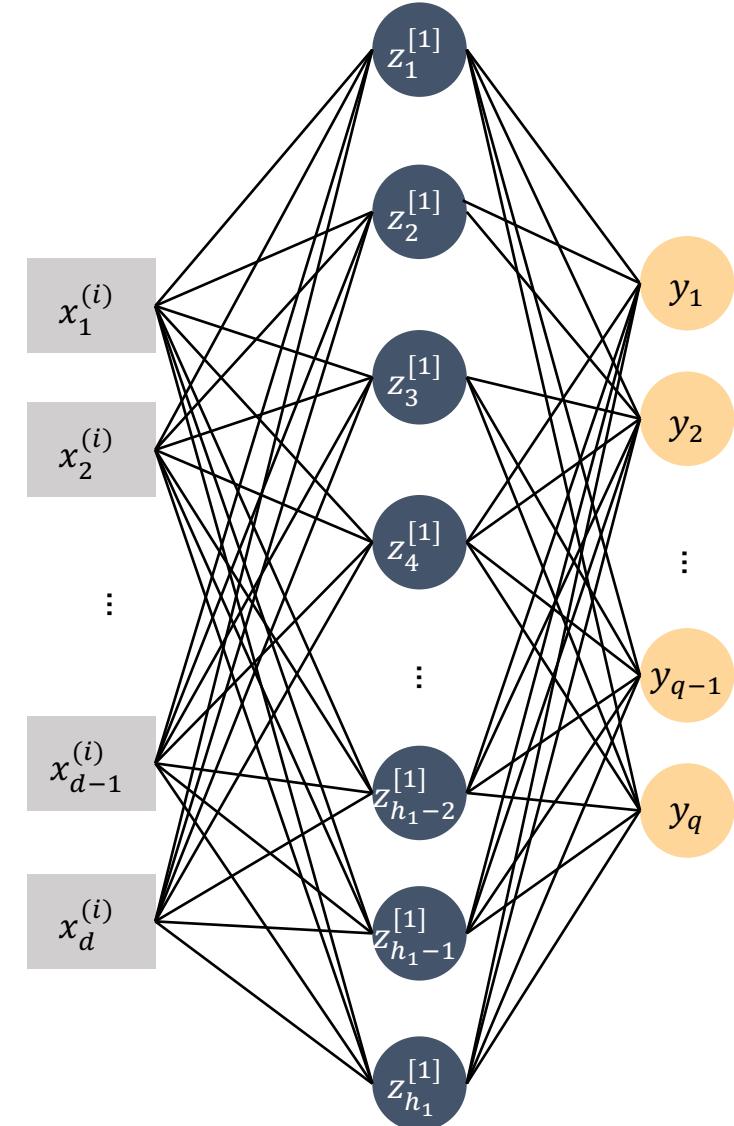
- Need a lot of data to train, subject to overfitting, uninterpretable, targeted-purpose architectures usually work better

## Artificial neural networks

- Some of the reasons why neural networks became later so popular "recently":
  1. Data accessibility
  2. Efficient algorithms and hardware evolution
  3. Automatic learning of features
- But they are also universal approximators (see [Cybenko 1989](#))

### Theorem (informal)

*A fully-connected neural network with a single hidden layer ( $L = 1$ ) with enough neurons ( $h_1$  large) can fit **any arbitrary smooth function.***



# ANN: example on the XNOR dataset

58

Introduction to ML

Our first model

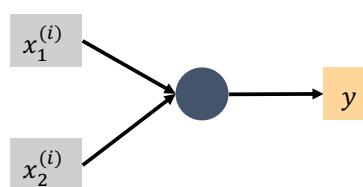
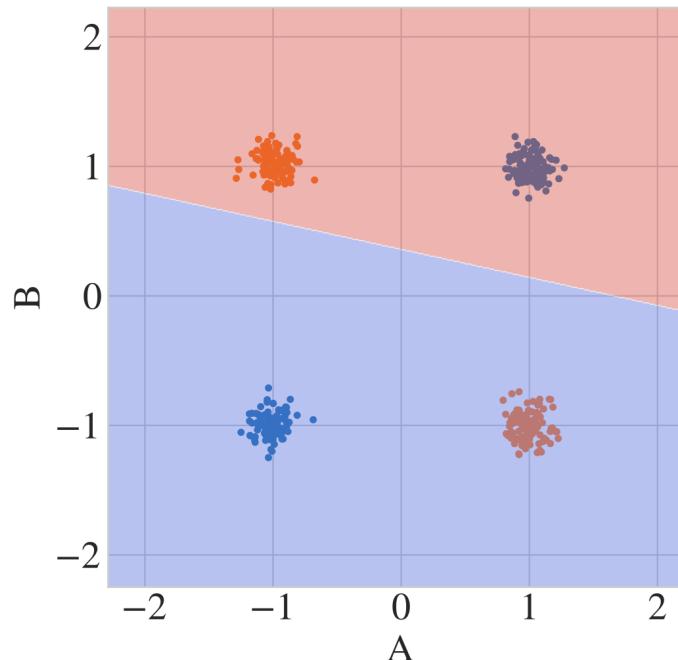
Supervised learning theory

Other basic models

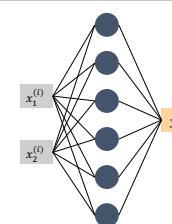
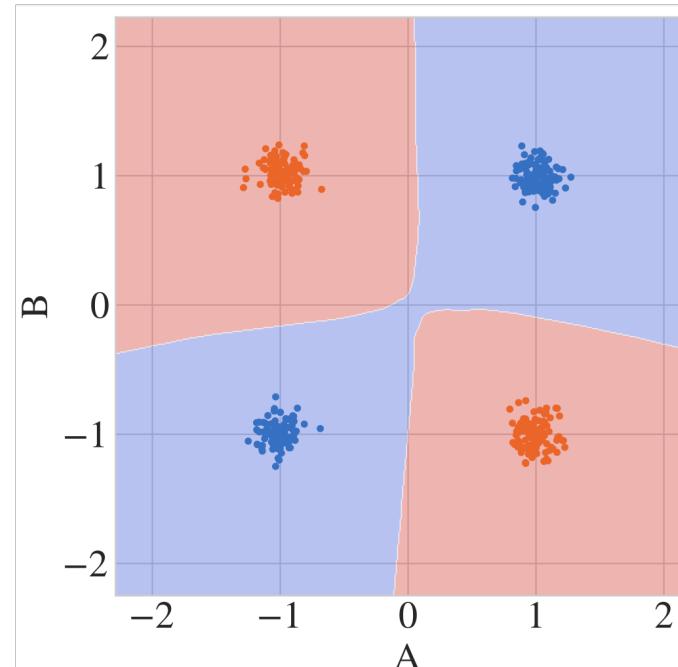
Deep learning models

## Artificial neural networks

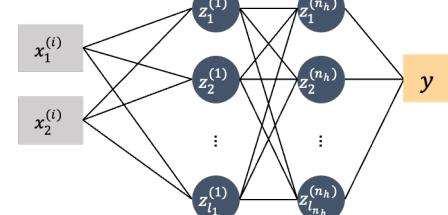
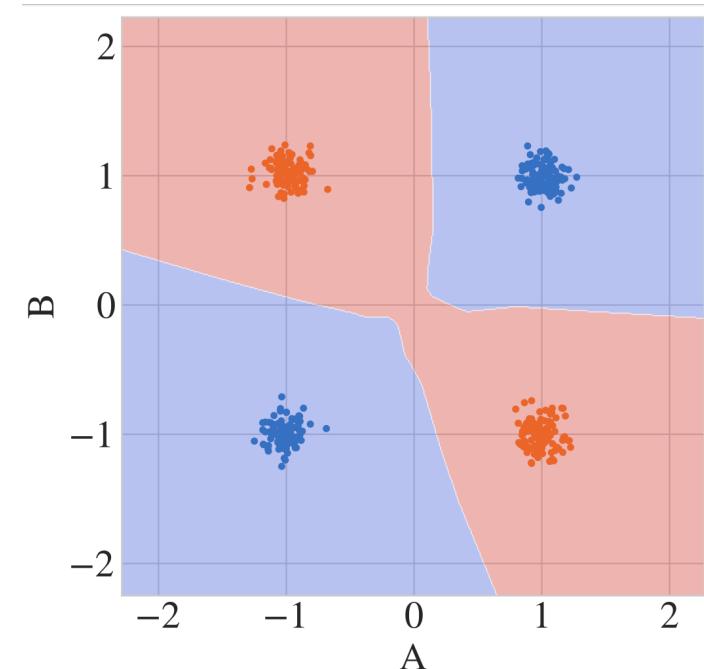
1 hidden layer of 1 unit



1 hidden layer of 100 units



2 hidden layers of 10 units



# Computing gradients in NNs: backpropagation

59

Introduction to ML

Our first model

Supervised learning theory

Other basic models

Deep learning models

- The problem is that neural networks are compositions of non-linear functions

$$y_j = s^{[L+1]} \left( \mathbf{W}^{[L+1]} s^{[L]} \left( \mathbf{W}^{[L]} \dots s^{[1]}(\mathbf{W}^{[1]} \mathbf{x}) \right) \right)$$

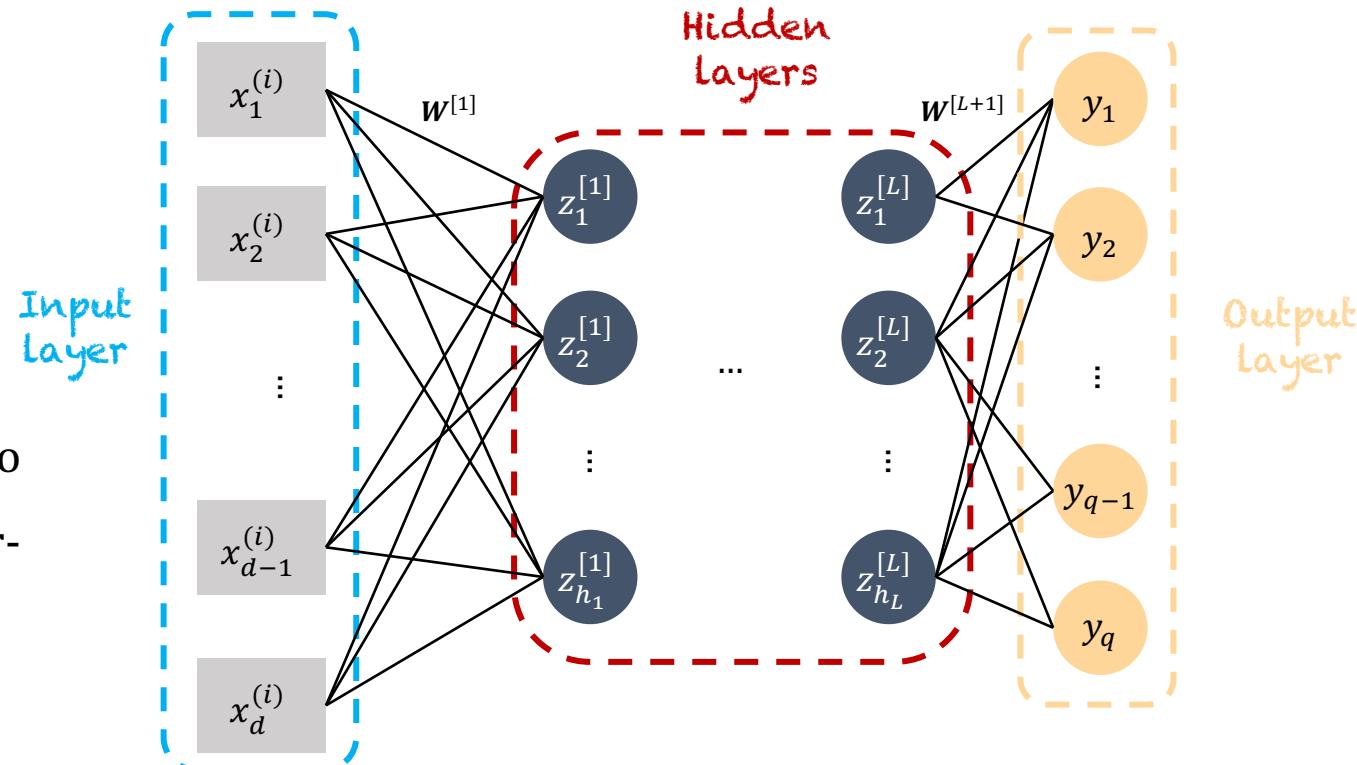
$$\mathbf{W}^{[l]} = [\mathbf{w}_1^{[l]}, \mathbf{w}_2^{[l]}, \dots, \mathbf{w}_{h_l}^{[l]}]^T \in \mathbb{R}^{h_l \times (h_{l-1}+1)}$$

$$\mathbf{w}_j^{[l]} = [b_j^{[l]}, w_{1j}^{[l]}, w_{2j}^{[l]}, \dots, w_{h_l j}^{[l]}]$$

- We however need to optimize the cost function to obtain the “best” values of  $\mathbf{W}$  producing the closer-to-optimal target values

- How to compute  $\frac{\partial \ell(\mathbf{W}^{[0]}, \dots, \mathbf{W}^{[L+1]})}{\partial w_{ij}^{[l]}}$ ?

- The backpropagation of errors:** an application of the **chain rule!**



# Computing gradients in NNs: backpropagation

60

Introduction to ML

Our first model

Supervised learning theory

Other basic models

Deep learning models

- Take the example of a fully-connected network with  $d = 3$ , one hidden layer, and two output neurons:

$$\hat{\mathbf{y}} = s(\mathbf{W}^{[2]}s(\mathbf{W}^{[1]}\mathbf{x})) \quad \text{with } \mathbf{x} = [1, x_1, x_2, \dots, x_d]^T$$

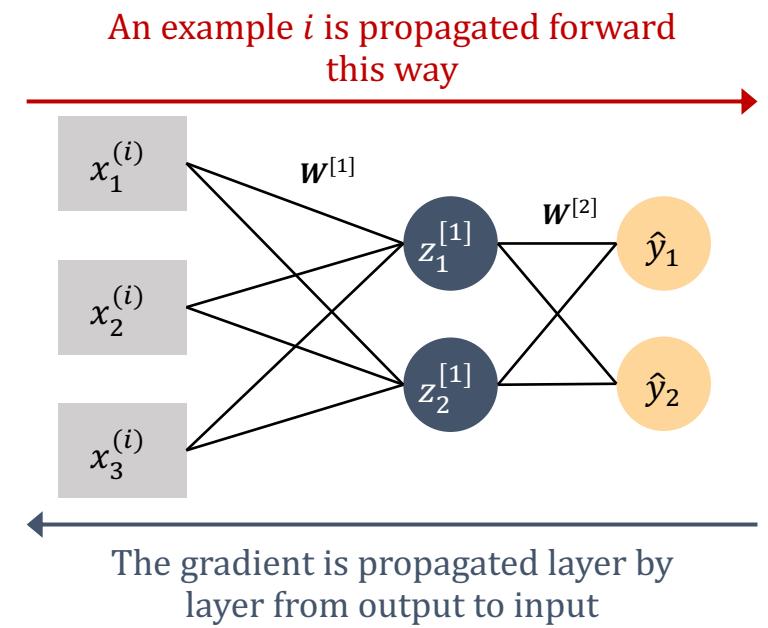
$$\mathbf{W}^{[1]} = \begin{bmatrix} b_1^{[1]} & w_{11}^{[1]} & w_{12}^{[1]} & w_{13}^{[1]} \\ b_2^{[1]} & w_{21}^{[1]} & w_{22}^{[1]} & w_{23}^{[1]} \end{bmatrix} \quad \mathbf{W}^{[2]} = \begin{bmatrix} b_1^{[2]} & w_{11}^{[2]} & w_{12}^{[2]} \\ b_2^{[2]} & w_{21}^{[2]} & w_{22}^{[2]} \end{bmatrix}$$

- Consider the squared error loss function for each training example  $i$

$$\ell(\mathbf{W}^{[1]}, \mathbf{W}^{[2]}) = \frac{1}{2} \sum_{j=1}^2 (y_j - \hat{y}_j)^2$$

- From the equation of  $\hat{\mathbf{y}}$ , we see the contribution of  $\mathbf{W}^{[2]}$  is “closer” to the output than  $\mathbf{W}^{[1]}$

- Let's compute  $\frac{\partial \ell(\mathbf{W}^{[1]}, \mathbf{W}^{[2]})}{\partial w_{11}^{[2]}}$



# Computing gradients in NNs: backpropagation

61

Introduction to ML

Our first model

Supervised learning theory

Other basic models

Deep learning models

- Using the chain rule

$$\frac{\partial \ell(\mathbf{W}^{[1]}, \mathbf{W}^{[2]})}{\partial w_{11}^{[2]}} = \frac{\partial \ell}{\partial u_1^{[2]}} \times \frac{\partial u_1^{[2]}}{\partial w_{11}^{[2]}}$$

where  $u_1^{[2]}$  is the pre-activation of the unit 1 in layer 2 (here output layer)

$$u_1^{[2]} = w_{11}^{[2]} z_1^{[1]} + w_{21}^{[2]} z_2^{[1]} + b_1^{[2]}$$

- The second term is then easy to compute as  $\frac{\partial u_1^{[2]}}{\partial w_{11}^{[2]}} = z_1^{[1]}$
- For the first term, use the chain rule again leads to

$$\delta_1^{[2]} = \frac{\partial \ell}{\partial u_1^{[2]}} = \frac{\partial \ell}{\partial \hat{y}_1} \frac{\partial \hat{y}_1}{\partial u_1^{[2]}}$$

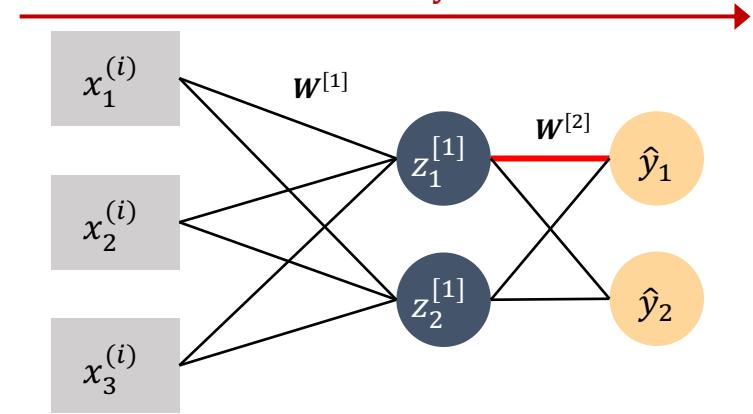
with  $\hat{y}_1 = s(u_1^{[2]}) = s(w_{11}^{[2]} z_1^{[1]} + w_{21}^{[2]} z_2^{[1]} + b_1^{[2]})$

Finally,

$$\frac{\partial \ell(\mathbf{W}^{[1]}, \mathbf{W}^{[2]})}{\partial w_{11}^{[2]}} = -(y_1 - \hat{y}_1) s'(u_1^{[2]}) z_1^{[1]}$$

$$\ell(\mathbf{W}^{[1]}, \mathbf{W}^{[2]}) = \frac{1}{2} \sum_{j=1}^2 (y_j - \hat{y}_j)^2$$

An example  $i$  is propagated forward this way



The gradient is propagated layer by layer from output to input

$$\hat{y} = s(\mathbf{W}^{[2]} s(\mathbf{W}^{[1]} \mathbf{x}))$$

$$\mathbf{W}^{[1]} = \begin{bmatrix} b_1^{[1]} & w_{11}^{[1]} & w_{12}^{[1]} & w_{13}^{[1]} \\ b_2^{[1]} & w_{21}^{[1]} & w_{22}^{[1]} & w_{23}^{[1]} \end{bmatrix}$$

$$\mathbf{W}^{[2]} = \begin{bmatrix} b_1^{[2]} & w_{11}^{[2]} & w_{12}^{[2]} \\ b_2^{[2]} & w_{21}^{[2]} & w_{22}^{[2]} \end{bmatrix}$$

# Computing gradients in NNs: backpropagation

62

Introduction to ML

Our first model

Supervised learning theory

Other basic models

Deep learning models

- You can proceed the same for all the weights linked to the output layer
- What about parameters in the hidden layers?
- Let's compute

$$\frac{\partial \ell(\mathbf{W}^{[1]}, \mathbf{W}^{[2]})}{\partial w_{11}^{[1]}} = \frac{\partial \ell}{\partial u_1^{[1]}} \times \frac{\partial u_1^{[1]}}{\partial w_{11}^{[1]}}$$

- The second term is still easy to compute
- For the first term, we have  $u_1^{[1]} = w_{11}^{[1]}x_1 + w_{21}^{[1]}x_2 + b_1^{[1]}$  and observe that

there are two paths to reach the weight  $w_{11}^{[1]}$ :

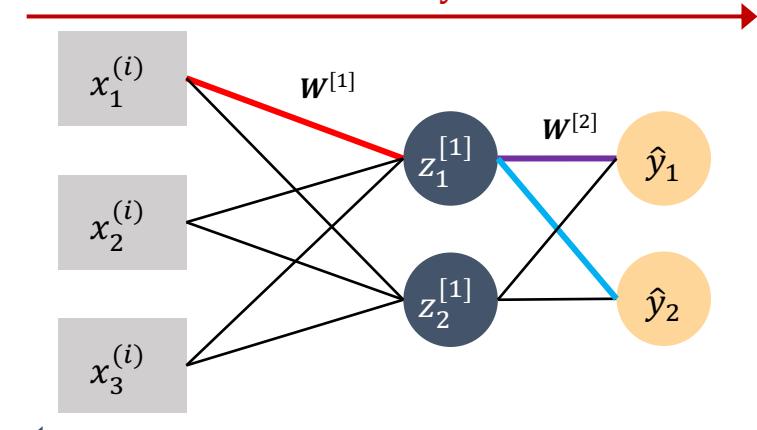
$$\delta_1^{[1]} = \frac{\partial \ell}{\partial u_1^{[2]}} \frac{\partial u_1^{[2]}}{\partial u_1^{[1]}} + \frac{\partial \ell}{\partial u_2^{[2]}} \frac{\partial u_2^{[2]}}{\partial u_1^{[1]}}$$



Already computed in the previous layer (hence the name "Backpropagation")

$$\ell(\mathbf{W}^{[1]}, \mathbf{W}^{[2]}) = \frac{1}{2} \sum_{j=1}^2 (y_j - \hat{y}_j)^2$$

An example  $i$  is propagated forward this way



The gradient is propagated layer by layer from output to input

$$\hat{\mathbf{y}} = s(\mathbf{W}^{[2]}s(\mathbf{W}^{[1]}\mathbf{x}))$$

$$\mathbf{W}^{[1]} = \begin{bmatrix} b_1^{[1]} & w_{11}^{[1]} & w_{12}^{[1]} & w_{13}^{[1]} \\ b_2^{[1]} & w_{21}^{[1]} & w_{22}^{[1]} & w_{23}^{[1]} \end{bmatrix}$$

$$\mathbf{W}^{[2]} = \begin{bmatrix} b_1^{[2]} & w_{11}^{[2]} & w_{12}^{[2]} \\ b_2^{[2]} & w_{21}^{[2]} & w_{22}^{[2]} \end{bmatrix}$$

# Computing gradients in NNs: backpropagation

63

Introduction to ML

Our first model

Supervised learning theory

Other basic models

Deep learning models

$$\delta_1^{[1]} = \frac{\partial \ell}{\partial u_1^{[2]}} \frac{\partial u_1^{[2]}}{\partial u_1^{[1]}} + \frac{\partial \ell}{\partial u_2^{[2]}} \frac{\partial u_2^{[2]}}{\partial u_1^{[1]}}$$



Already computed in the previous layer (hence the name "Backpropagation")

- To compute  $\frac{\partial u_1^{[2]}}{\partial u_1^{[1]}}$  and  $\frac{\partial u_2^{[2]}}{\partial u_1^{[1]}}$ , we can use the same chain rule again giving

$$\frac{\partial u_1^{[2]}}{\partial u_1^{[1]}} = \frac{\partial u_1^{[2]}}{\partial z_1^{[1]}} \frac{\partial z_1^{[1]}}{\partial u_1^{[1]}} = w_{11}^{[2]} s'(u_1^{[1]})$$

$$\frac{\partial u_2^{[2]}}{\partial u_1^{[1]}} = \frac{\partial u_2^{[2]}}{\partial z_1^{[1]}} \frac{\partial z_1^{[1]}}{\partial u_1^{[1]}} = w_{21}^{[2]} s'(u_1^{[1]})$$

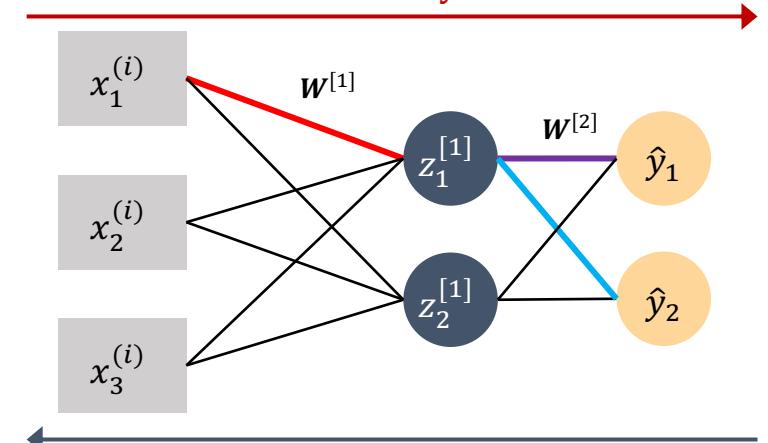
- Hence,

$$\frac{\partial \ell(\mathbf{W}^{[1]}, \mathbf{W}^{[2]})}{\partial w_{11}^{[1]}} = (\delta_1^{[2]} w_{11}^{[2]} + \delta_2^{[2]} w_{21}^{[2]}) s'(u_1^{[1]}) \mathbf{x}_1$$

- Where everything is known!

$$\ell(\mathbf{W}^{[1]}, \mathbf{W}^{[2]}) = \frac{1}{2} \sum_{j=1}^2 (y_j - \hat{y}_j)^2$$

An example  $i$  is propagated forward this way



The gradient is propagated layer by layer from output to input

$$\hat{\mathbf{y}} = s(\mathbf{W}^{[2]} s(\mathbf{W}^{[1]} \mathbf{x}))$$

$$\mathbf{W}^{[1]} = \begin{bmatrix} b_1^{[1]} & w_{11}^{[1]} & w_{12}^{[1]} & w_{13}^{[1]} \\ b_2^{[1]} & w_{21}^{[1]} & w_{22}^{[1]} & w_{23}^{[1]} \end{bmatrix}$$

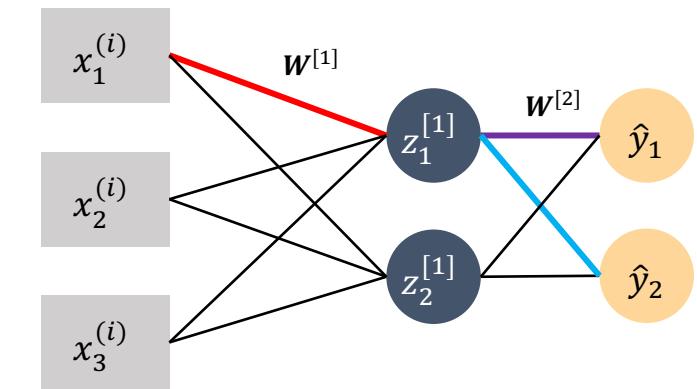
$$\mathbf{W}^{[2]} = \begin{bmatrix} b_1^{[2]} & w_{11}^{[2]} & w_{12}^{[2]} \\ b_2^{[2]} & w_{21}^{[2]} & w_{22}^{[2]} \end{bmatrix}$$

$$\ell(\mathbf{W}^{[1]}, \mathbf{W}^{[2]}) = \frac{1}{2} \sum_{j=1}^2 (y_j - \hat{y}_j)^2$$



- Writing those equations **propagating the errors from output to input recursively** for both weights and biases leads to the **backpropagation algorithm**
- This method grants an **efficient computation of the gradient** in neural networks
- Together with the SGD algorithm they allow to train efficiently networks with many parameters

An example  $i$  is propagated forward this way



The gradient is propagated layer by layer from output to input

$$\hat{\mathbf{y}} = s(\mathbf{W}^{[2]} s(\mathbf{W}^{[1]} \mathbf{x}))$$

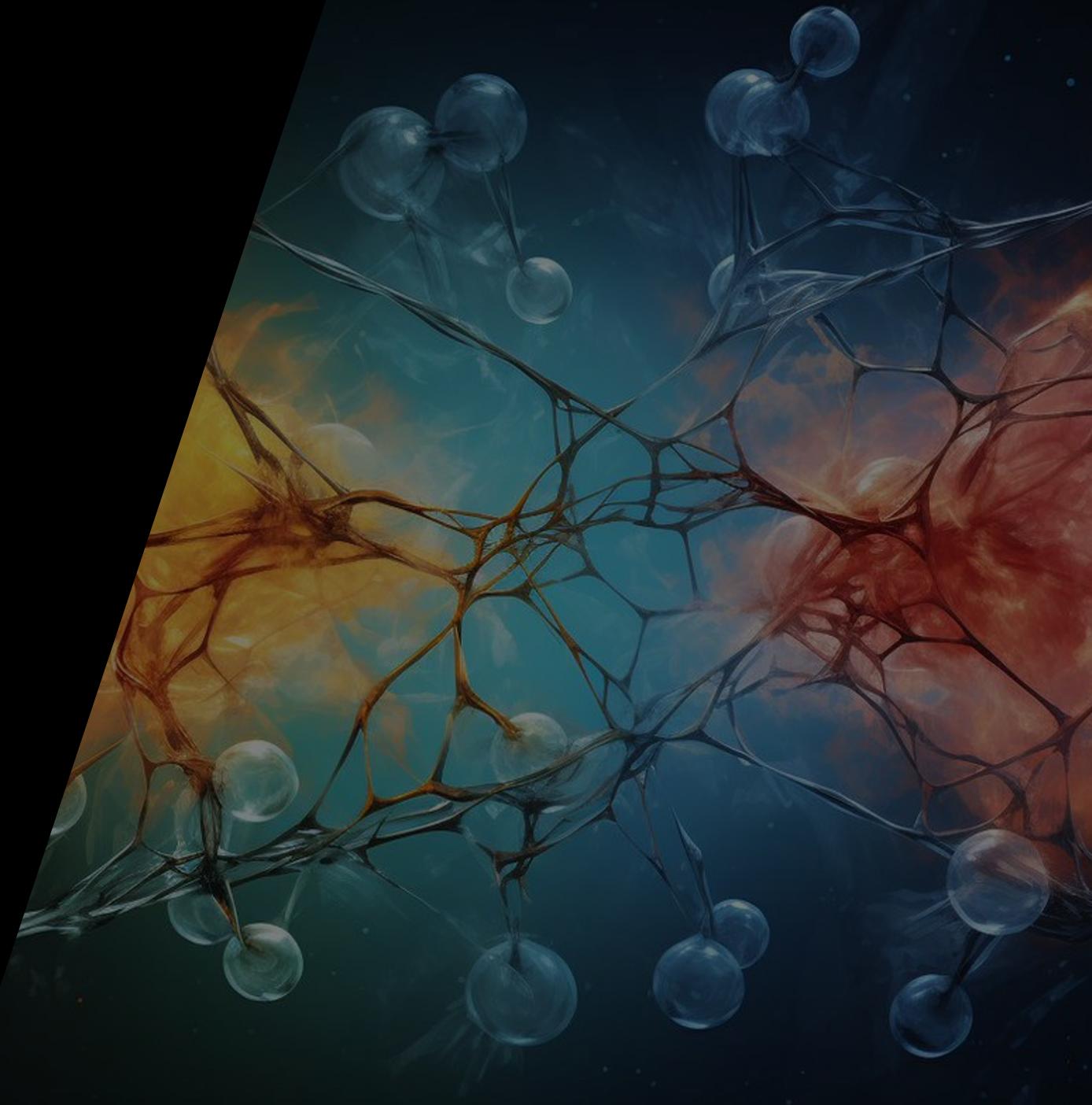
$$\mathbf{W}^{[1]} = \begin{bmatrix} b_1^{[1]} & w_{11}^{[1]} & w_{12}^{[1]} & w_{13}^{[1]} \\ b_2^{[1]} & w_{21}^{[1]} & w_{22}^{[1]} & w_{23}^{[1]} \end{bmatrix}$$

$$\mathbf{W}^{[2]} = \begin{bmatrix} b_1^{[2]} & w_{11}^{[2]} & w_{12}^{[2]} \\ b_2^{[2]} & w_{21}^{[2]} & w_{22}^{[2]} \end{bmatrix}$$

# Encoding translation symmetry: Convolutional neural networks

*Contents :*

- *Computer vision motivations*
- *Convolution operation*
- *Convolutional layers*
- *Pooling layers*
- *Convolutional neural nets*



# Computer vision tasks and challenges

66

Introduction to ML

Our first model

Supervised learning theory

Other basic models

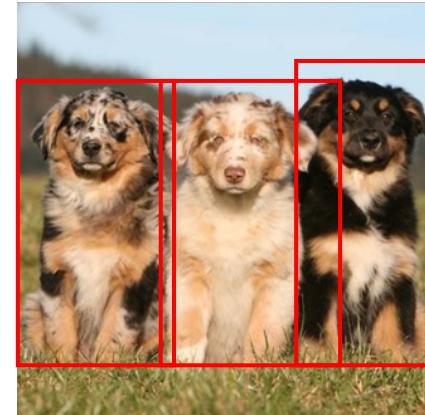
Deep learning models

## Image classification

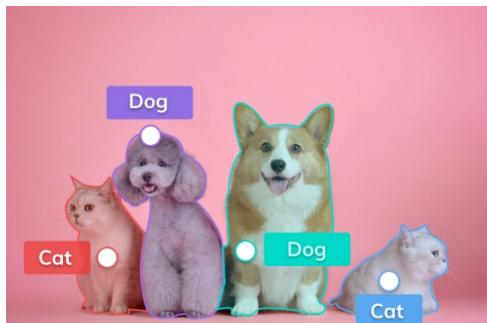


Dog?

## Object detection



## Instance segmentation



## Caption generation



"A dog and a little boy playing with a basketball in the grass"

## Image classification

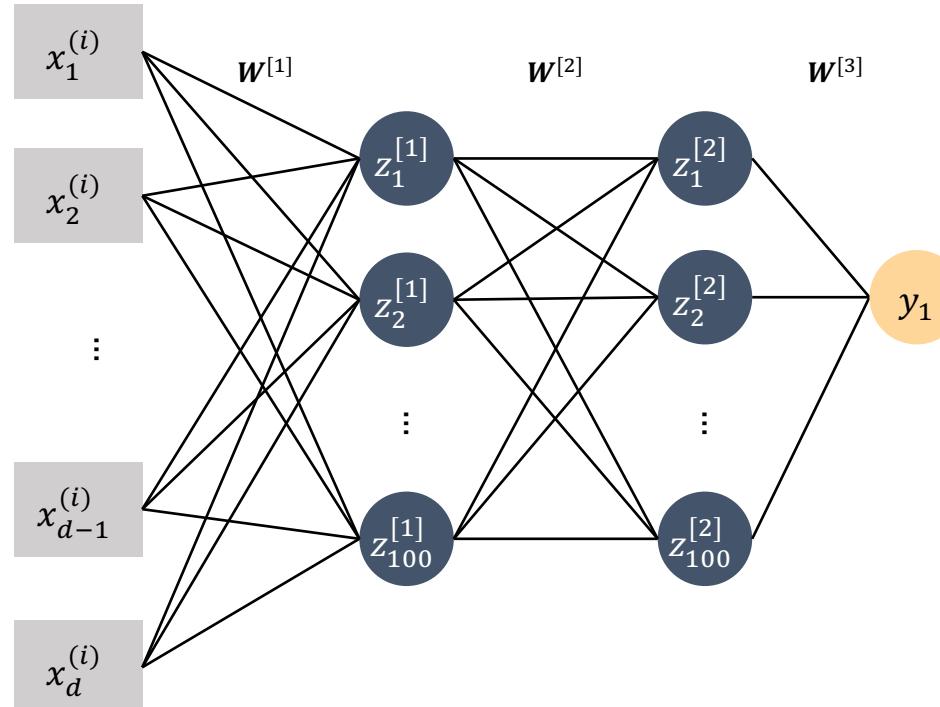


Dog?

12M pixels  
 $x^{(i)} = [x_1^{(i)}, x_2^{(i)}, \dots, x_{12192768}^{(i)}]$   
 $d = 12192768$

## Why going beyond fully-connected neural network?

- Consider the simplest task with **classification**
- Standard images taken by a current smartphone are of size  $4032 \times 3024 = 12\text{M}$  pixels
- A 2 layer-FCNN with just 100 (!) units has  $> 1\text{B}$  parameters



## Image classification



A dog

Translation



Still a dog!



In principle, an FCNN *could* learn that but at the cost of:  
1) A huge database,  
2) A lot of redundant weights to encode the same patterns at different locations.



They are all zeros!

# The convolution operation

70

Introduction to ML

Our first model

Supervised learning theory

Other basic models

Deep learning models

## Convolution

1D

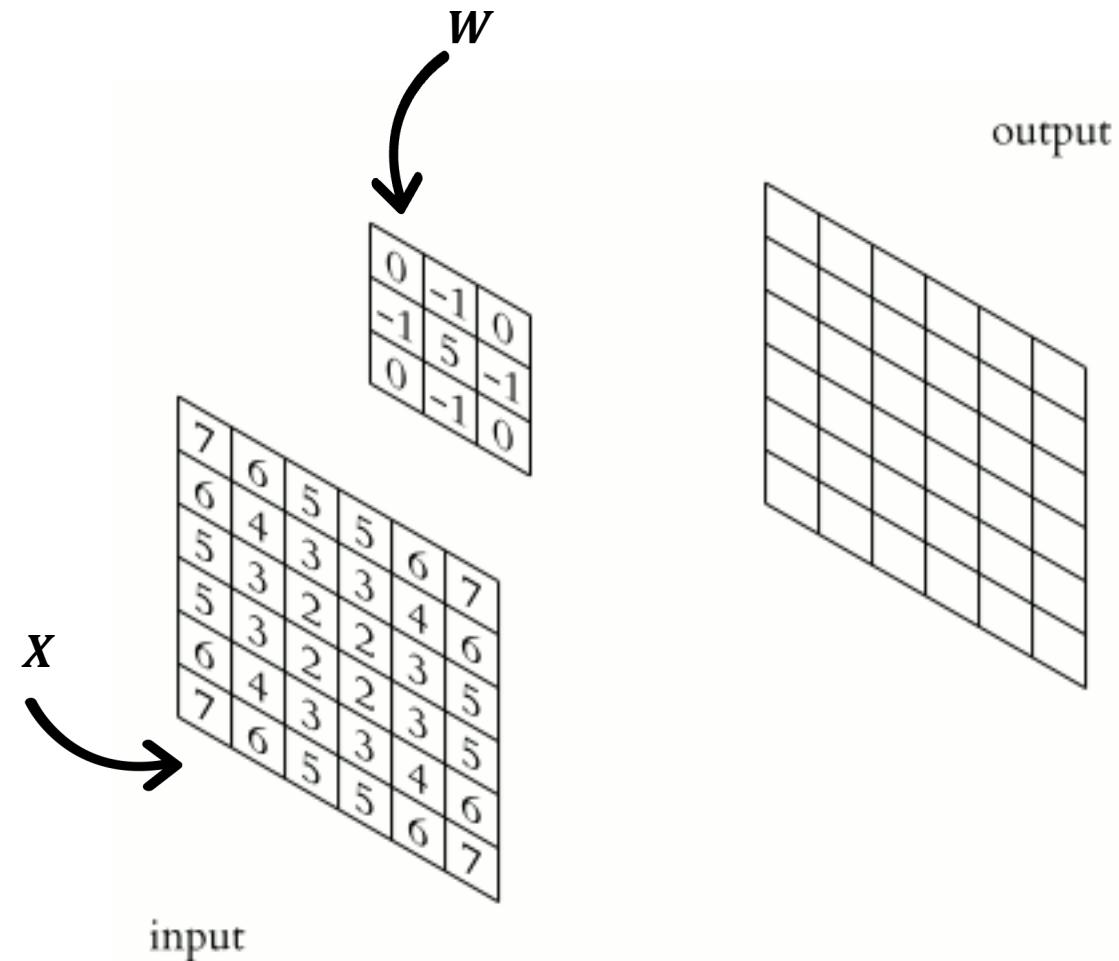
$$(x * w)(t) = \sum_{k=-K}^K x(k)w(t-k)$$

2D

$$(\mathbf{X} * \mathbf{W})_{ij} = \sum_{l=1}^L \sum_{k=1}^K x_{lk} w_{i-l, j-k}$$

Convolution operation is a **linear operation**  
**equivariant** to translation

$$\begin{aligned} \sum_{l=1}^L \sum_{k=1}^K x_{l+c, k+c} w_{i-l, j-k} &= \sum_{l=1}^L \sum_{k=1}^K x_{lk} w_{i-l+c, j-k+c} \\ &= (\mathbf{X} * \mathbf{W})_{i+c, j+c} \end{aligned}$$



Animation from Wikipédia

# Example: edge detection

71

Introduction to ML

Our first model

Supervised learning theory

Other basic models

Deep learning models

## Convolution

Example of a kernel for vertical edge detection

 $X$ 

1	1	1	0	0	0
1	1	1	0	0	0
1	1	1	0	0	0
1	1	1	0	0	0
1	1	1	0	0	0
1	1	1	0	0	0

 $W$ 

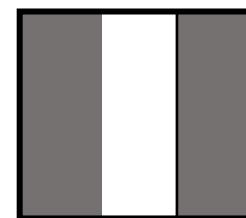
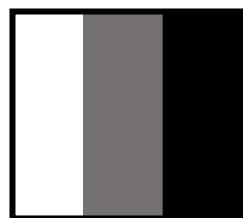
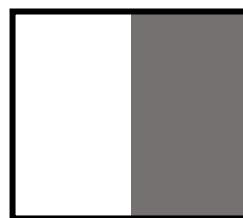
\*

1	0	-1
1	0	-1
1	0	-1

 $X * W$ 

=

0	3	3	0
0	3	3	0
0	3	3	0
0	3	3	0



The border is enhanced!

# Example: edge detection

72

Introduction to ML

Our first model

Supervised learning theory

Other basic models

Deep learning models

## Convolution

Example of a kernel for vertical edge detection

 $X$ 

0	0	0	0	0	0	0	0	0
0	1	1	1	1	0	0	0	0
0	1	1	1	1	0	0	0	0
0	1	1	1	1	0	0	0	0
0	1	1	1	1	0	0	0	0
0	1	1	1	1	0	0	0	0
0	1	1	1	1	0	0	0	0
0	1	1	1	1	0	0	0	0
0	0	0	0	0	0	0	0	0

 $W$ 

1	0	-1
1	0	-1
1	0	-1

\*

=

 $X * W$ 

-2	0	2	2	0	0
-2	0	3	3	0	0
-3	0	3	3	0	0
-3	0	3	3	0	0
-2	0	3	3	0	0
-2	0	2	2	0	0

$X$  and  $X * W$  have the same size



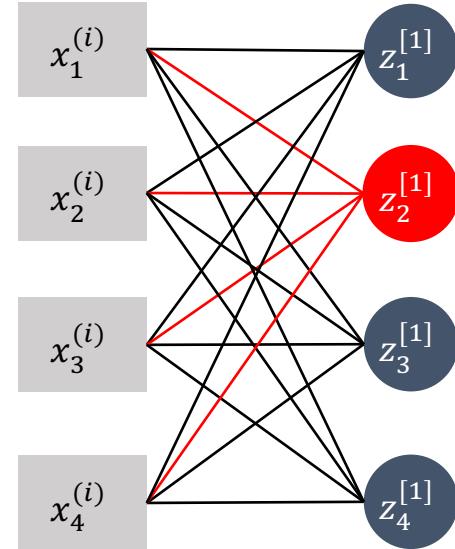
In practice, we use **padding** to add zeros around the image  $X$  so that the border of the image are seen as much as other pixels.

## Convolutional layers

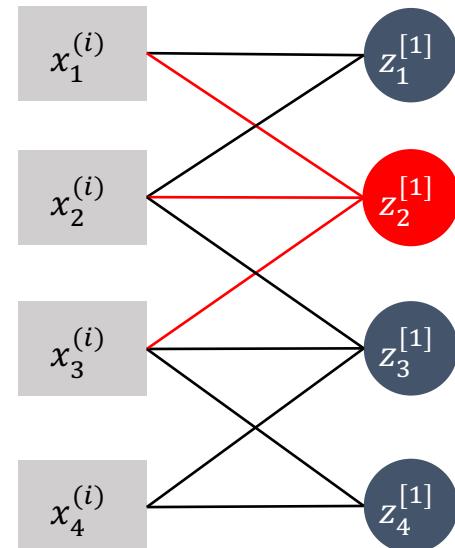
- A **convolutional layer** in a neural network is simply implementing the convolution operation with a filter (or kernel)  $W$  **shared across all the inputs**
- In 2 dimensions with a first layer of the same size as the input, a square weight matrix has  $K \times K$  elements, while a fully connected layer has  $d \times d$
- Typically,  $K \ll d$ . For instance, in MNIST  $d = 28 \times 28$ , while  $K \sim O(1)$ , usually 3 or 5
- Units of a given layer  $\ell$  only sees a subset of the activations of the previous layer  $\ell - 1$ : **local receptive field**
- **Deeper units** in the network **are influenced by more inputs**, hence learning higher-order features



The **convolutional layer** has three properties: **sparse interactions**, **parameters sharing** and **equivariant representation**.



A fully-connected layer



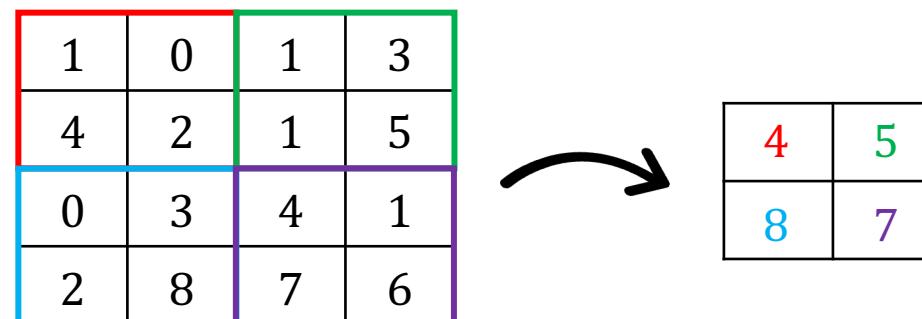
A convolutional layer

## Pooling layers

- To make the network **invariant to translation**, we need an additional element
- Idea: local invariance to translations of features in a given window
- In 2D

$$z_{ij}^{[\ell+1]} = \max_{(r,s) \in K_1 \times K_2} z_{rs}^{[\ell]}$$

$$z_{ij}^{[\ell+1]} = \frac{1}{K_1 \times K_2} \sum_r^K \sum_s^{K_2} z_{rs}^{[\ell]}$$



A max non-overlapping pooling layer with window of size 2x2



The **pooling layer** grants the network **some local invariance to translation** and **reduces the image size**

# Convolutional neural network

Introduction to ML

Our first model

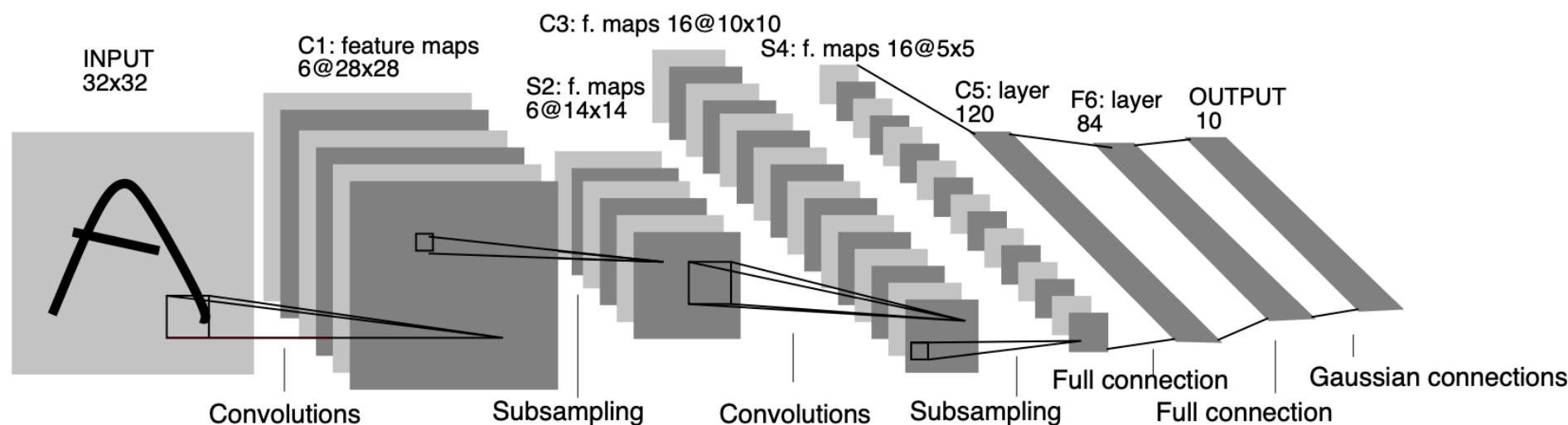
Supervised learning theory

Other basic models

Deep learning models

## CNNs

A convolutional neural network (CNN) is a cascade of **convolutional layers** and **pooling layers** to ensure **invariance to small (local) shifting, scaling and distortions** through local receptive field, shared weights and spatial sub-sampling.



Architecture of LeNet-5 used for digit recognition in [LeCun+98](#). It has 60,000 trainable parameters.

# Features learnt by a CNN

76

Introduction to ML

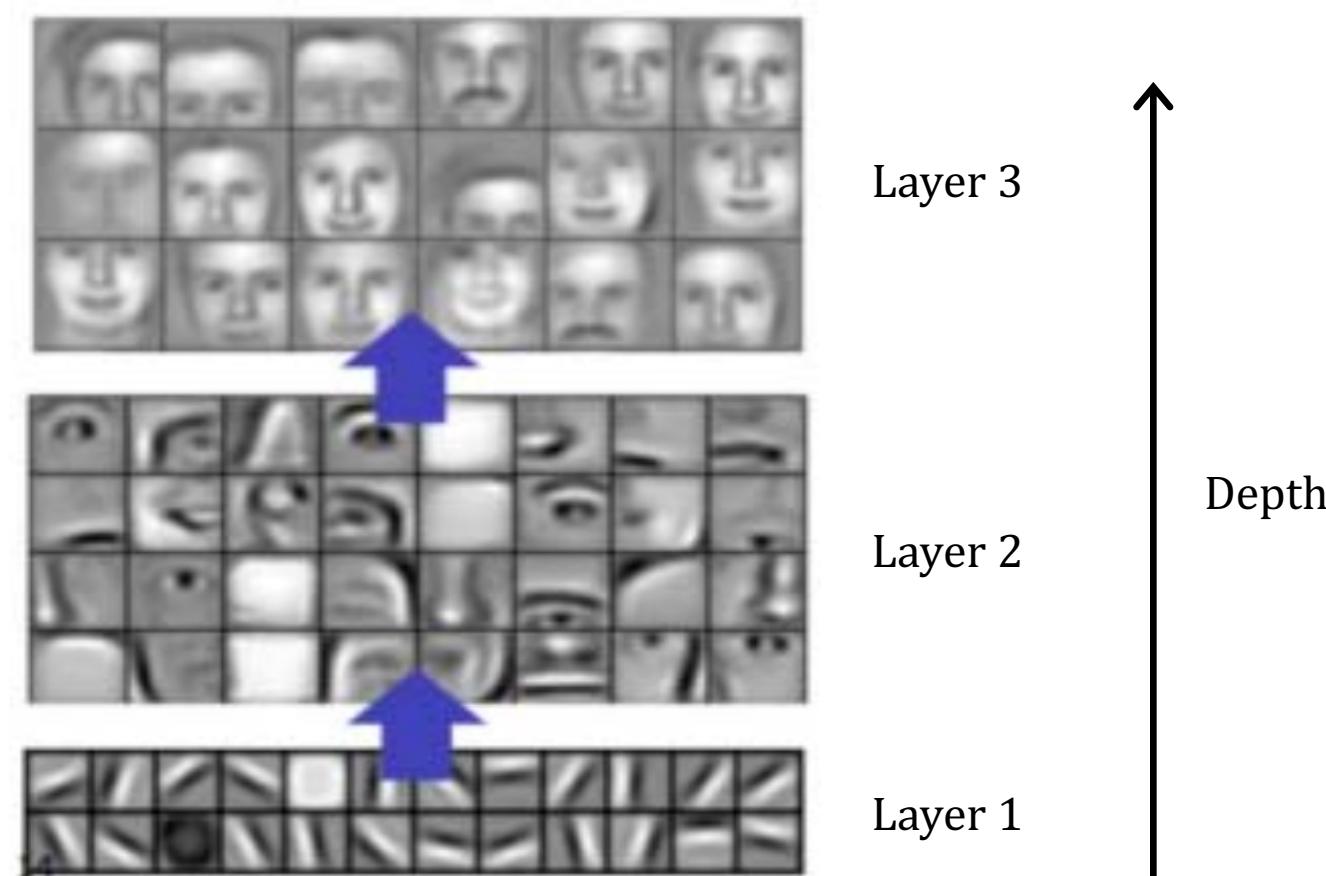
Our first model

Supervised learning theory

Other basic models

Deep learning models

More specific, large-scale,  
high-order, features



Local, first-order,  
features (edges)

Image from [Albawi et al., 2017](#)

## Organisation

- L'examen a lieu à la fin du semestre (date à définir)
- Il consiste en un challenge de science de données sur un sujet précis hébergé sur la plateforme de Challenge Data ENS <https://challengedata.ens.fr/challenges/122> (créez un compte participant !)

## Consignes

- Travail en groupe (3 maximum)
- Interactions hebdomadaires avec Pablo (ou moi selon le nombre de groupes)
- Deux soumissions par jours maximum (restriction du site)

## Travail attendu et notation

- Présentation orale de 15+5 minutes de la/les solution(s) retenue(s)
- Il ne s'agit pas que de battre le benchmark: il faut analyser les performances de l'algorithme, le comprendre, le présenter, le justifier
- Analyser les données : plusieurs chemins sont possibles
- Vous pouvez bien sûr vous inspirer de méthodes trouvées sur internet, dans des articles, etc.
- Expliquez les résultats numériques: avis sur les sources d'erreurs, signes de sur-apprentissage, les avantages et inconvénients de votre approche, etc.