

Machine learning principles with applications in physics

A crash course on
Machine Learning

Tony Bonnaire

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

Organisation of the course

Generalities

Introduction to ML

Basic supervised models

Optimisation

Generalization



Given by: Myself and Jorge Fernandez-de-Cossío Díaz.



Format: 3 lectures + projects.



Exam: Oral presentation of your “solution” to the projects.



Aim: Introduce you to the basics of ML through applications in physics.

Quick syllabus:

1. **(today 05/09, Tony) Crash course on ML: introduction, basic models, optimisation, generalisation**
2. (next week 12/09, Jorge) More on DL (VAEs, GANs, CNNs)
3. (in two weeks 19/09, Jorge + Tony) Introduction to Python for ML and DL (hands-on with PyTorch)

Some references:

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

Generalities

Generalities

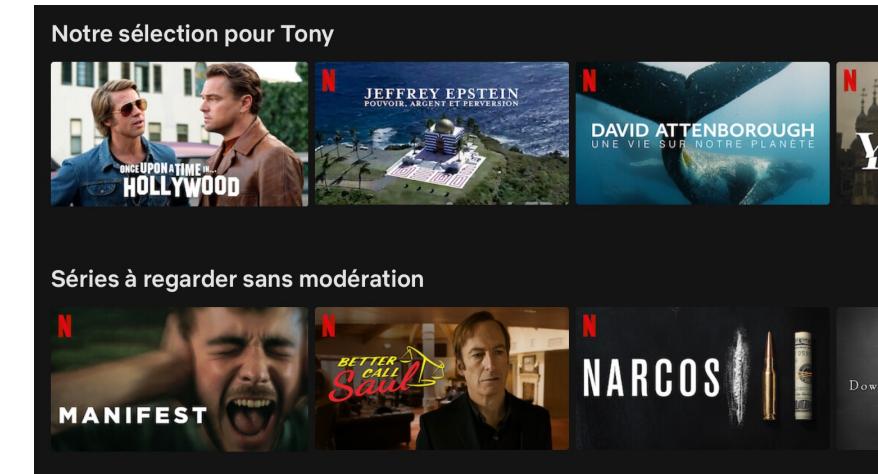
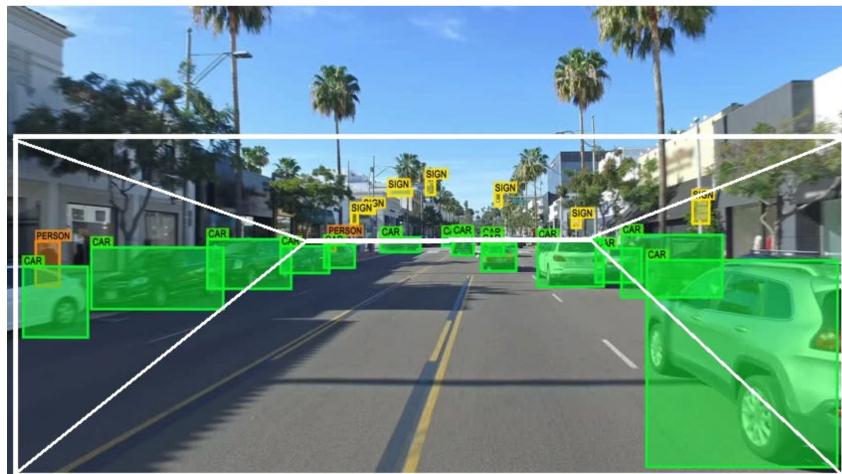
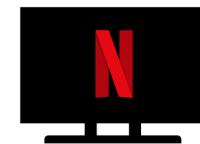
Introduction to ML

Basic supervised models

Optimisation

Generalization

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



ML in science

Generalities

Introduction to ML

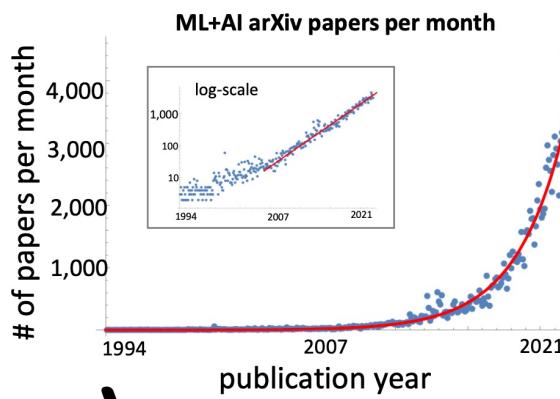
Basic supervised models

Optimisation

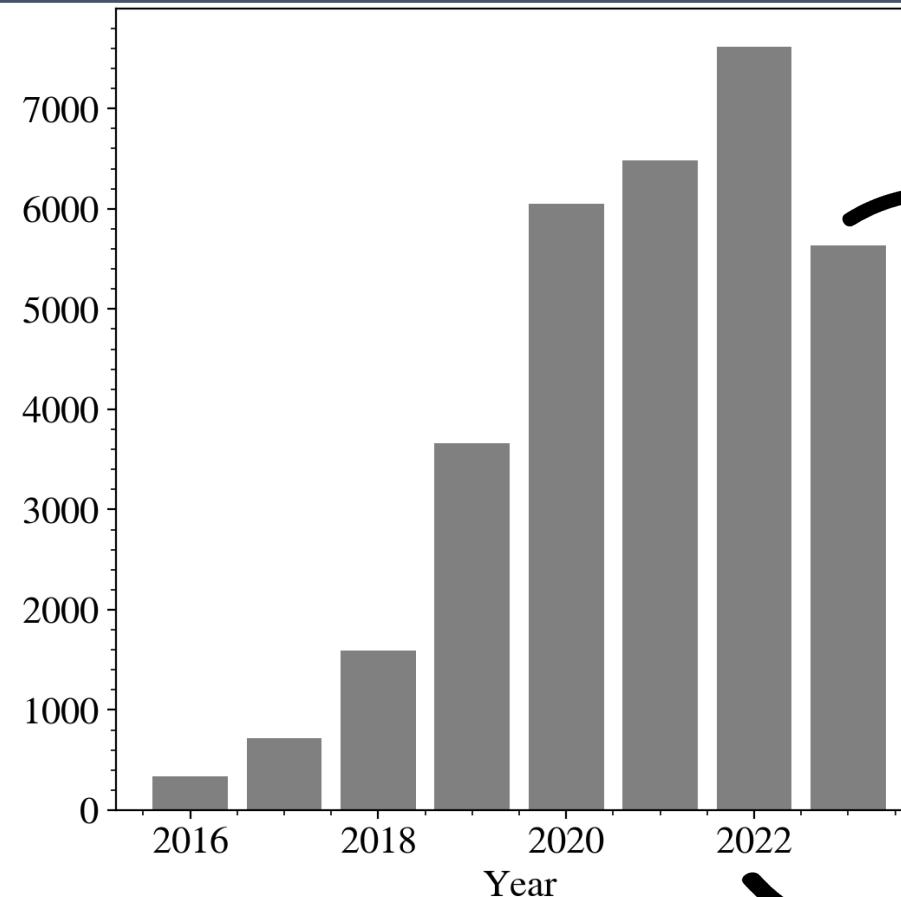
Generalization

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

Generalities

Introduction to ML

Basic supervised models

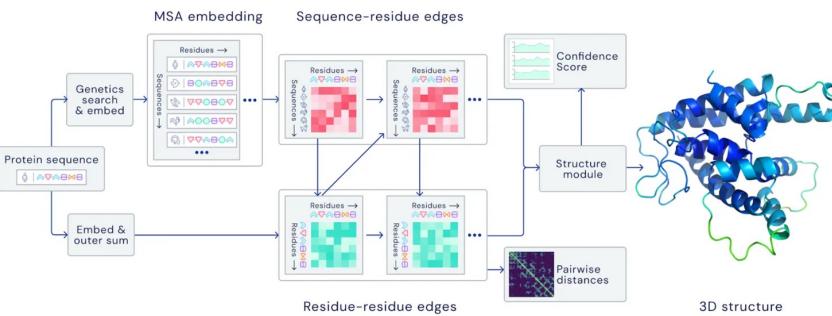
Optimisation

Generalization

But also in science

Healthcare

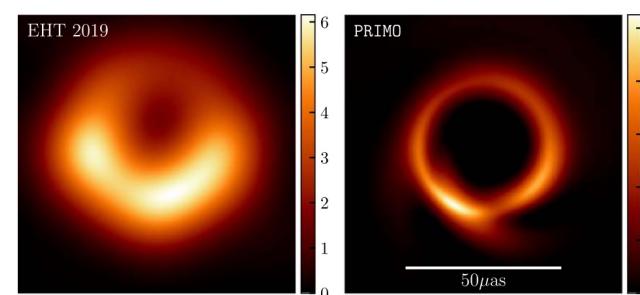
- Drug discovery
- Protein structure reconstruction



[Jumper et al., 2021](#)

Astrophysics and cosmology

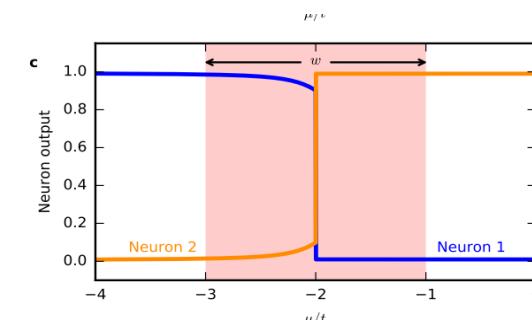
- Galaxy deblending
- Image restoration
- Source separation



[Medeiros et al., 2023](#)

Theoretical physics

- Study phase transitions
- Discover experiments and equations



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

What is “learning”?

5

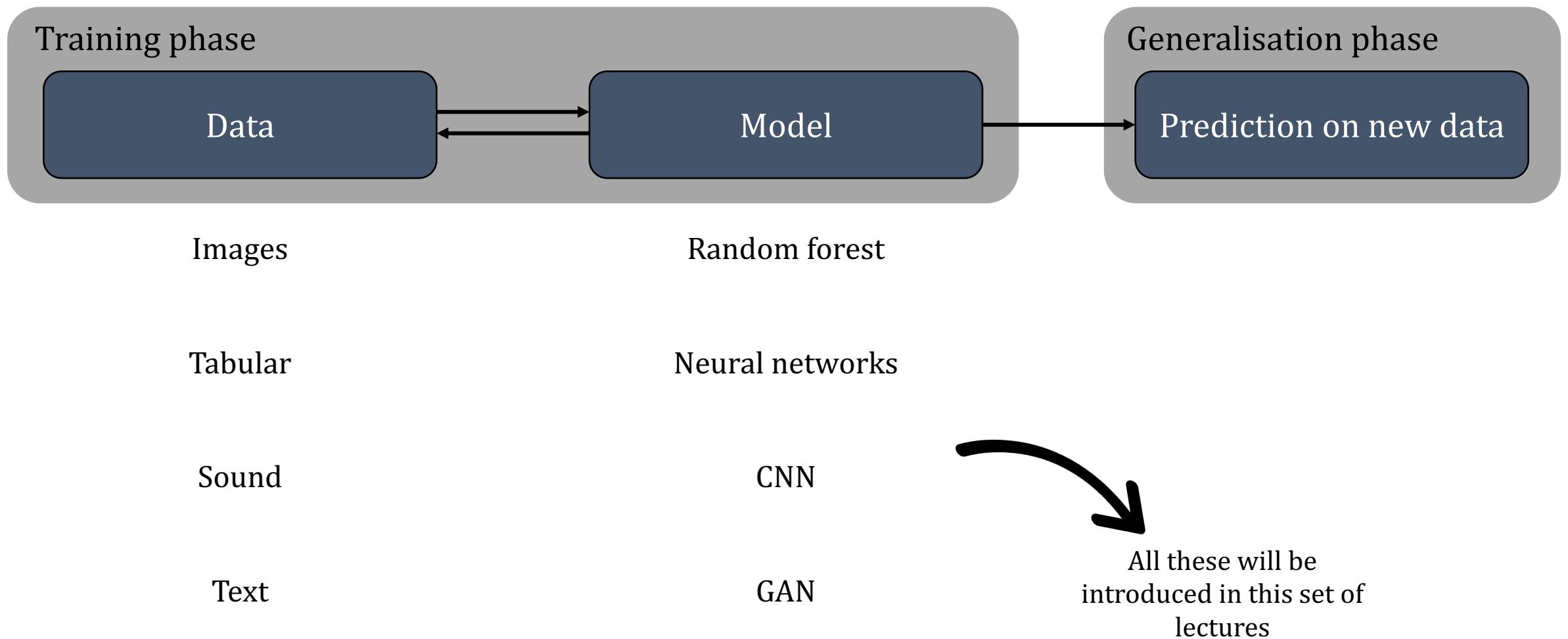
Generalities

Introduction to ML

Basic supervised models

Optimisation of ML models

Generalization



What is “learning”? An example

Generalities

Introduction to ML

Basic supervised models

Optimisation of ML models

Generalization

Training phase

Data

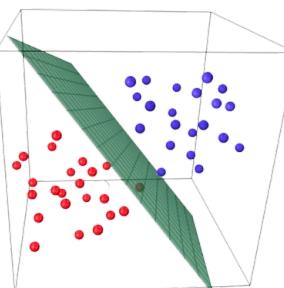
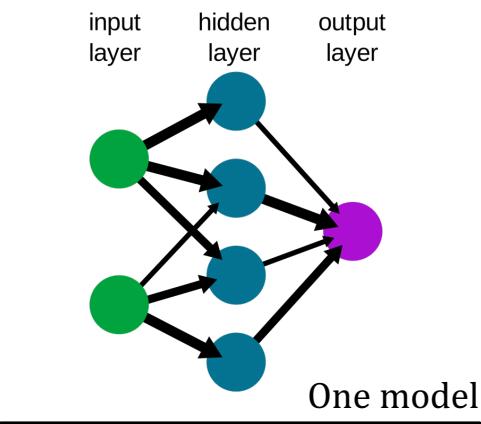
Model

Generalisation phase

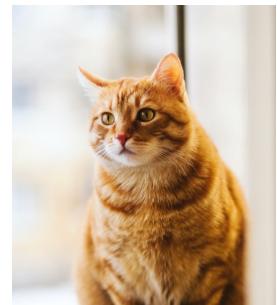
Prediction on new data



Images of a “cat” or “dog”



Another model



“cat” or “dog” ?

ML and the scientific method

Generalities

Introduction to ML

Basic supervised models

Optimisation of ML models

Generalization

Training phase

Data

Model

Generalisation phase

Prediction on new data

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

The scientific method

Measurements from an experiment

Build a model or propose a theory

Test with new experiments

Validate/invalidate your model

ML building blocks

Generalities

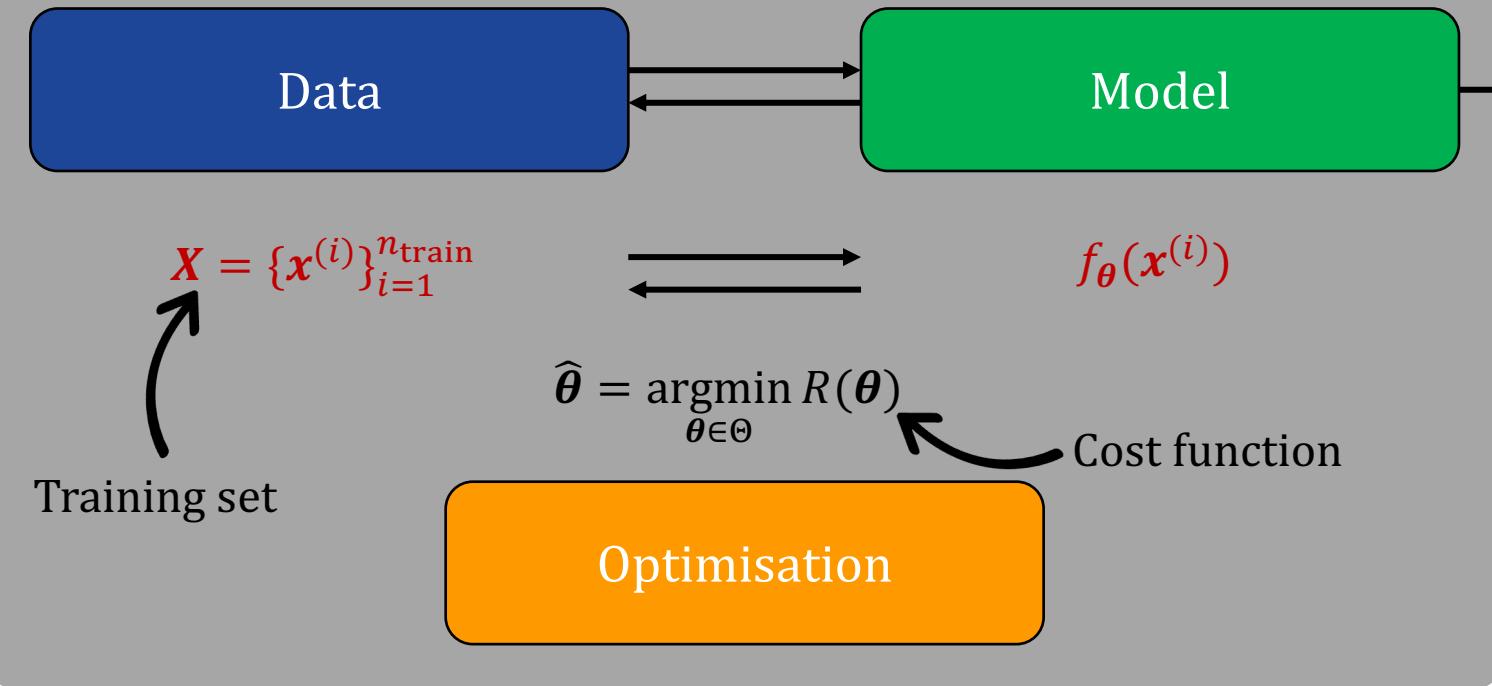
Introduction to ML

Basic supervised models

Optimisation of ML models

Generalization

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

Generalities

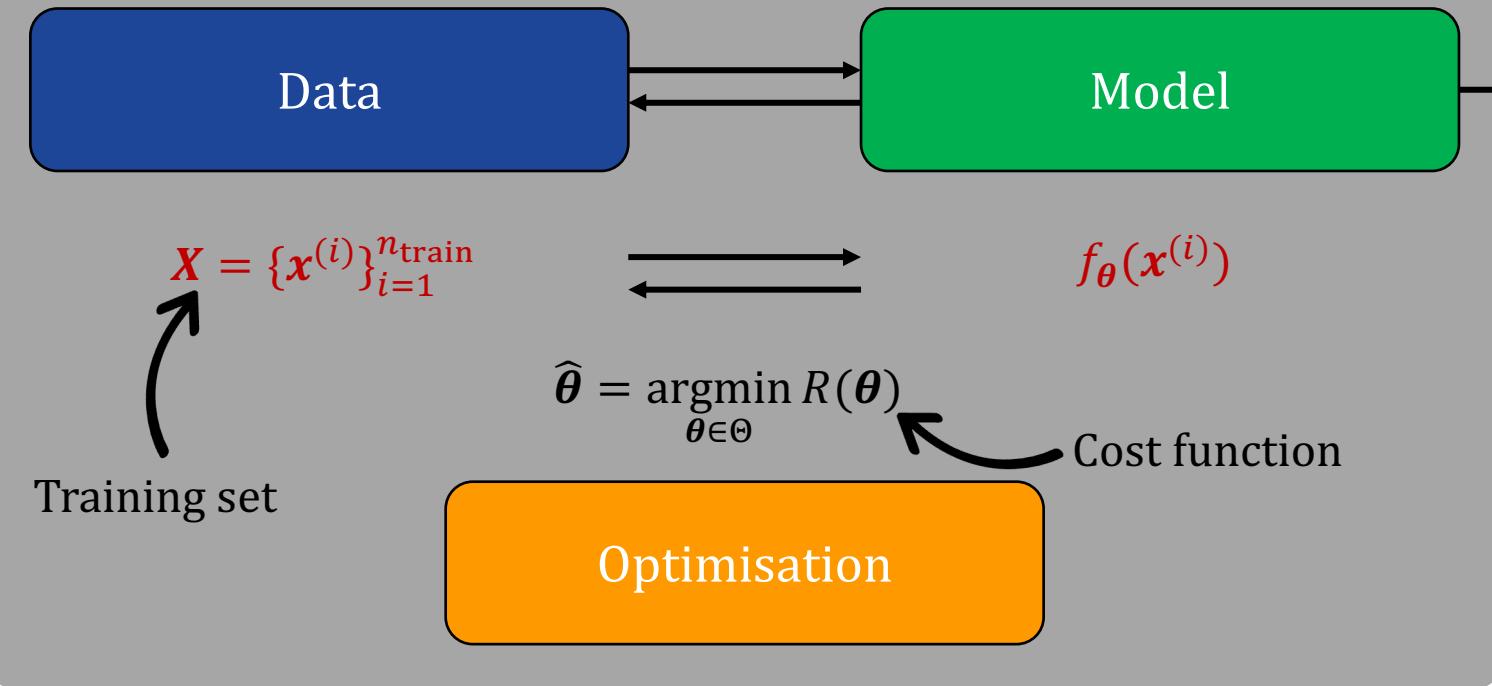
Introduction to ML

Basic supervised models

Optimisation of ML models

Generalization

Training phase



Generalisation phase

Prediction on new data

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

\tilde{X} Test set

Some notations:

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

Usually,

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



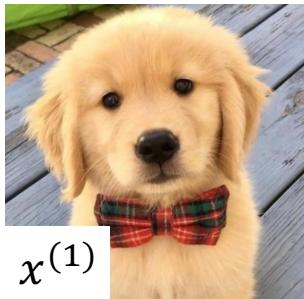
Three main building blocks of *most* ML systems: some **data** feed a **parametrised model**. The best set of parameters are obtained through the **optimisation** of a cost function. The model can then be assessed on **new data**.

Supervised learning

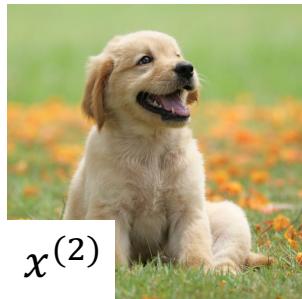
- Data are actually X and Y coming as pairs

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

- 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)

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

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



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

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

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

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

Data

Model

Optimisation

Families of ML problems

Generalities

Introduction to ML

Basic supervised models

Optimisation of ML models

Generalization

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

Generalities

Introduction to ML

Basic supervised models

Optimisation of ML models

Generalization

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)

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



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

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

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

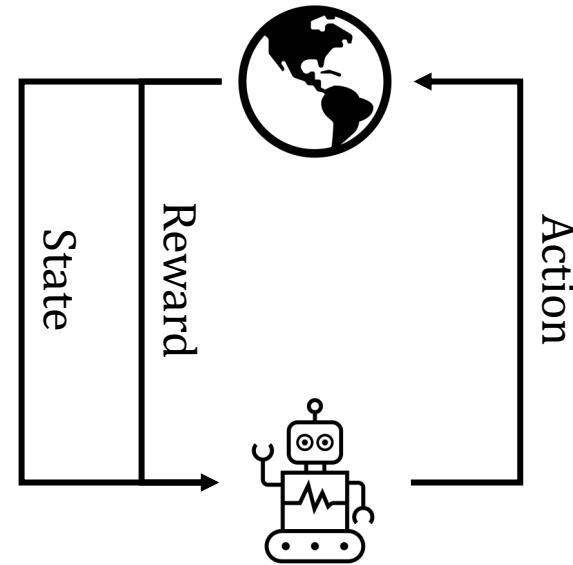
Boltzmann Machine

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

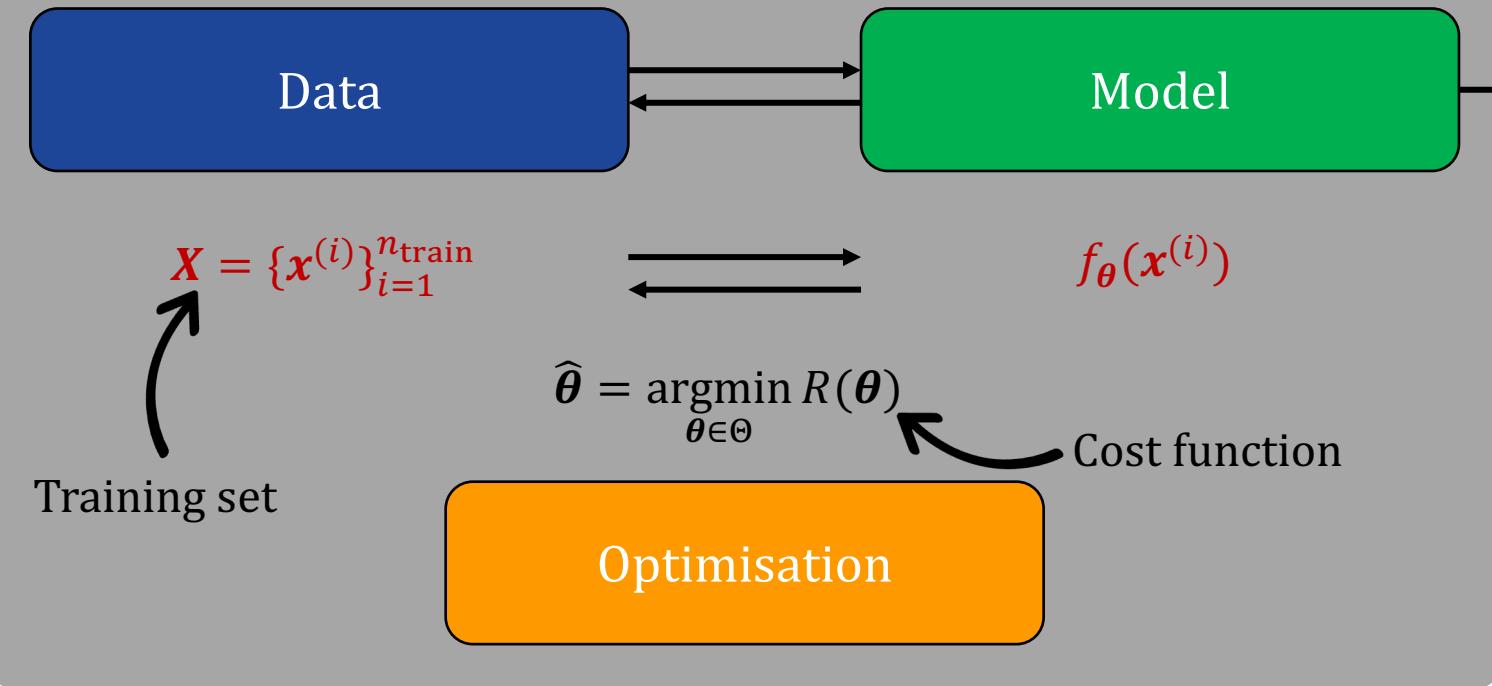
Markov decision processes

Q-networks

Deep policy gradient

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

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**.

The bias-variance trade-off

Generalities

Introduction to ML

Basic supervised models

Optimisation of ML models

Generalization

- 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_{θ} 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

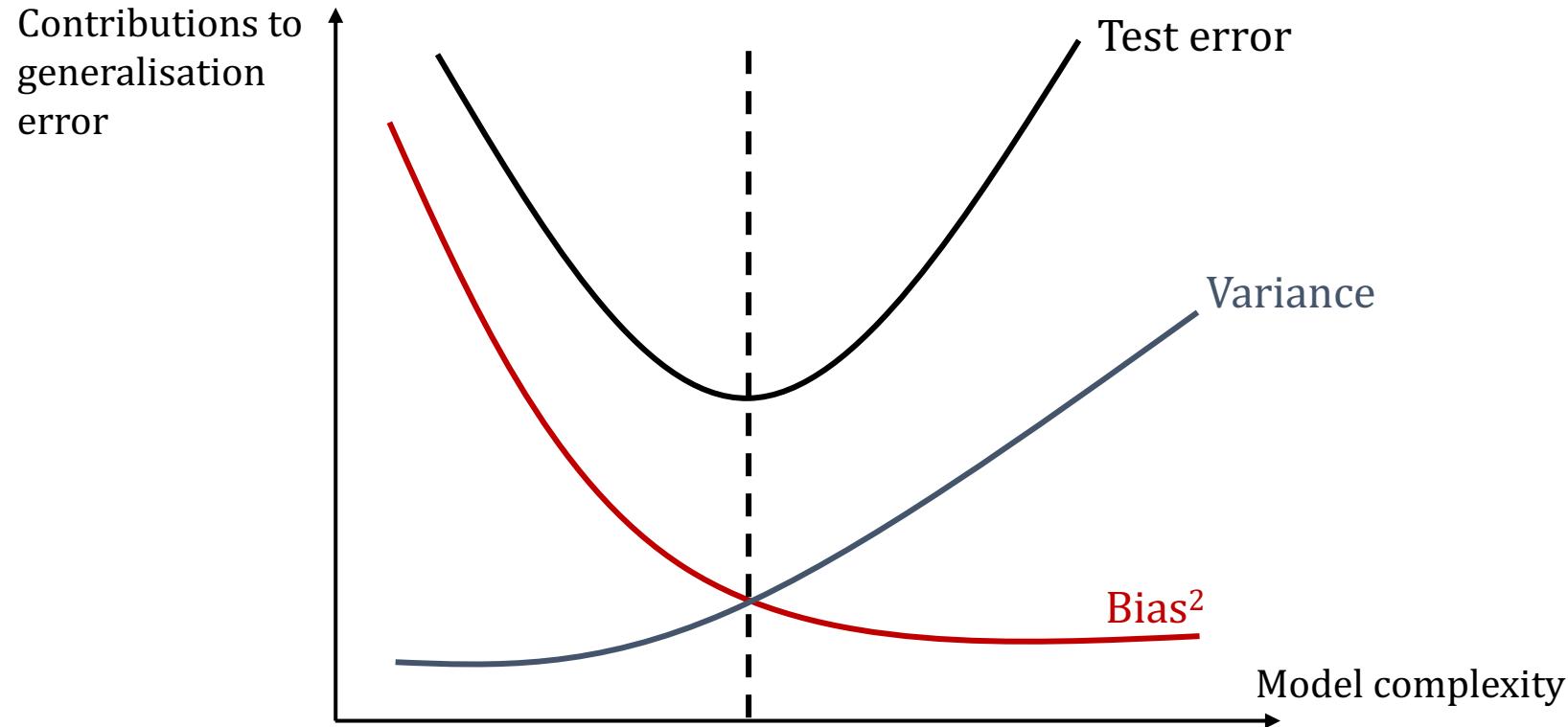
\downarrow
**Modelling
error**

\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



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

Model complexity and generalisation error

18

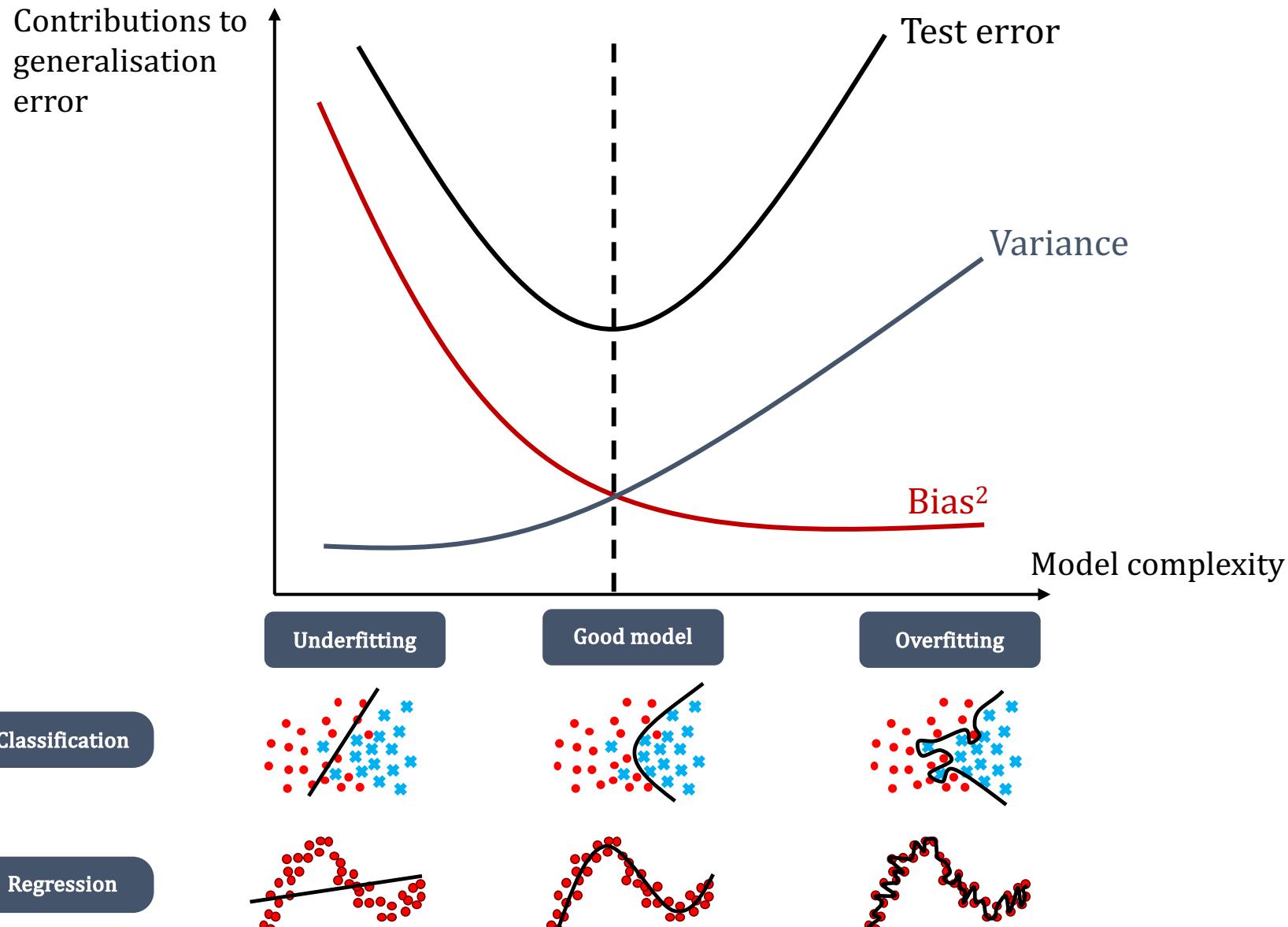
Generalities

Introduction to ML

Basic supervised models

Optimisation of ML models

Generalization



Linear regression: the model

19

Generalities

Introduction to ML

Basic supervised models

Optimisation of ML models

Generalization

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

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

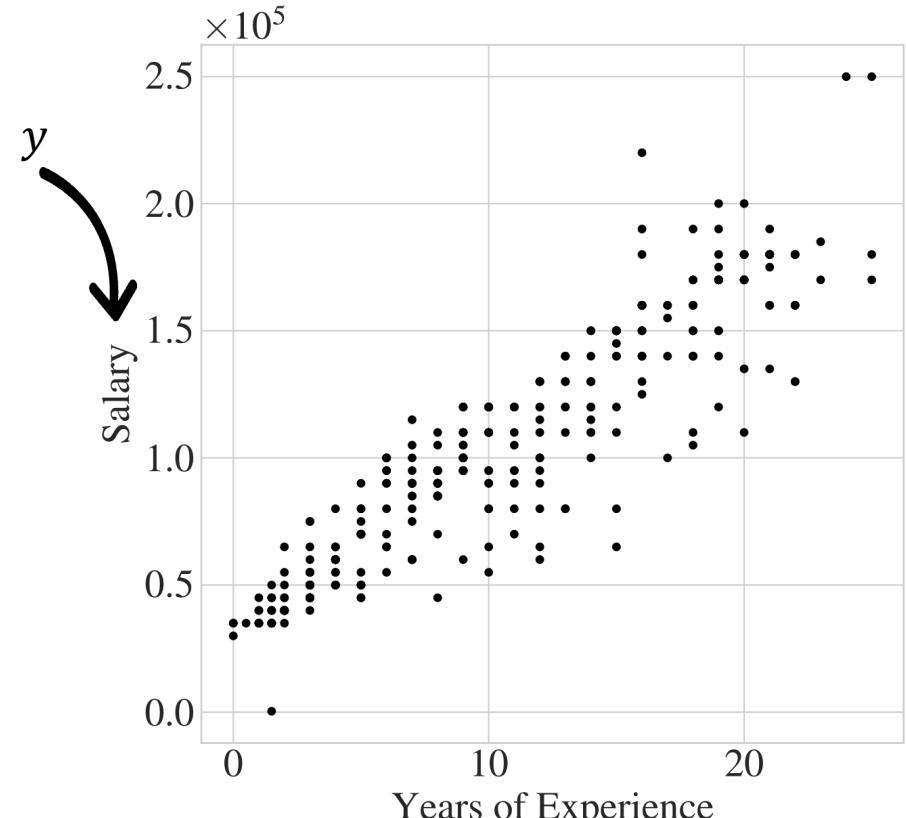
where $x_1^{(i)}$ is the nb. of years of experience of the i^{th} training example

- 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**

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

- Note that this is the maximum likelihood estimator, assuming Gaussian error between $\hat{y}^{(i)} = f_{\theta}(x^{(i)})$ and $y^{(i)}$
- Here**, the optimisation problem can be solved analytically in closed-

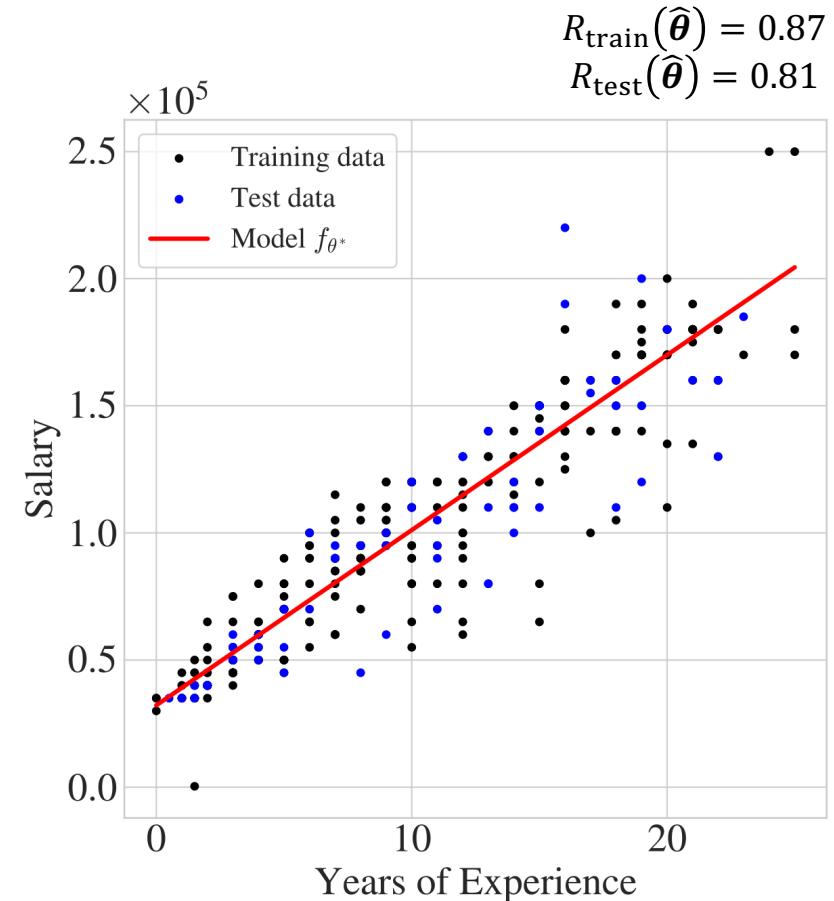
form giving (usually not the case): $\hat{\theta} = (X^T X)^{-1} X^T y$



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

A classification task: the XOR dataset

Generalities

Introduction to ML

Basic supervised models

Optimisation of ML models

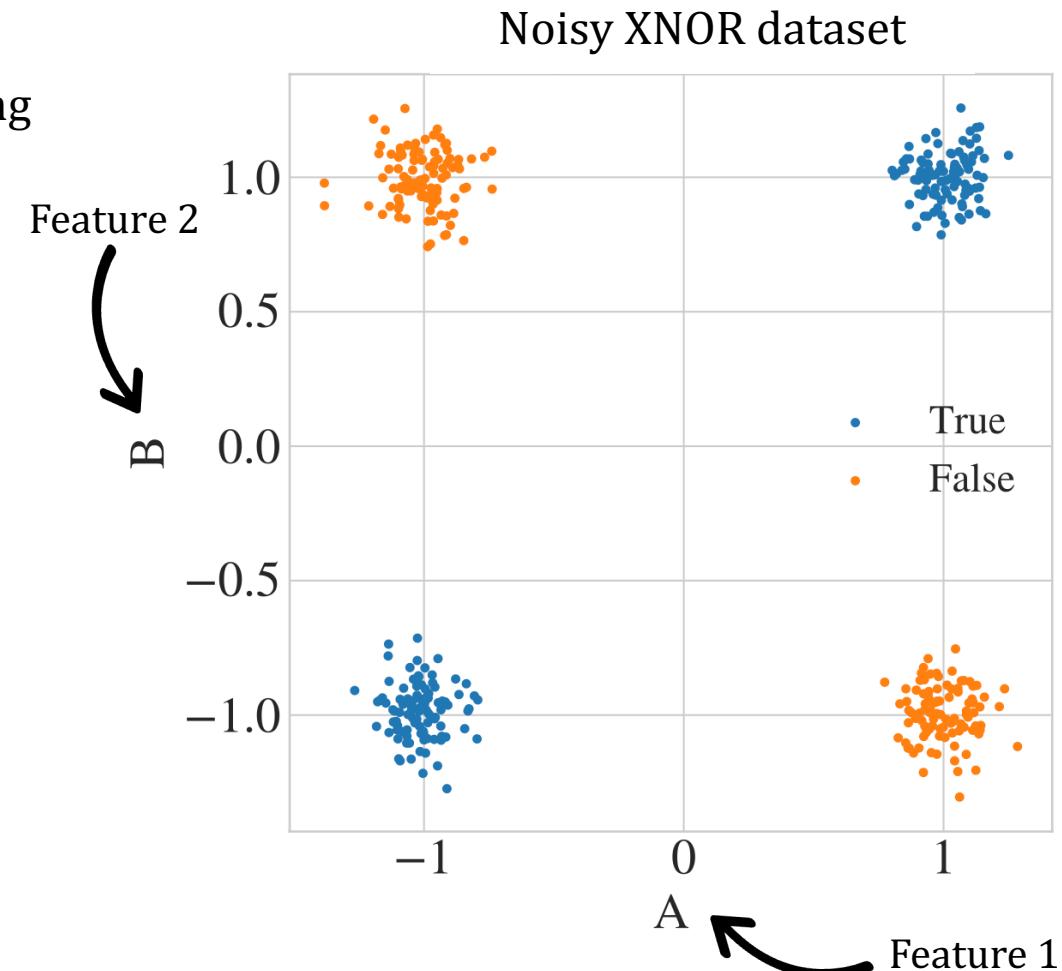
Generalization

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

Generalities

Introduction to ML

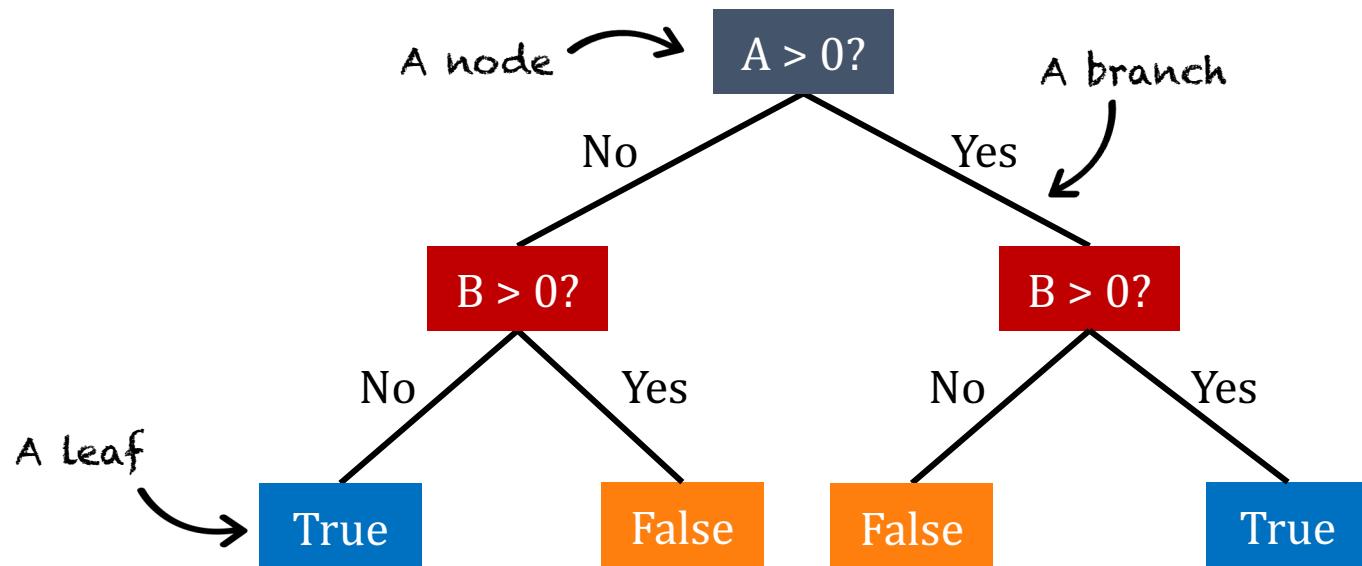
Basic supervised models

Optimisation of ML models

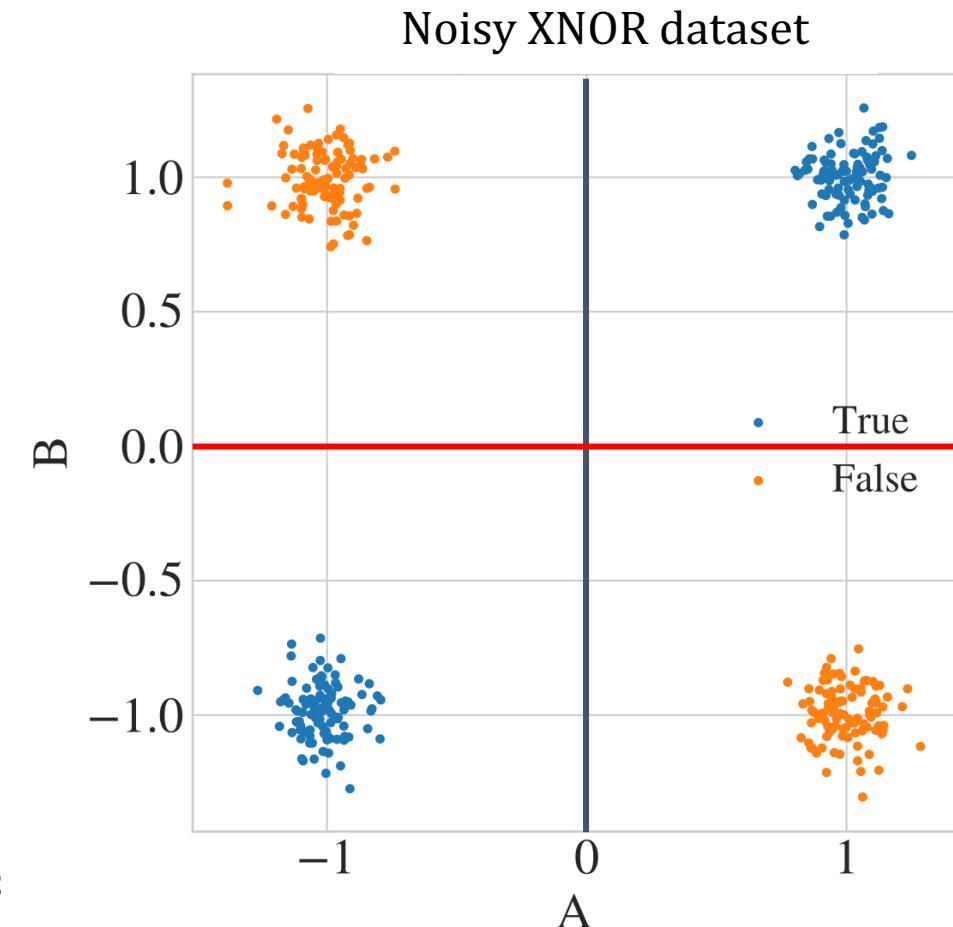
Generalization

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

Generalities

Introduction to ML

Basic supervised models

Optimisation of ML models

Generalization

Decision trees

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

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

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

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

$$L(R_i) = - \sum_{k=1}^q \rho_k^i \log_2 \rho_k^i, \quad \rho_k^i = \frac{|\{\mathbf{x}^{(j)} \mid \mathbf{x}^{(j)} \in R_i, y^{(j)} = k\}|}{|\{\mathbf{x}^{(j)} \mid \mathbf{x}^{(j)} \in R_i\}|}$$

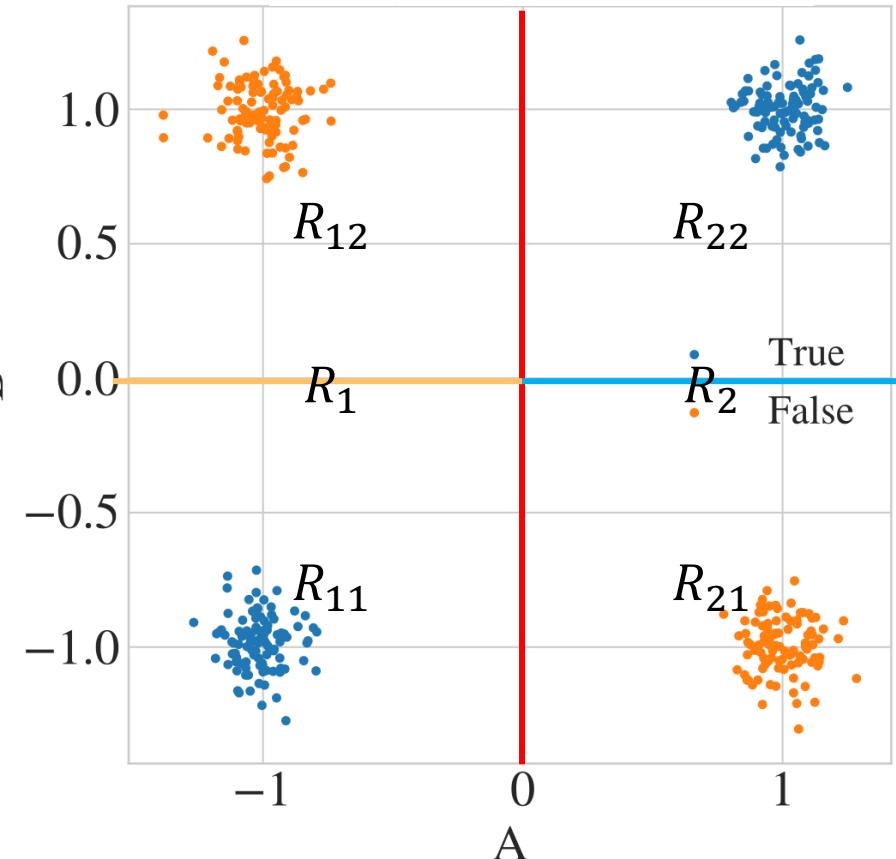


can handle categorical values, easy to interpret, fast to compute, both regression and classification



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

Noisy XNOR dataset



A more formal view of decision trees

Generalities

Introduction to ML

Basic supervised models

Optimisation of ML models

Generalization

Decision trees

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

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

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

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

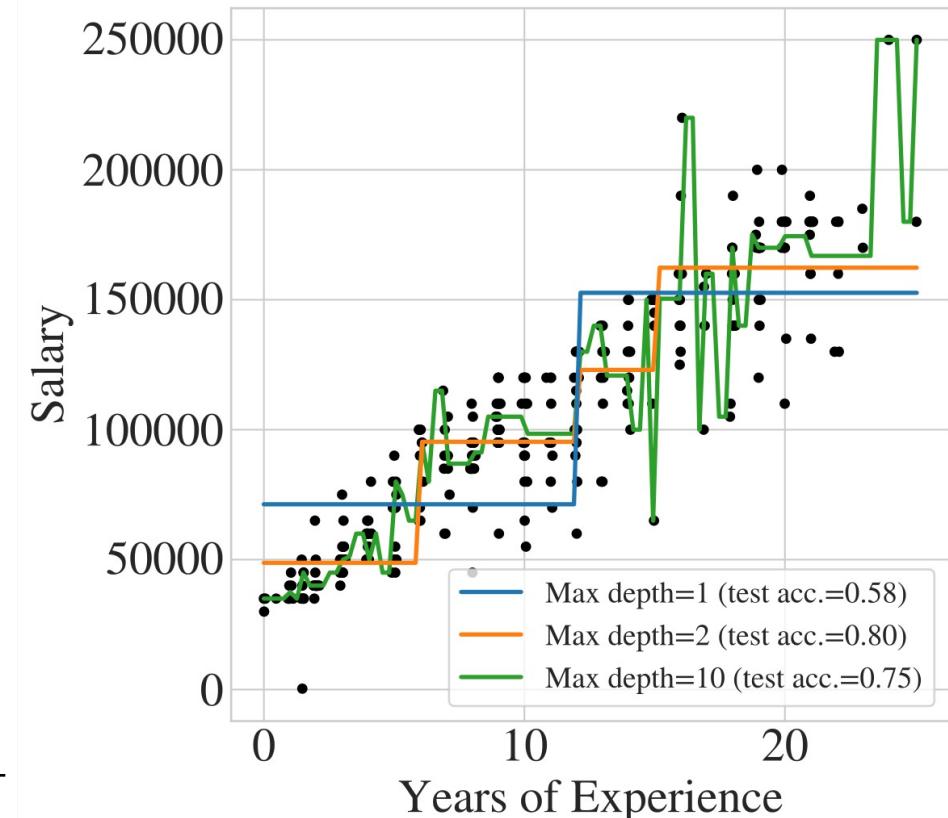
$$L(R_i) = - \sum_{k=1}^q \rho_k^i \log_2 \rho_k^i, \quad \rho_k^i = \frac{|\{x^{(j)} \mid x^{(j)} \in R_i, y^{(j)} = k\}|}{|\{x^{(j)} \mid x^{(j)} \in R_i\}|}$$



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

25

Generalities

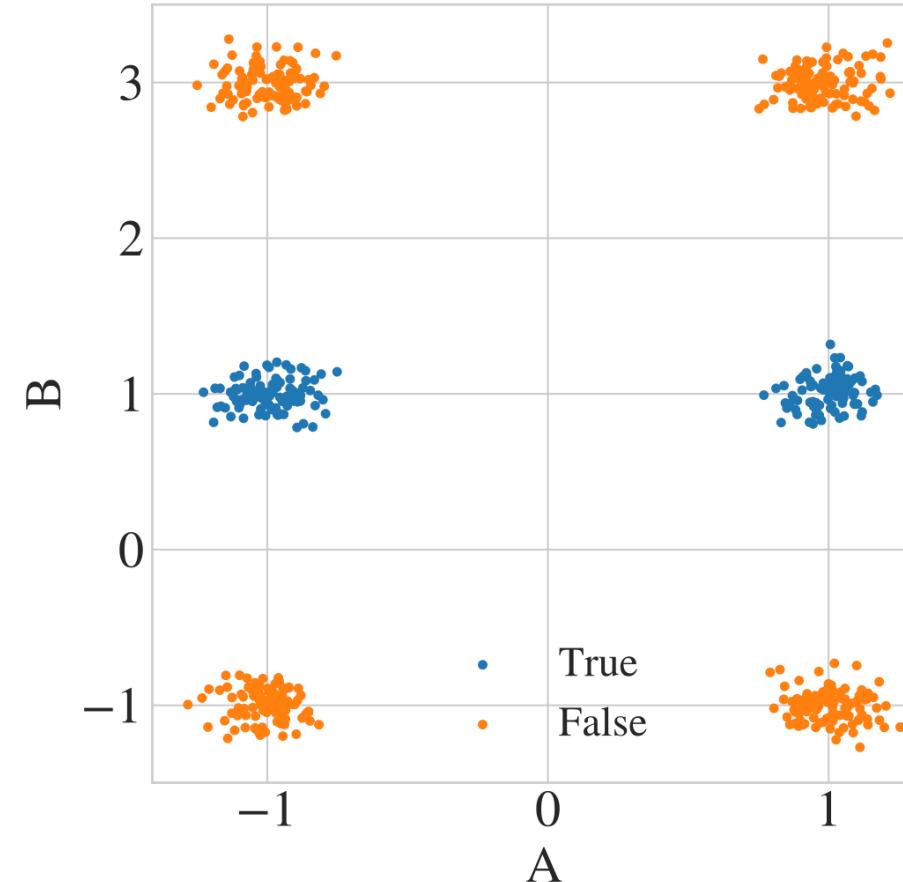
Introduction to ML

Basic supervised models

Optimisation of ML models

Generalization

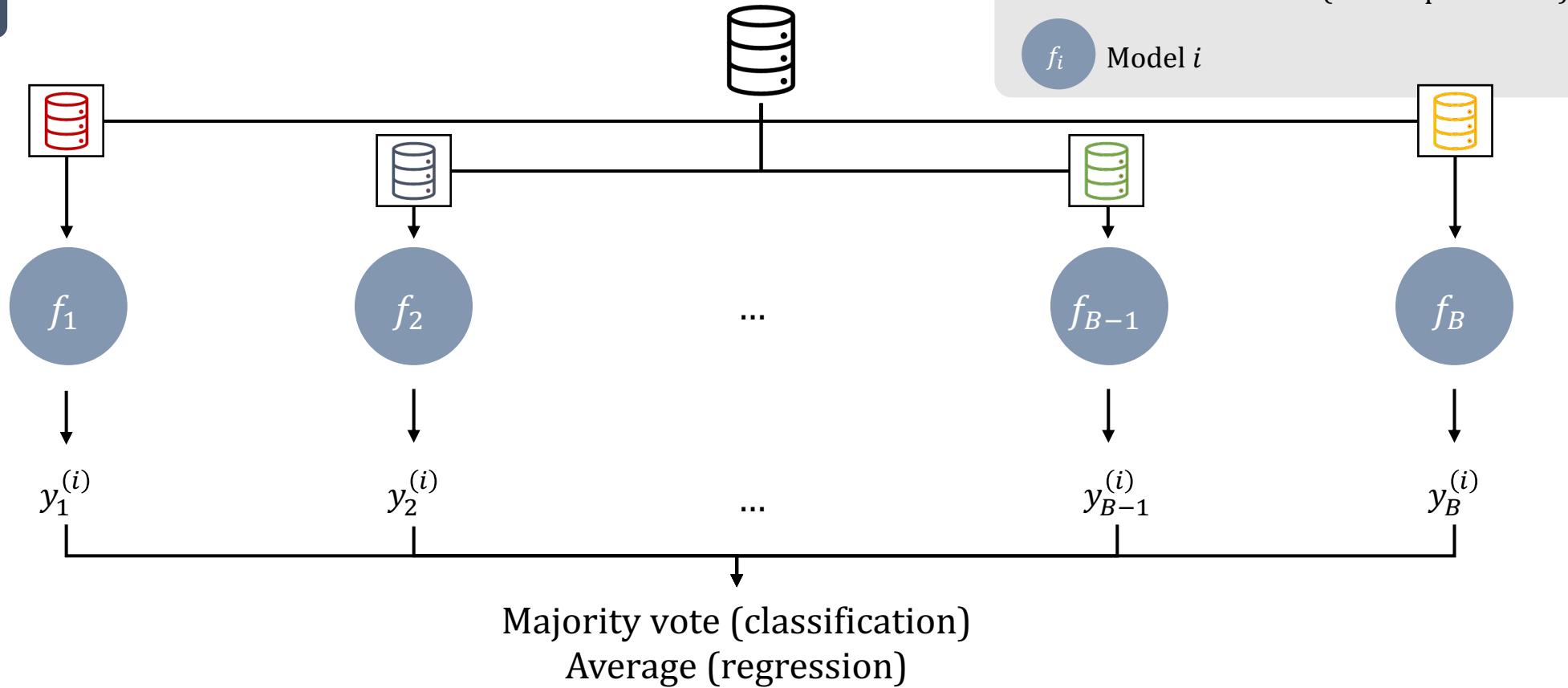
Decision trees



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

27

Generalities

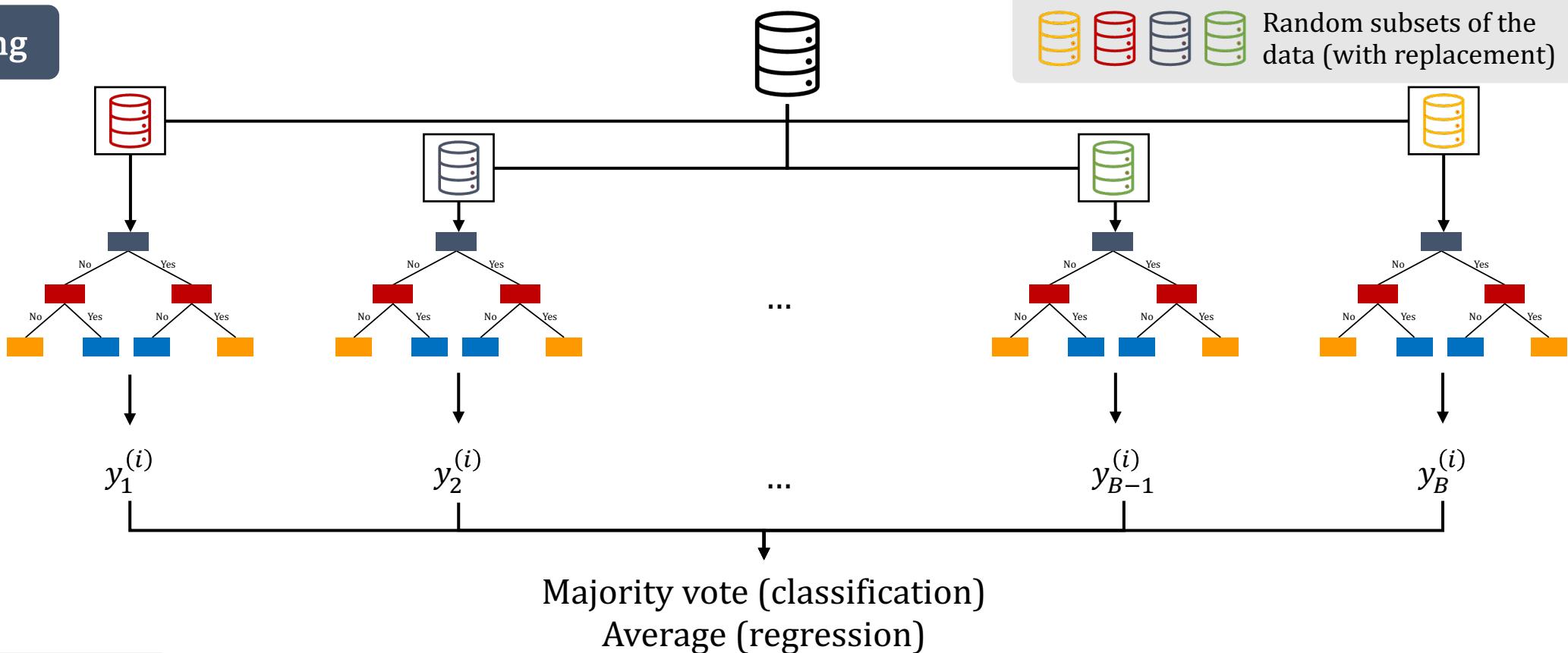
Introduction to ML

Basic supervised models

Optimisation of ML models

Generalization

Bagging



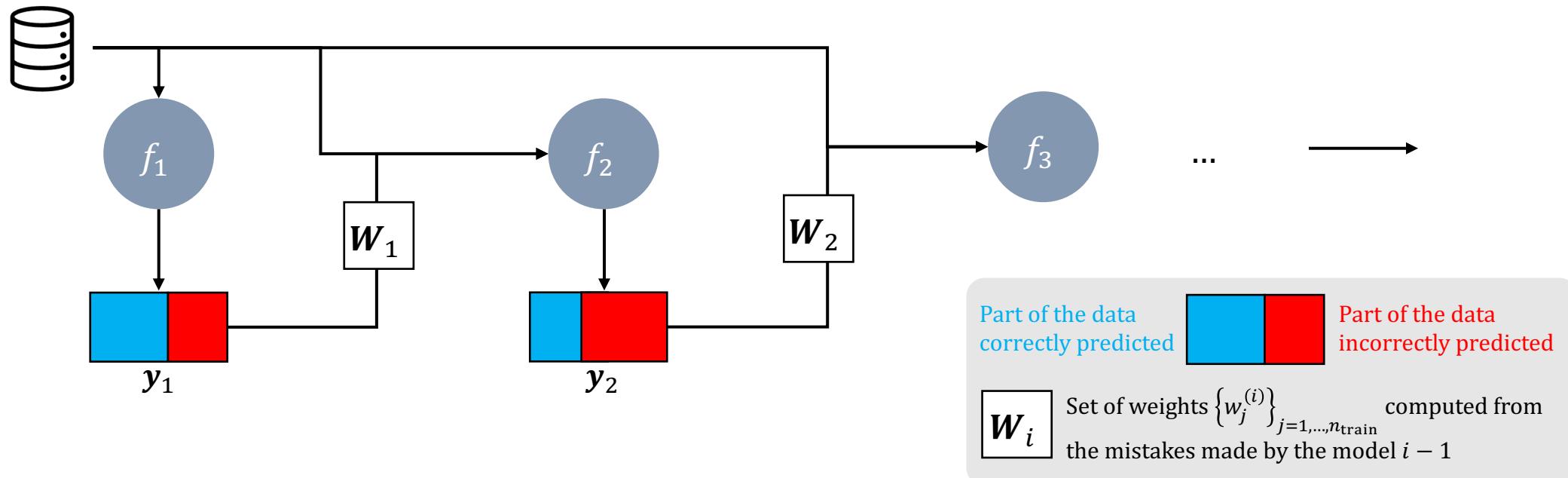
Random Forest

= 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

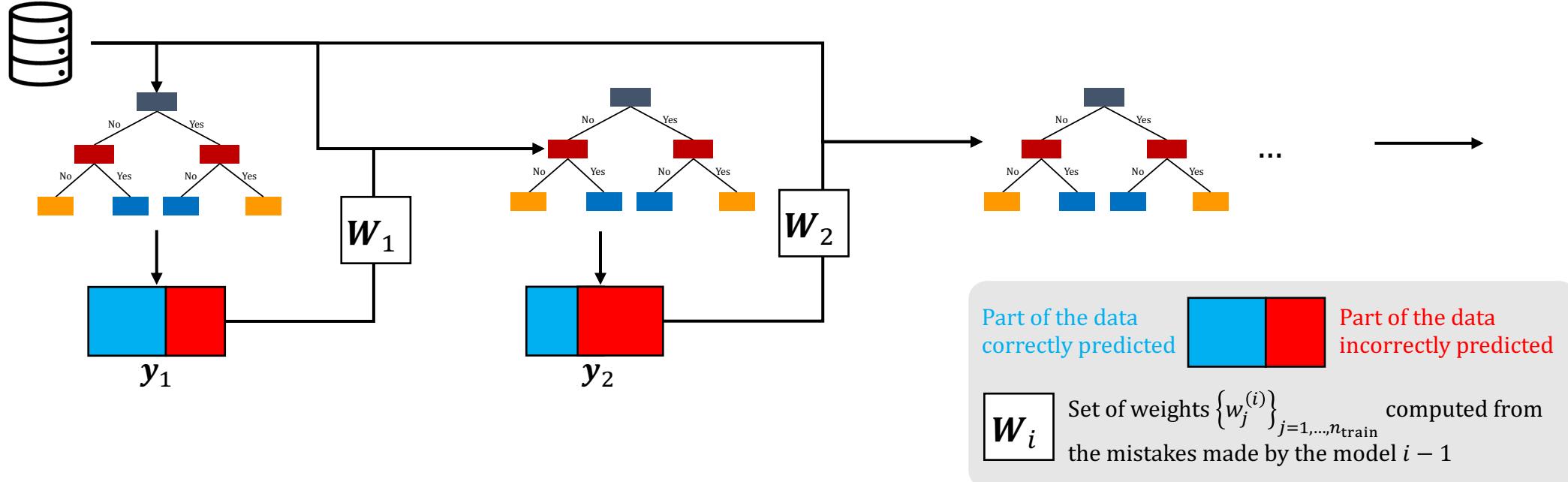
- While bagging trains high-variance models in parallel to reduce the variance of the combined estimate, **boosting trains high-bias models sequentially to reduce the overall bias**

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

30

Generalities

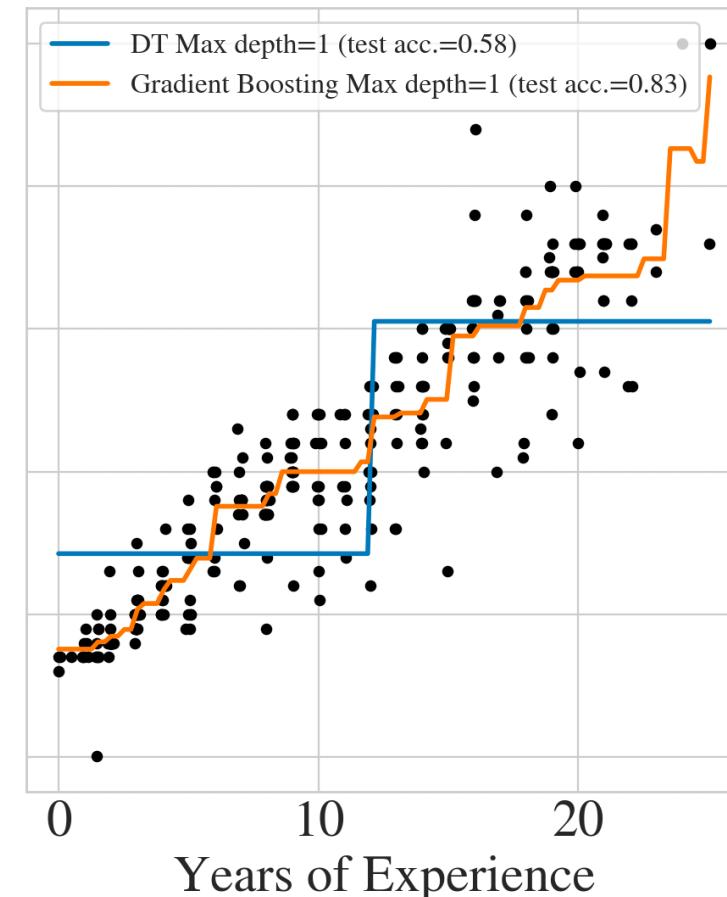
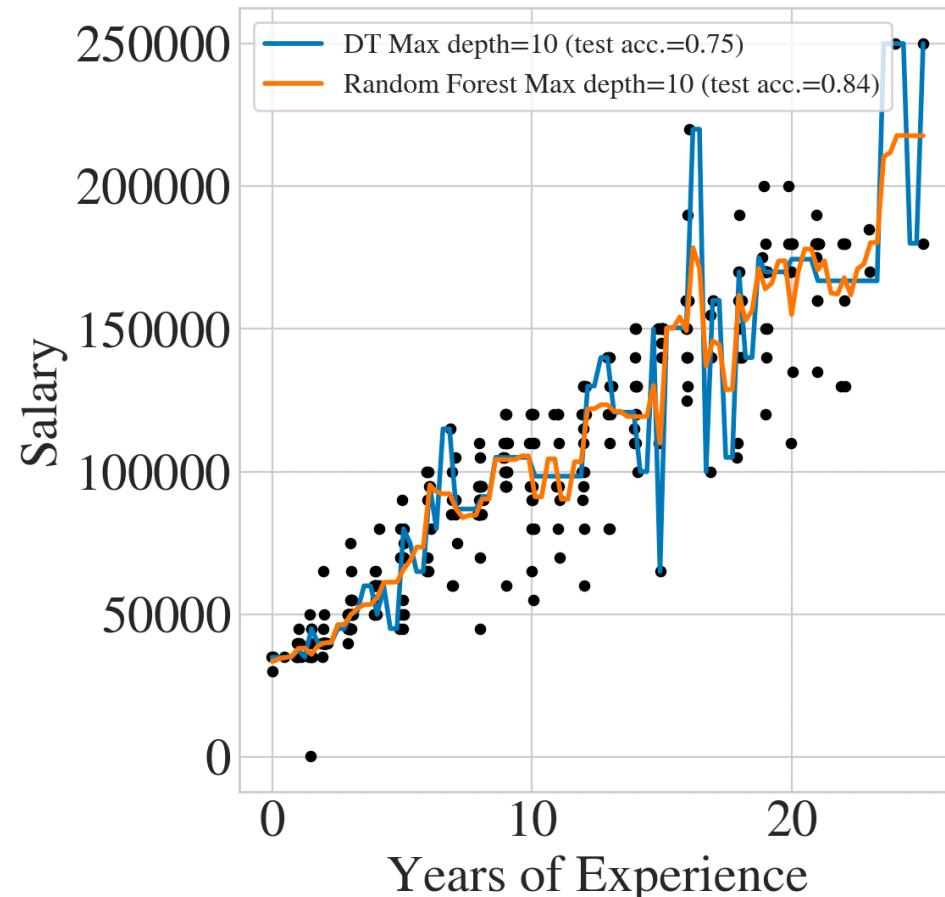
Introduction to ML

Basic supervised models

Optimisation of ML models

Generalization

Comparison on our regression problem



Artificial neural networks: neurons

Generalities

Introduction to ML

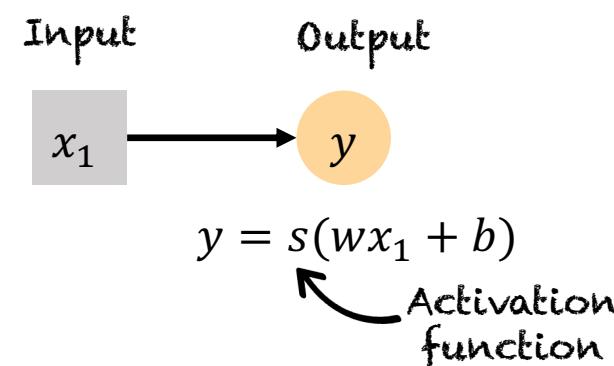
Basic supervised models

Optimisation of ML models

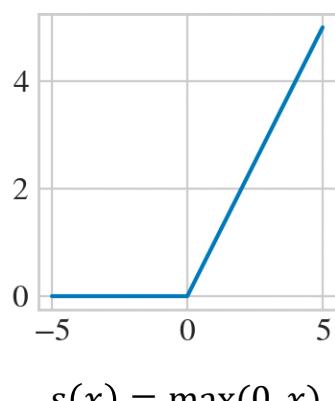
Generalization

Artificial neural networks

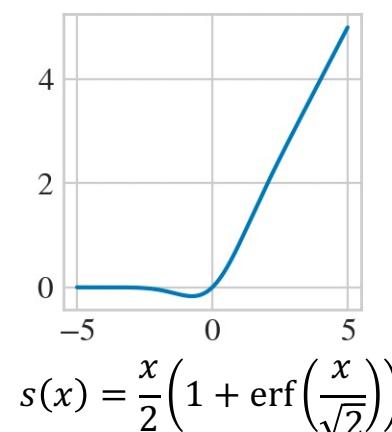
- Building block of neural networks: the **neuron**
- Made of three operations: it first multiplies the input by a **weight**, then adds a **bias**, and finally applies an **activation function** s to the result
- If s is the identity, you recognize the linear regression model with $\theta_0 = b$ and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



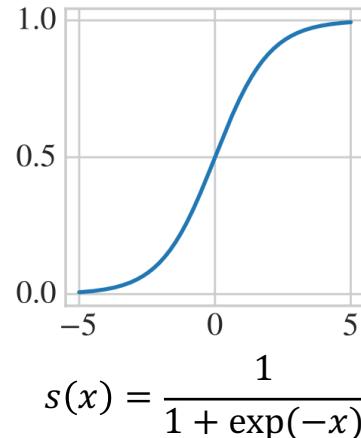
ReLU



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

32

Generalities

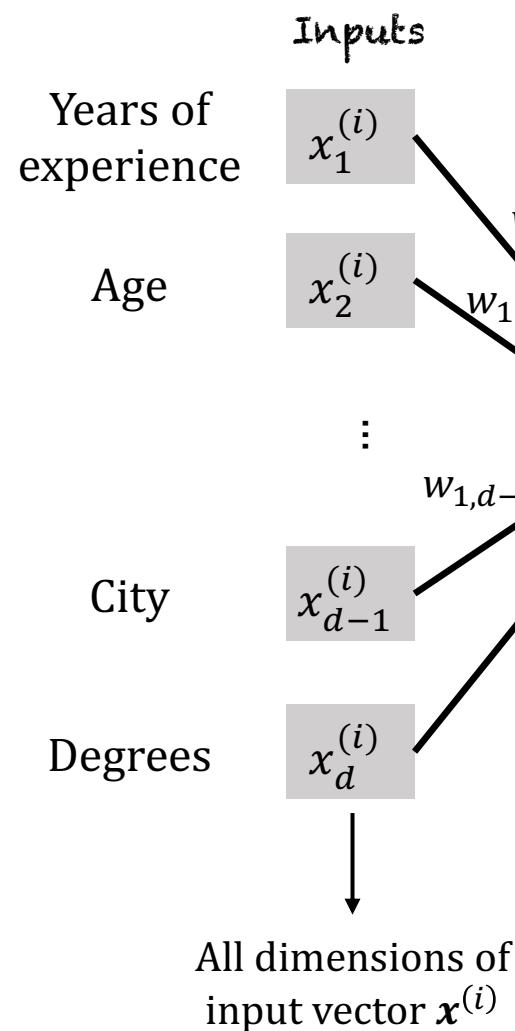
Introduction to ML

Basic supervised models

Optimisation of ML models

Generalization

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

33

Generalities

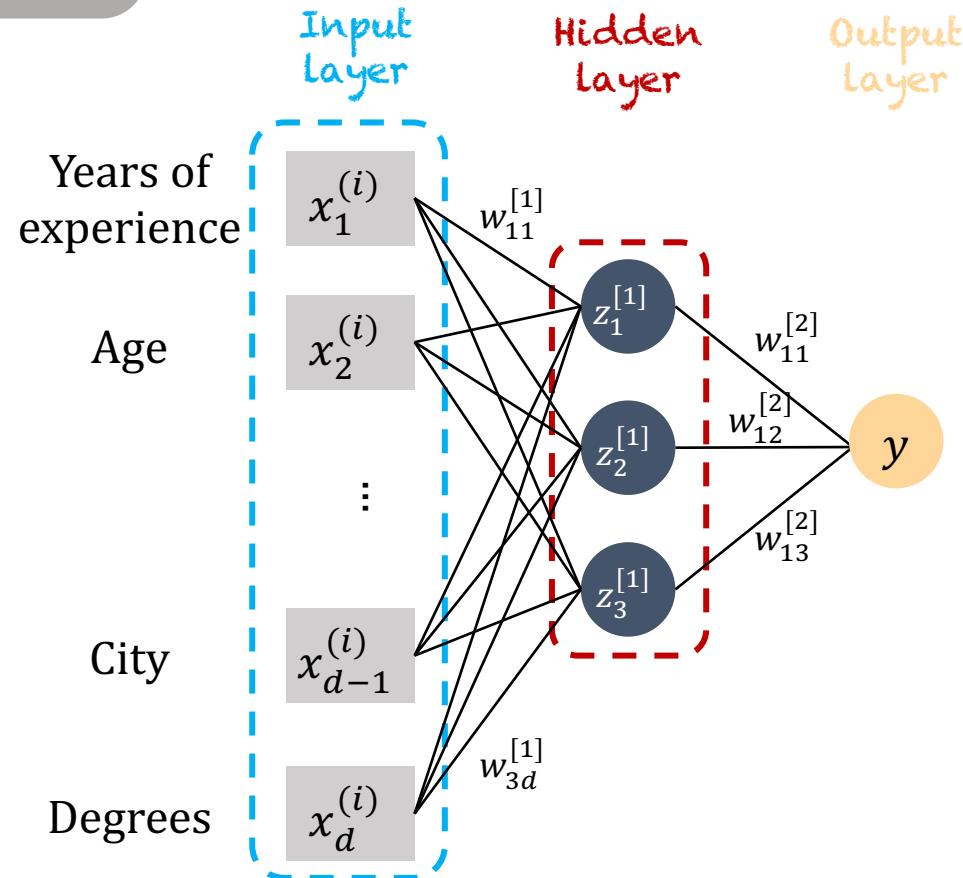
Introduction to ML

Basic supervised models

Optimisation of ML models

Generalization

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

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

Annotations for the equation:

- Layer 1**: Points to the first term $\mathbf{w}_j^{[1]}$.
- Weights of Layer 1 for unit j (vectorized)**: Points to the term $\mathbf{w}_j^{[1]}$.
- Bias of the unit j in layer 1**: Points to the term $b_j^{[1]}$.
- Unit j** : Points to the index j in $z_j^{[1]}$.

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

$$u_j^{[1]} = \mathbf{w}_j^{[1]} \cdot \mathbf{x} + b_j^{[1]}$$

More layers: the art of learning features

34

Generalities

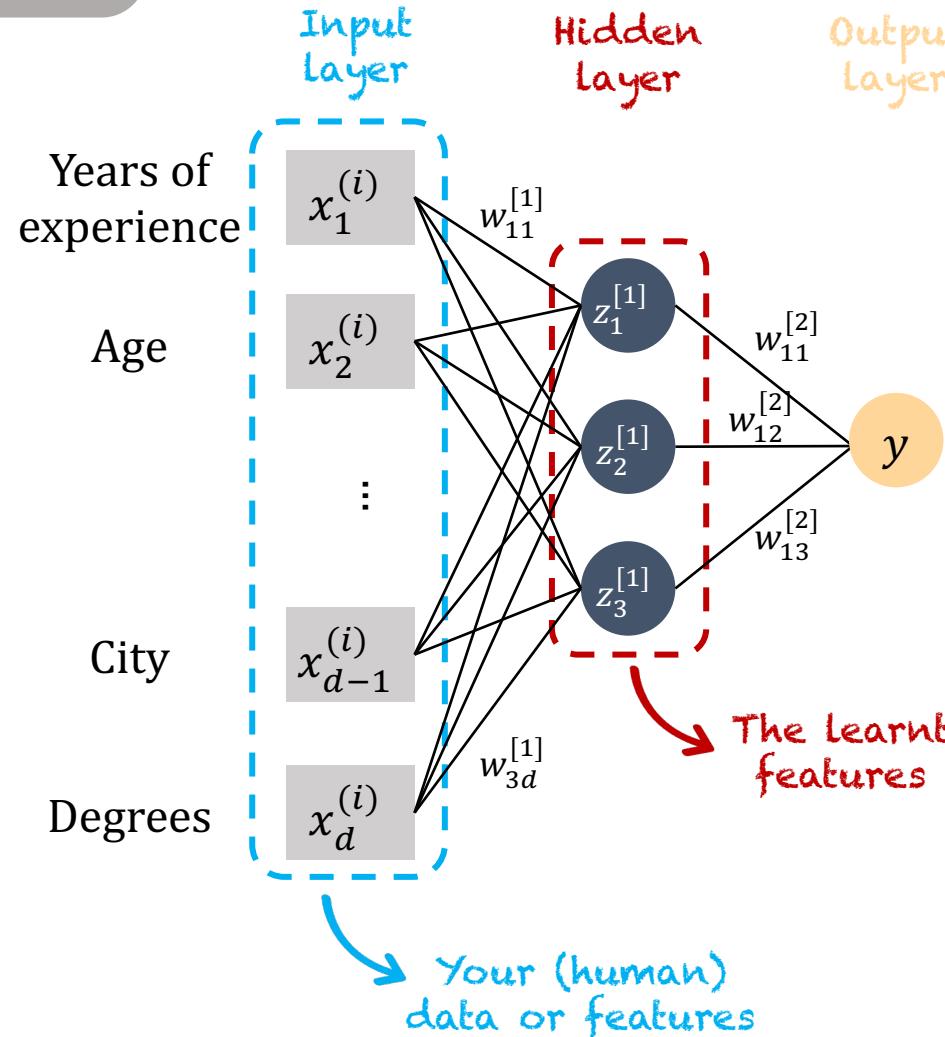
Introduction to ML

Basic supervised models

Optimisation of ML models

Generalization

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

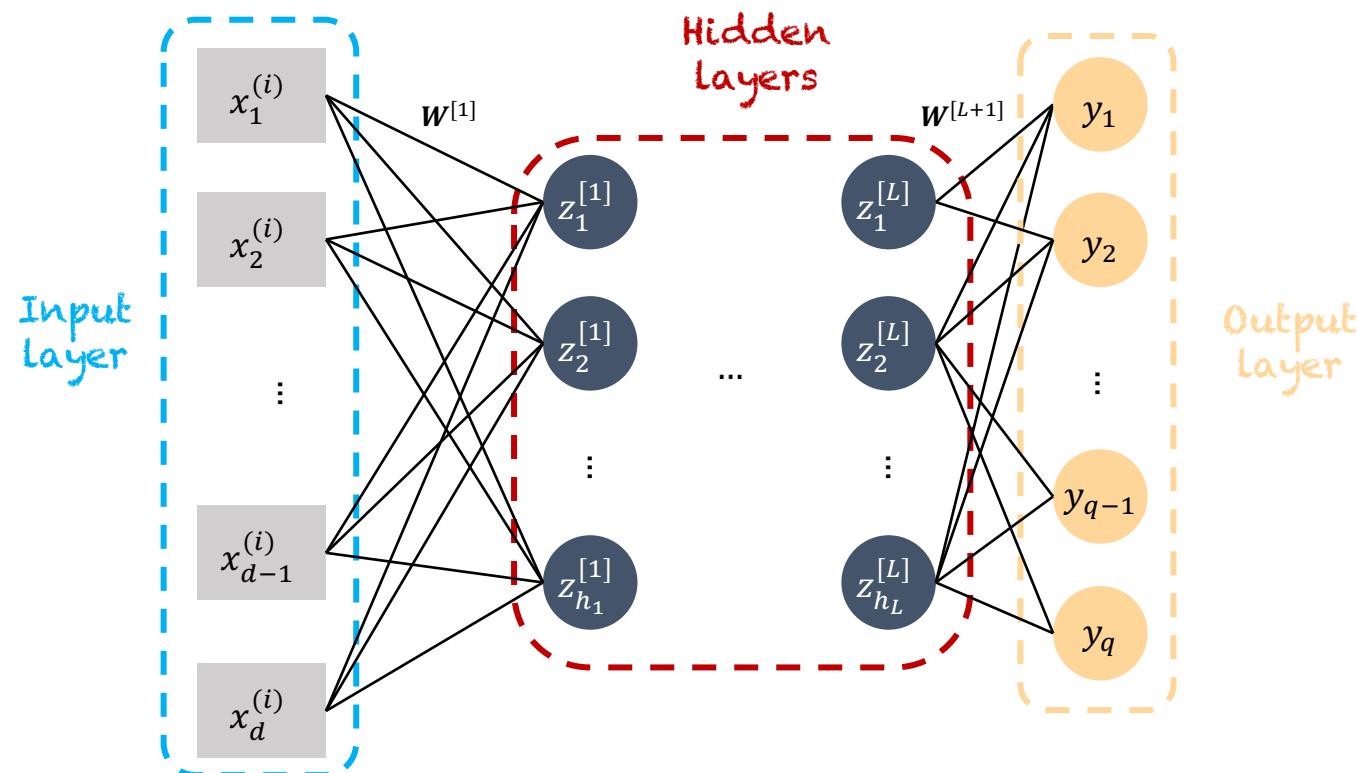
- A fully-connected neural network with d inputs, L hidden layers of width h_1, h_2, \dots, h_L and an output layer of size q
- The j th output is computed as

$$y_j = s^{[L+1]}(\mathbf{W}^{[L+1]}\mathbf{z}^{[L]})$$

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

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

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



- In the end, a **neural network** is a function $f_\theta: X \rightarrow Y$ of some parameters ($\theta = \text{weights and biases}$)
- f_θ 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

36

Generalities

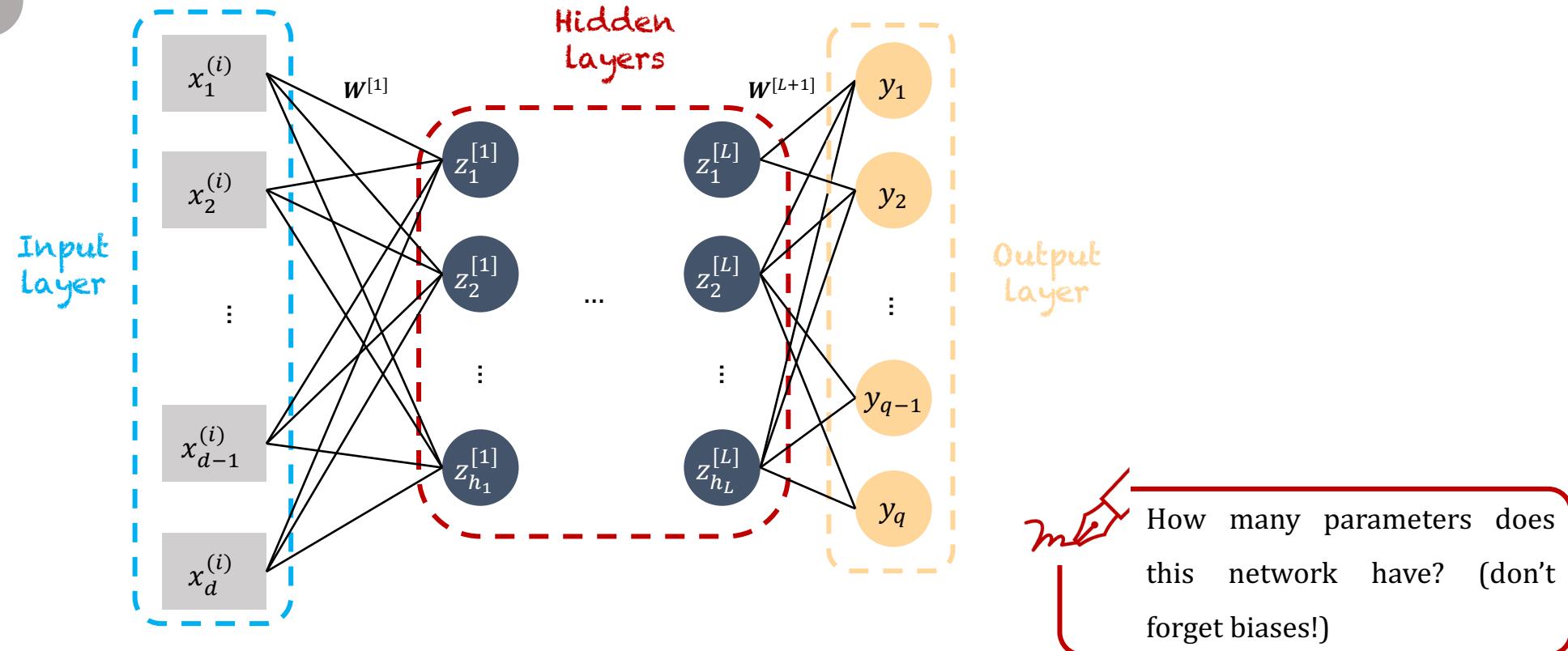
Introduction to ML

Basic supervised models

Optimisation of ML models

Generalization

Artificial neural networks



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

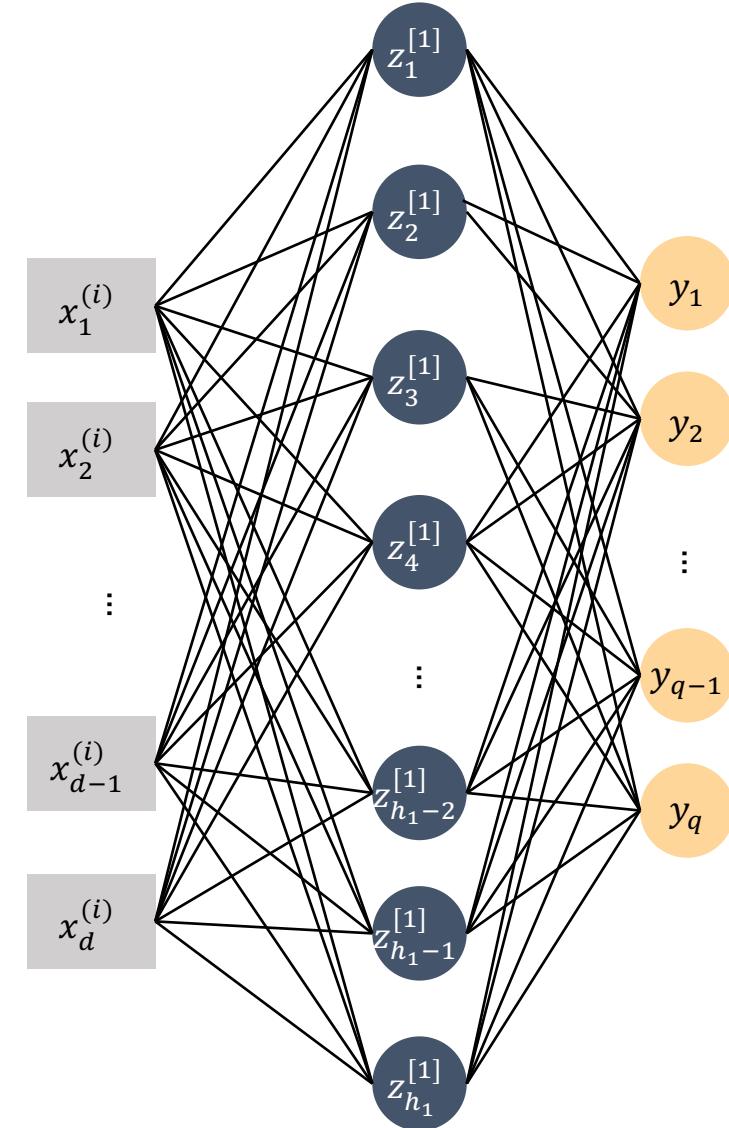
- Need a lot of data to train, subject to overfitting, uninterpretable, targeted-purpose architectures usually work better

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 XOR dataset

38

Generalities

Introduction to ML

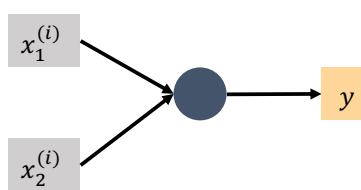
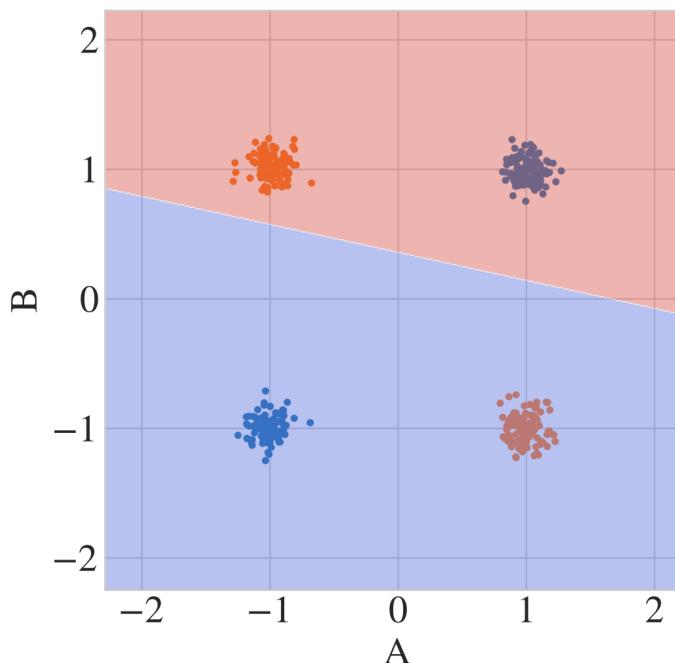
Basic supervised models

Optimisation of ML models

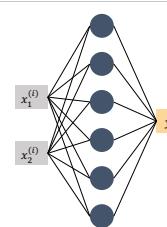
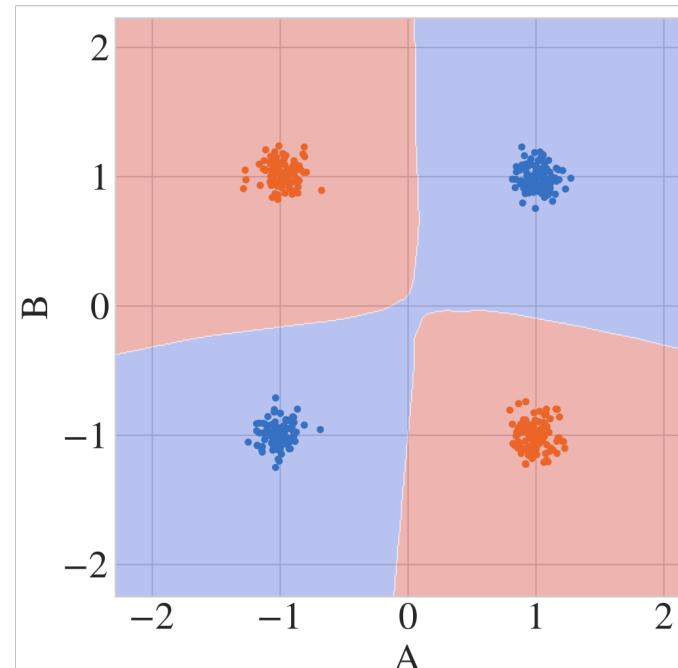
Generalization

Artificial neural networks

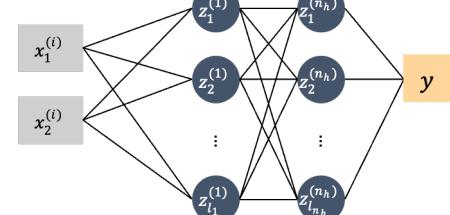
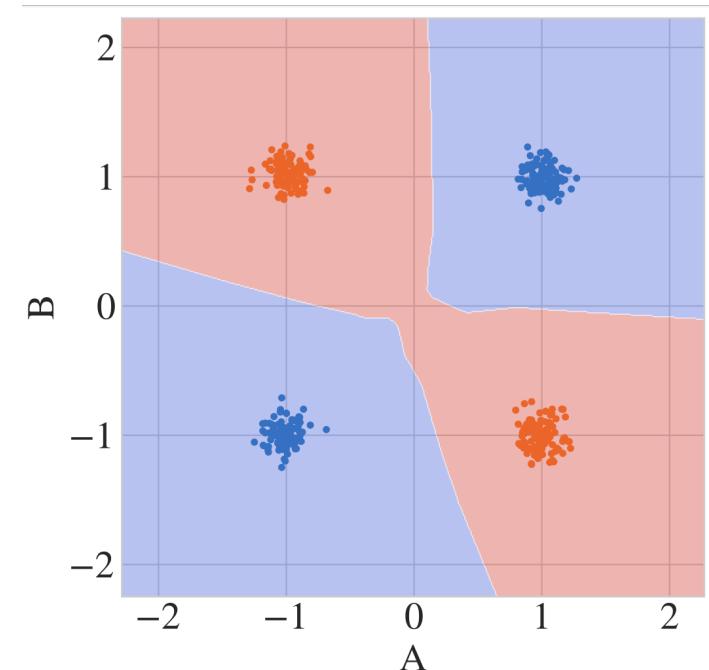
1 hidden layer of 1 unit



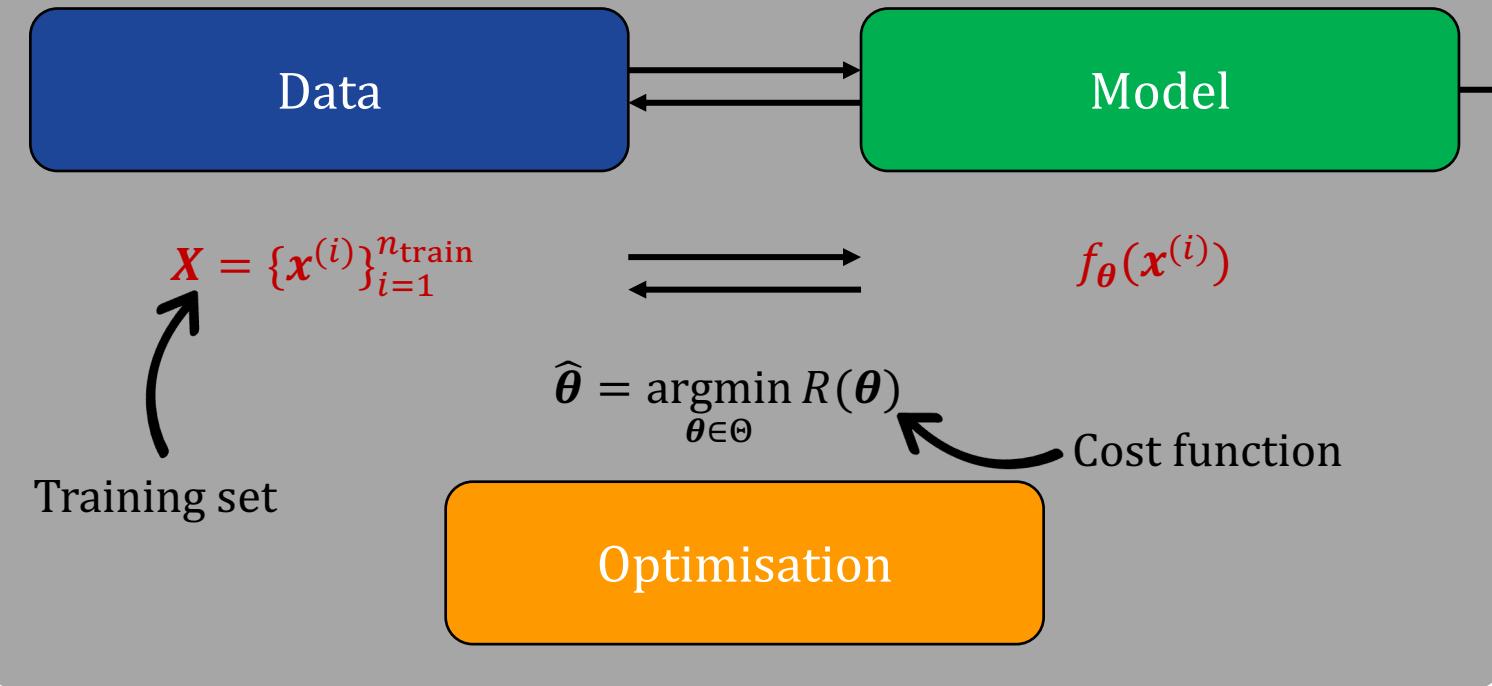
1 hidden layer of 100 units



2 hidden layers of 10 units



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**.

How to perform ERM in general?

Generalities

Introduction to ML

Basic supervised models

Optimisation of ML models

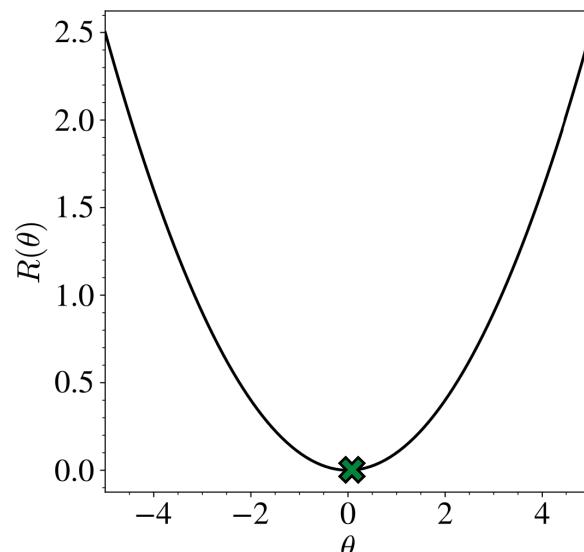
Generalization

- **Empirical risk minimisation** (ERM) provides a generic way to choose $\hat{\theta}$ in supervised learning

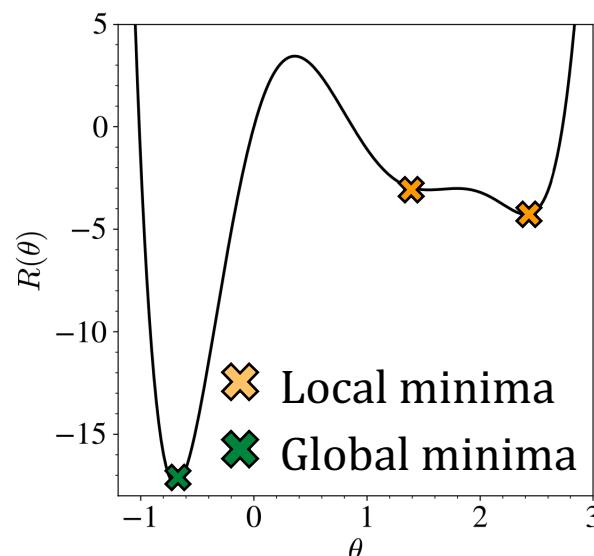
$$\hat{\theta} = \operatorname{argmin}_{\theta} R(\theta) = \operatorname{argmin}_{\theta} \frac{1}{n} \sum_{i=1}^n \ell(f_{\theta}(x^{(i)}), y^{(i)})$$

- The aim is to minimize a function of (many) parameters under a loss function encoding how well the model approximates the target variable
- In linear regression, we could directly optimise the empirical risk in closed-form, but **in general it is not possible**

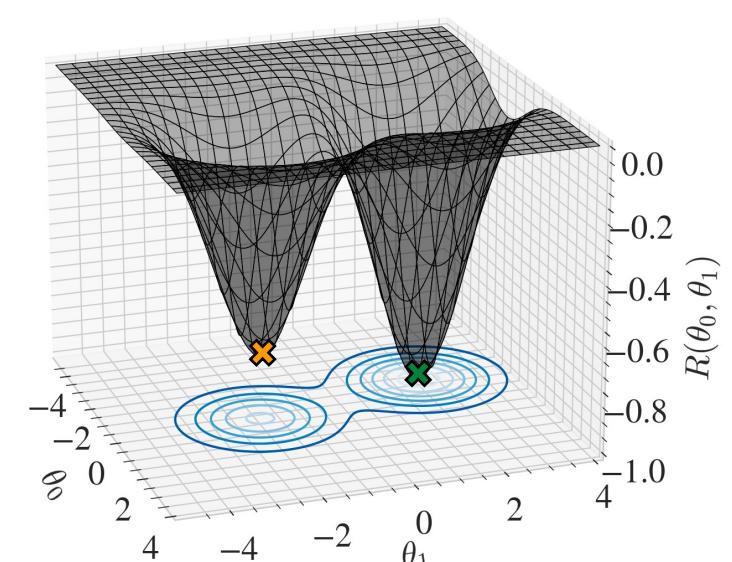
1D convex



1D non-convex



2D non-convex



Naïve view and curse of dimensionality

Generalities

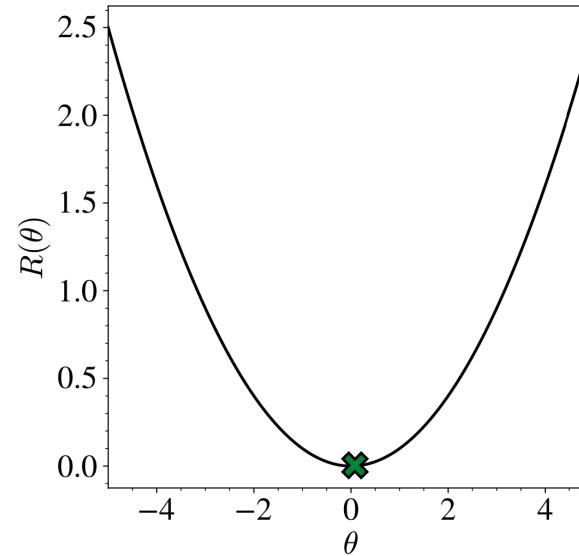
Introduction to ML

Basic supervised models

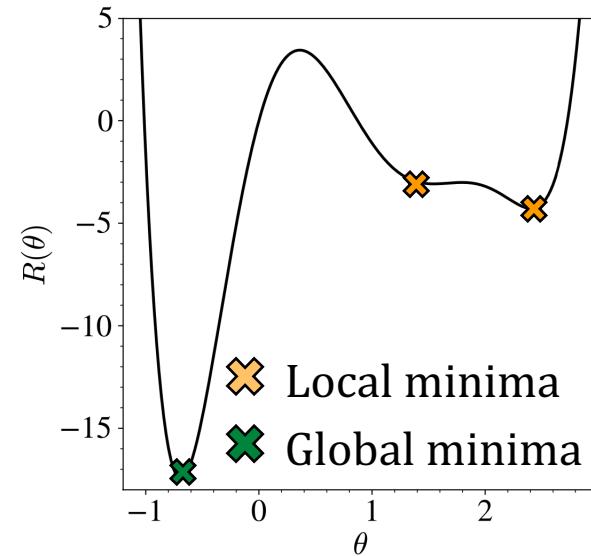
Optimisation of ML models

Generalization

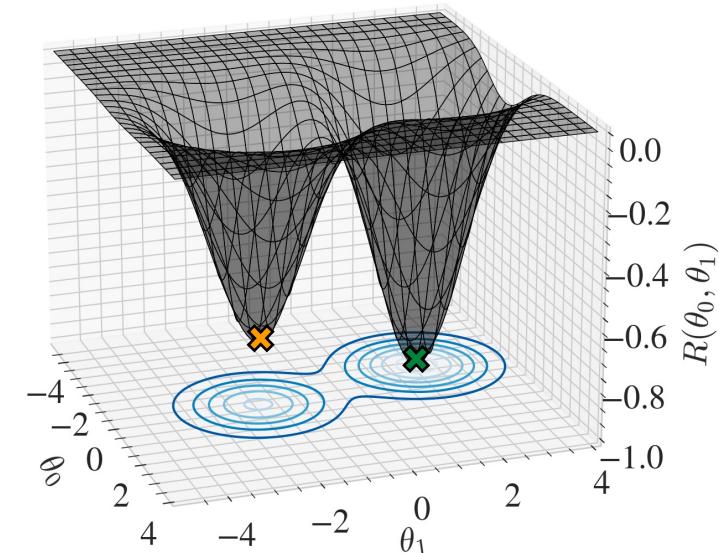
1D convex



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

Generalities

Introduction to ML

Basic supervised models

Optimisation of ML models

Generalization

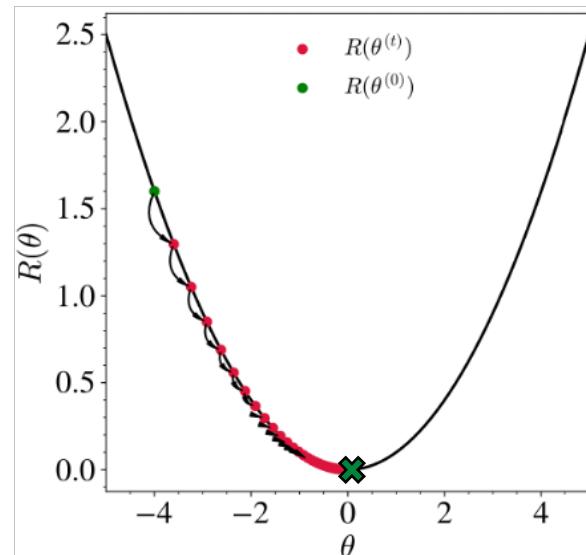
- A solution: local, directed search to navigate through the landscape
 - Numerical optimisation by **gradient descent**

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

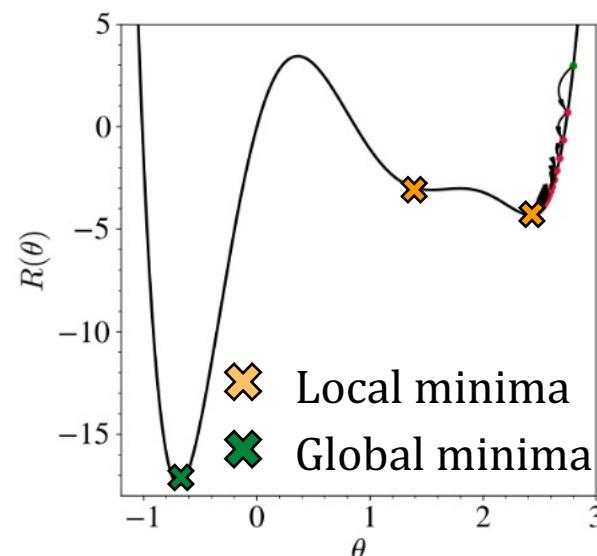
Algorithm: Gradient descent

- Initialise θ_0 randomly
- Compute $\theta^{(t+1)} = \theta^{(t)} - \eta \nabla_{\theta} R(\theta^{(t)})$
- Repeat 2 until $\|\theta^{(t+1)} - \theta^{(t)}\|_2^2 \leq \epsilon$

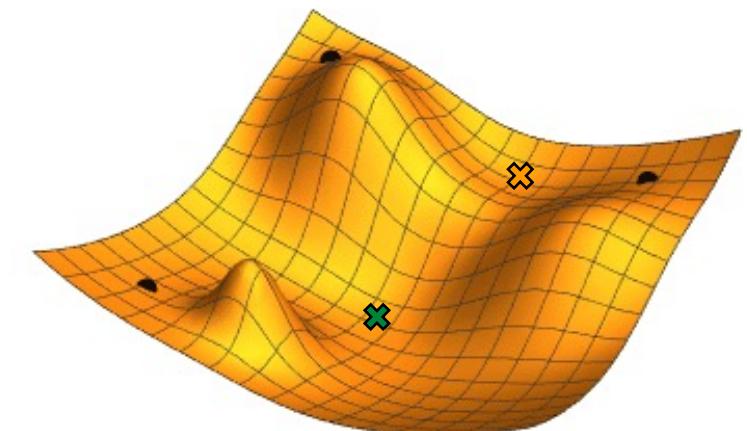
1D convex



1D non-convex



2D non-convex



A word about hyperparameters

43

Generalities

Introduction to ML

Basic supervised models

Optimisation of ML models

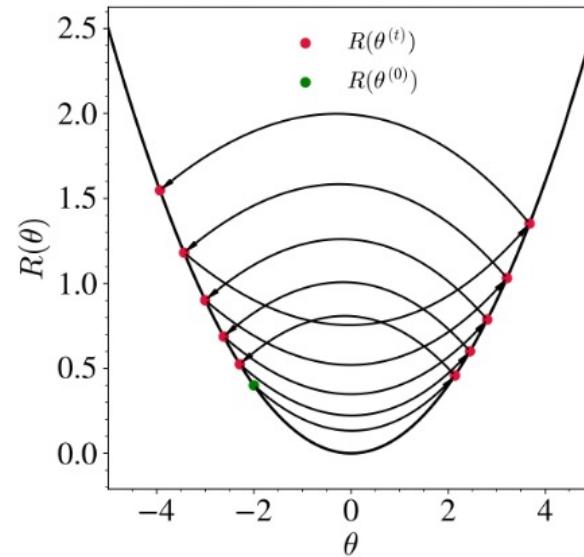
Generalization

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

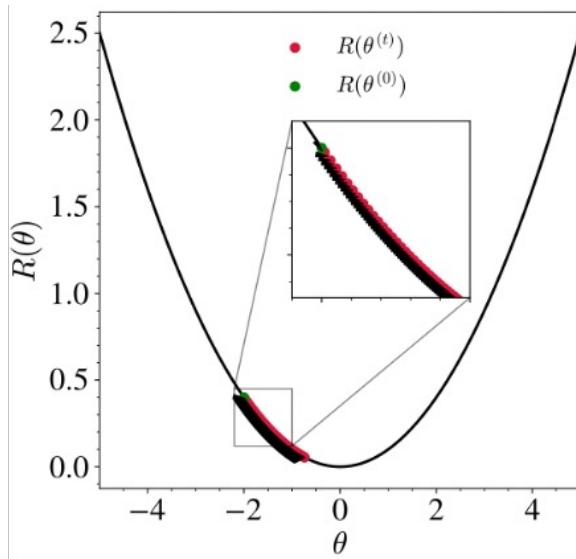


Hyperparameter: parameter that is **not learned** during the optimisation.
Ex: depth of tree in DTs, # of trees in RFs, learning rates, etc.

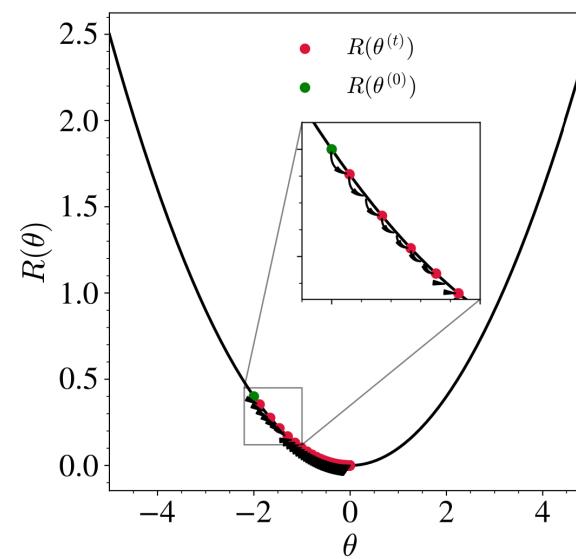
η too large



η too small



appropriate η



- 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

44

Generalities

Introduction to ML

Basic supervised models

Optimisation of ML models

Generalization

- 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 $\nabla_{\boldsymbol{\theta}} R(\boldsymbol{\theta})$

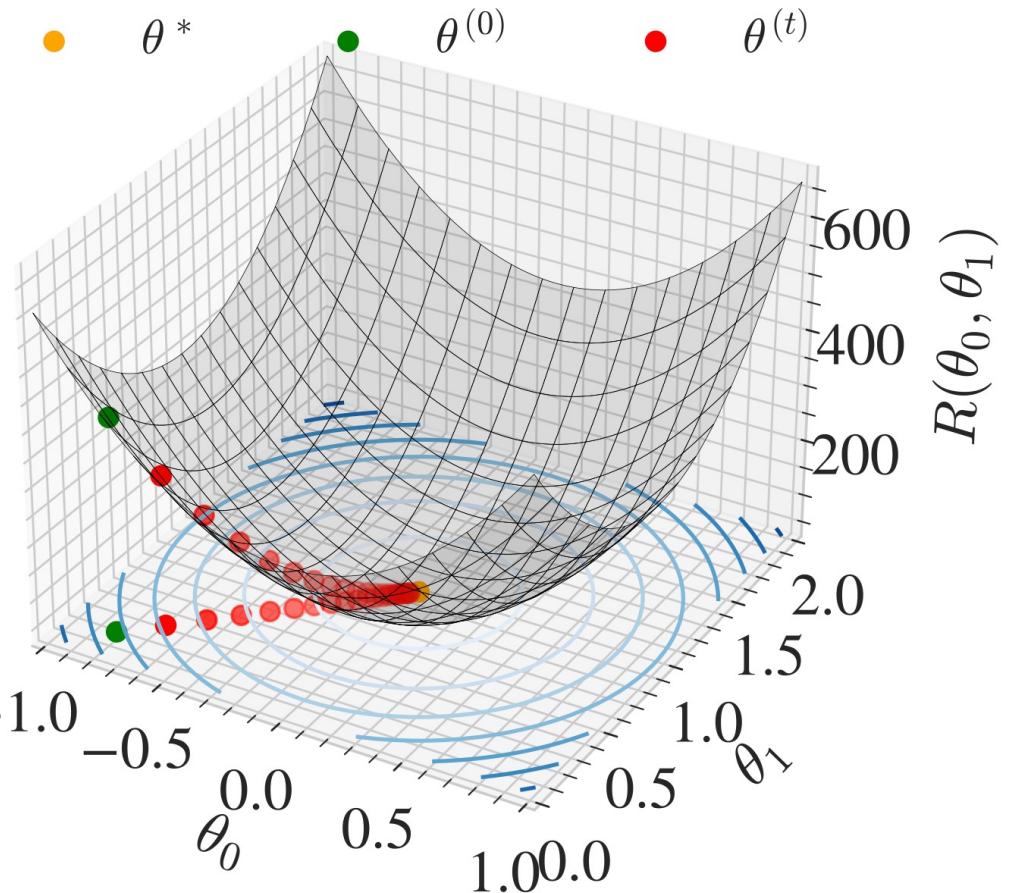
$$\frac{\partial R(\boldsymbol{\theta})}{\partial \theta_0} = \frac{2}{n} \sum_{i=1}^n (\theta_0 + \theta_1 x_1^{(i)} - y^{(i)}),$$

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

- 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.

Stochastic gradient descent algorithm

45

Generalities

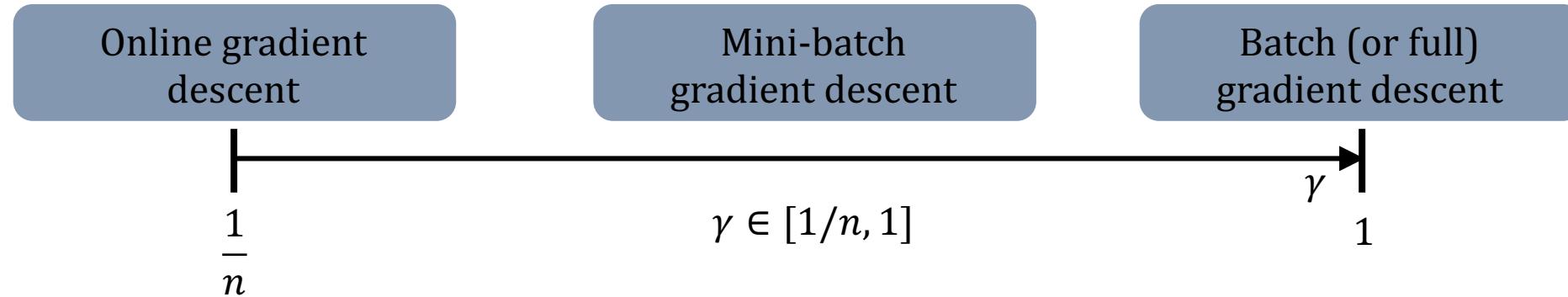
Introduction to ML

Basic supervised models

Optimisation of ML models

Generalization

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



Algorithm: Stochastic gradient descent

1. Initialise θ_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 γn 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

46

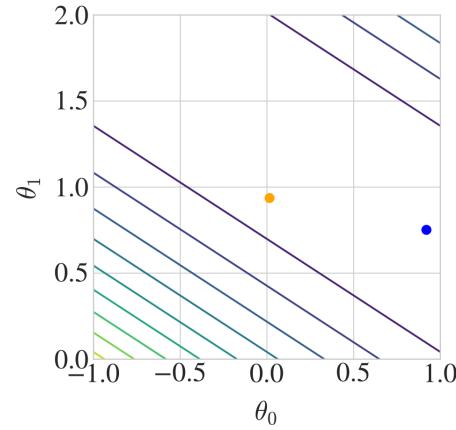
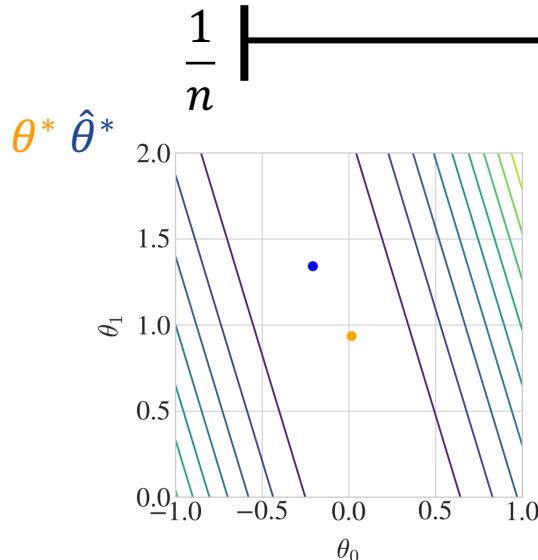
Generalities

Introduction to ML

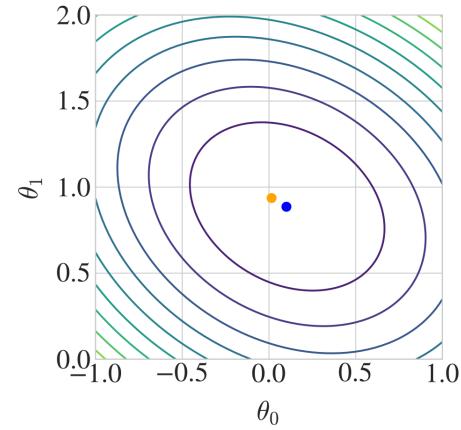
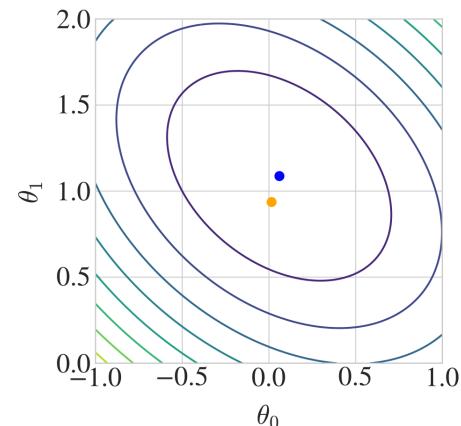
Basic supervised models

Optimisation of ML models

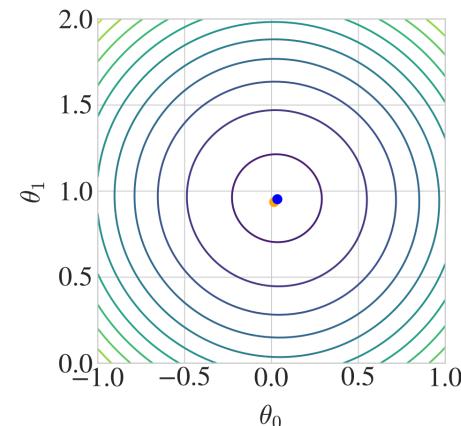
Generalization

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

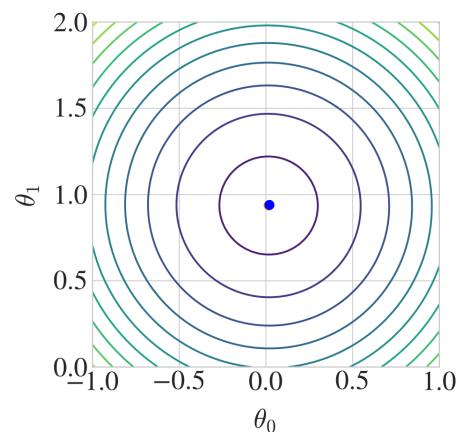
$$\gamma = 1/n$$



$$\gamma = 0.1$$



$$\gamma = 0.5$$



$$\gamma = 1$$

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

1

SGD on our linear regression

47

Generalities

Introduction to ML

Basic supervised models

Optimisation of ML models

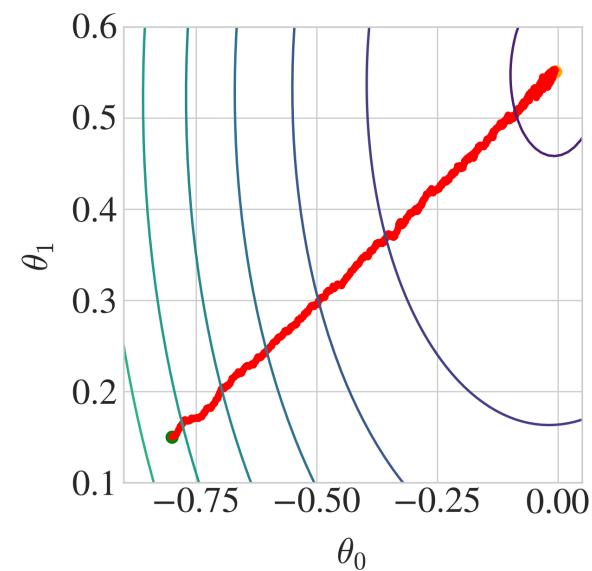
Generalization

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

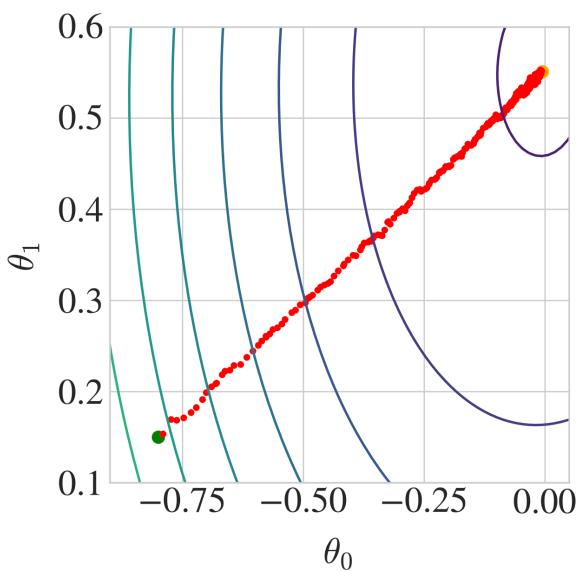
$$\frac{1}{n}$$

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

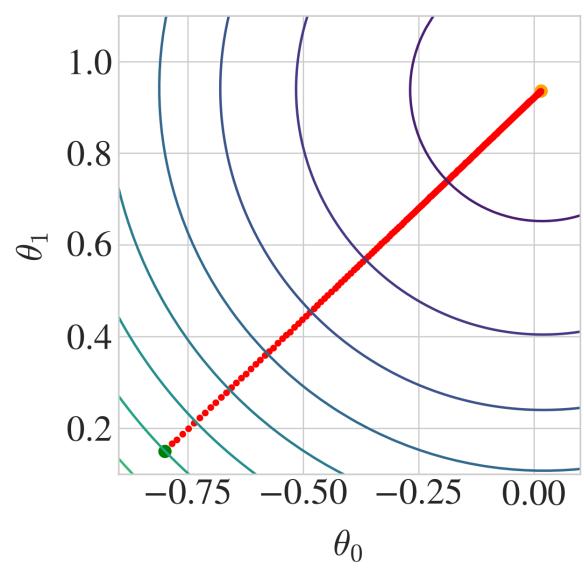
$$1$$



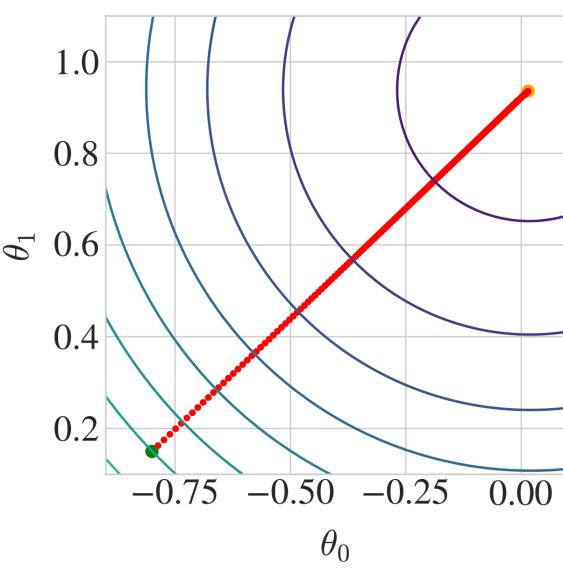
$$\gamma = 1/n$$



$$\gamma = 0.1$$



$$\gamma = 0.5$$



$$\gamma = 1$$

Computing gradients in NNs: backpropagation

48

Generalities

Introduction to ML

Basic supervised models

Optimisation of ML models

Generalization

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

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

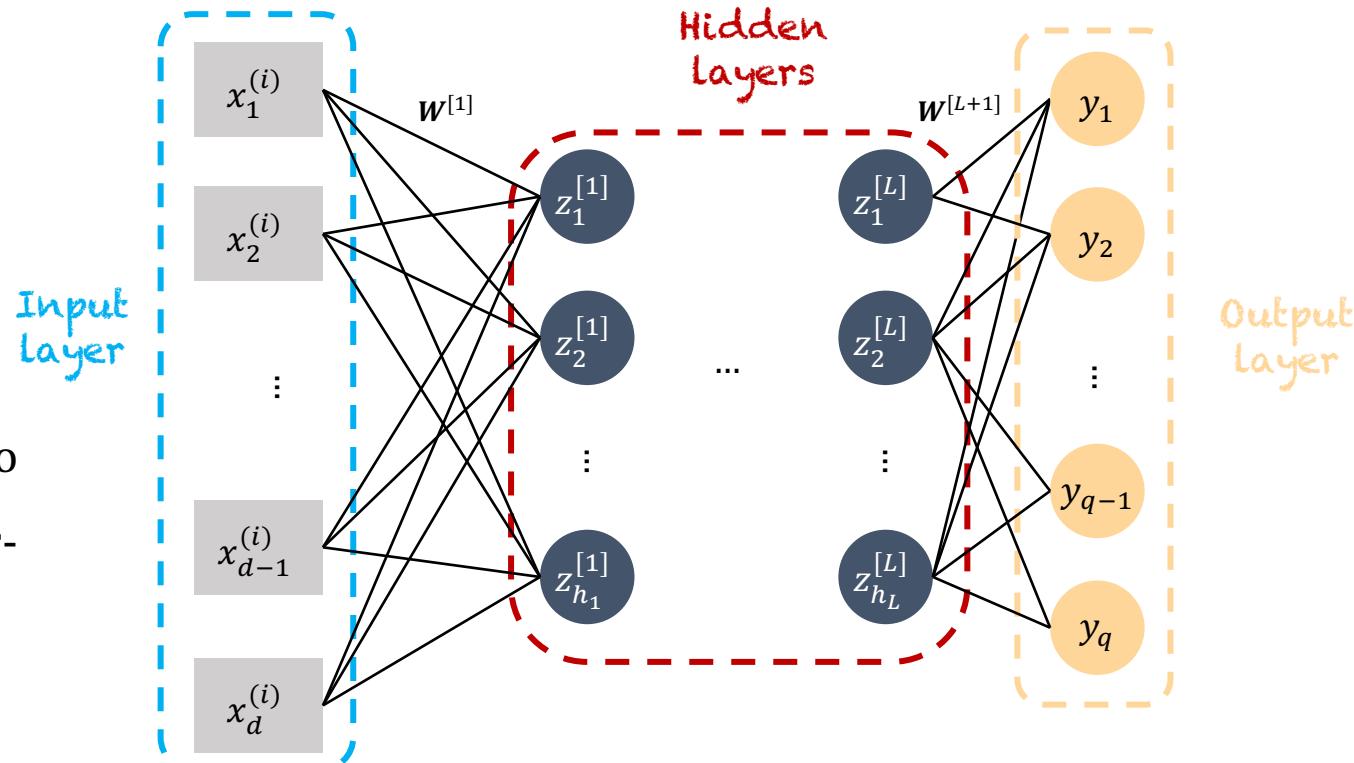
$$\mathbf{W}^{(l)} = [\mathbf{w}_1^{(l)}, \mathbf{w}_2^{(l)}, \dots, \mathbf{w}_{h_l}^{(l)}] \in \mathbb{R}^{h_l \times h_{l-1}}$$

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

- 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

49

Generalities

Introduction to ML

Basic supervised models

Optimisation of ML models

Generalization

- Take the example of a fully-connected network with $d = 3$, one hidden layer,

and two output neurons:

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

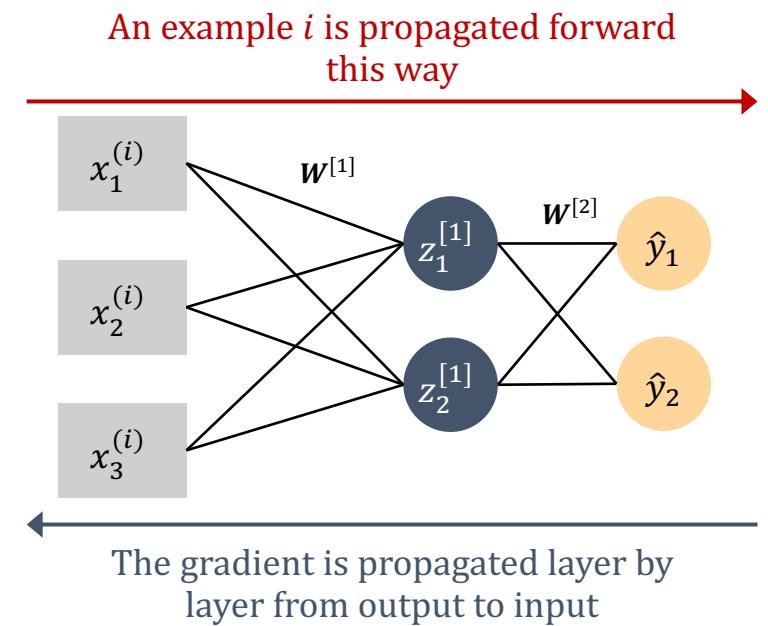
$$\mathbf{W}^{[1]} = \begin{bmatrix} b_1^{[1]} & w_{11}^{[1]} & w_{21}^{[1]} & w_{31}^{[1]} \\ b_2^{[1]} & w_{12}^{[1]} & w_{22}^{[1]} & w_{32}^{[1]} \end{bmatrix} \quad \mathbf{W}^{[2]} = \begin{bmatrix} b_1^{[2]} & w_{11}^{[2]} & w_{21}^{[2]} \\ b_2^{[2]} & w_{12}^{[2]} & w_{22}^{[2]} \end{bmatrix}$$

- 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

50

Generalities

Introduction to ML

Basic supervised models

Optimisation of ML models

Generalization

- 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

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

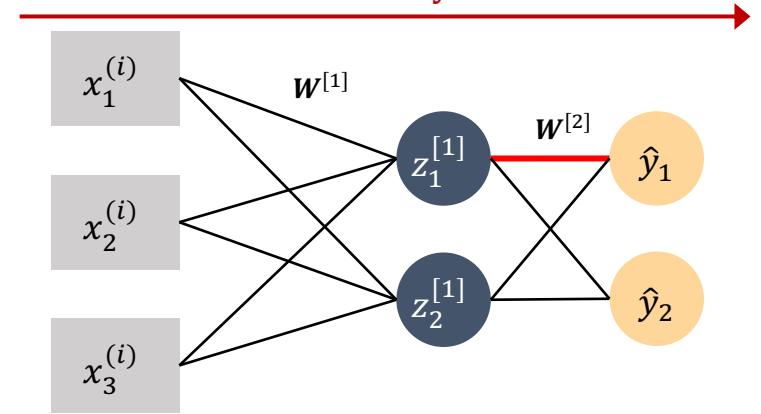
with $\hat{y}_1 = s(u_1^{[2]}) = s(w_{11}^{[2]} z_1^{[1]} + w_{21}^{[2]} z_2^{[1]} + b_1^{[2]})$

Finally,

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

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

An example i is propagated forward this way



The gradient is propagated layer by layer from output to input

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

$$\mathbf{W}^{[1]} = \begin{bmatrix} b_1^{[1]} & w_{11}^{[1]} & w_{12}^{[1]} & w_{31}^{[1]} \\ b_2^{[1]} & w_{12}^{[1]} & w_{22}^{[1]} & w_{32}^{[1]} \end{bmatrix}$$

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

Computing gradients in NNs: backpropagation

51

Generalities

Introduction to ML

Basic supervised models

Optimisation of ML models

Generalization

- 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 $w_{11}^{[1]}$:

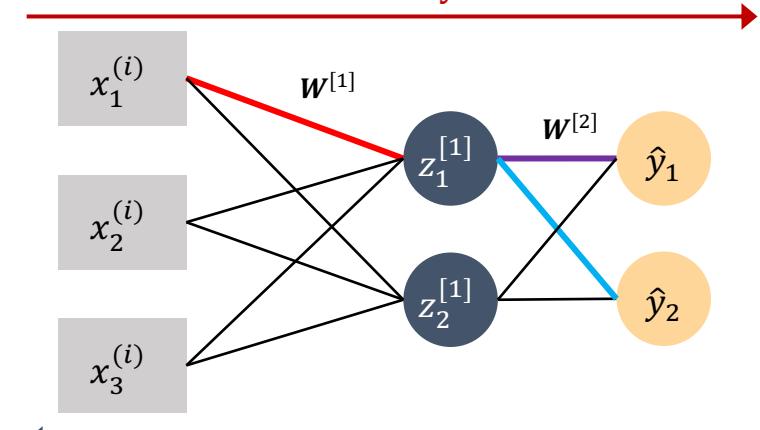
$$\delta_1^{[1]} = \frac{\partial \ell}{\partial u_1^{[2]}} \frac{\partial u_1^{[2]}}{\partial u_1^{[1]}} + \frac{\partial \ell}{\partial u_2^{[2]}} \frac{\partial u_2^{[2]}}{\partial u_1^{[1]}}$$



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

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

An example i is propagated forward this way



The gradient is propagated layer by layer from output to input

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

$$\mathbf{W}^{[1]} = \begin{bmatrix} b_1^{[1]} & w_{11}^{[1]} & w_{11}^{[1]} & w_{31}^{[1]} \\ b_2^{[1]} & w_{12}^{[1]} & w_{22}^{[1]} & w_{32}^{[1]} \end{bmatrix}$$

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

Computing gradients in NNs: backpropagation

52

Generalities

Introduction to ML

Basic supervised models

Optimisation of ML models

Generalization

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



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

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

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

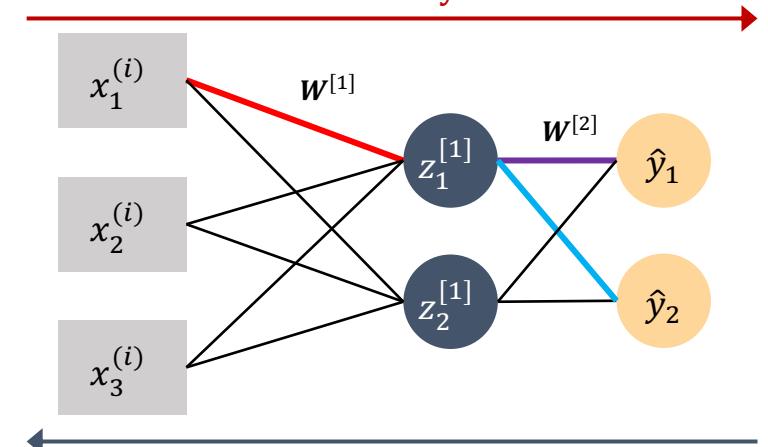
- Hence,

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

- Where everything is known!

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

An example i is propagated forward this way



The gradient is propagated layer by layer from output to input

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

$$\mathbf{w}^{[1]} = \begin{bmatrix} b_1^{[1]} & w_{11}^{[1]} & w_{11}^{[1]} & w_{31}^{[1]} \\ b_2^{[1]} & w_{12}^{[1]} & w_{22}^{[1]} & w_{32}^{[1]} \end{bmatrix}$$

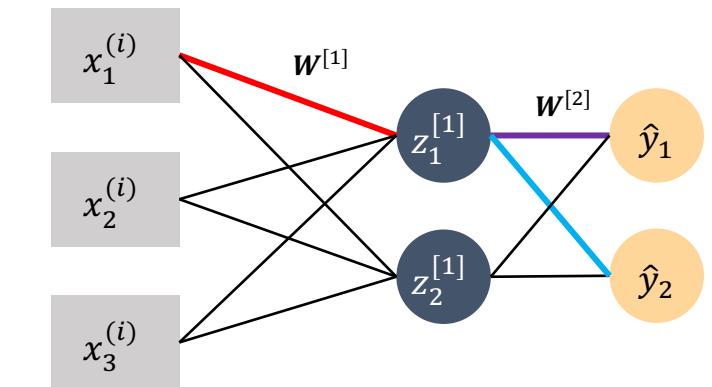
$$\mathbf{w}^{[2]} = \begin{bmatrix} b_1^{[2]} & w_{11}^{[2]} & w_{21}^{[2]} \\ b_2^{[2]} & w_{12}^{[2]} & w_{22}^{[2]} \end{bmatrix}$$

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



- 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



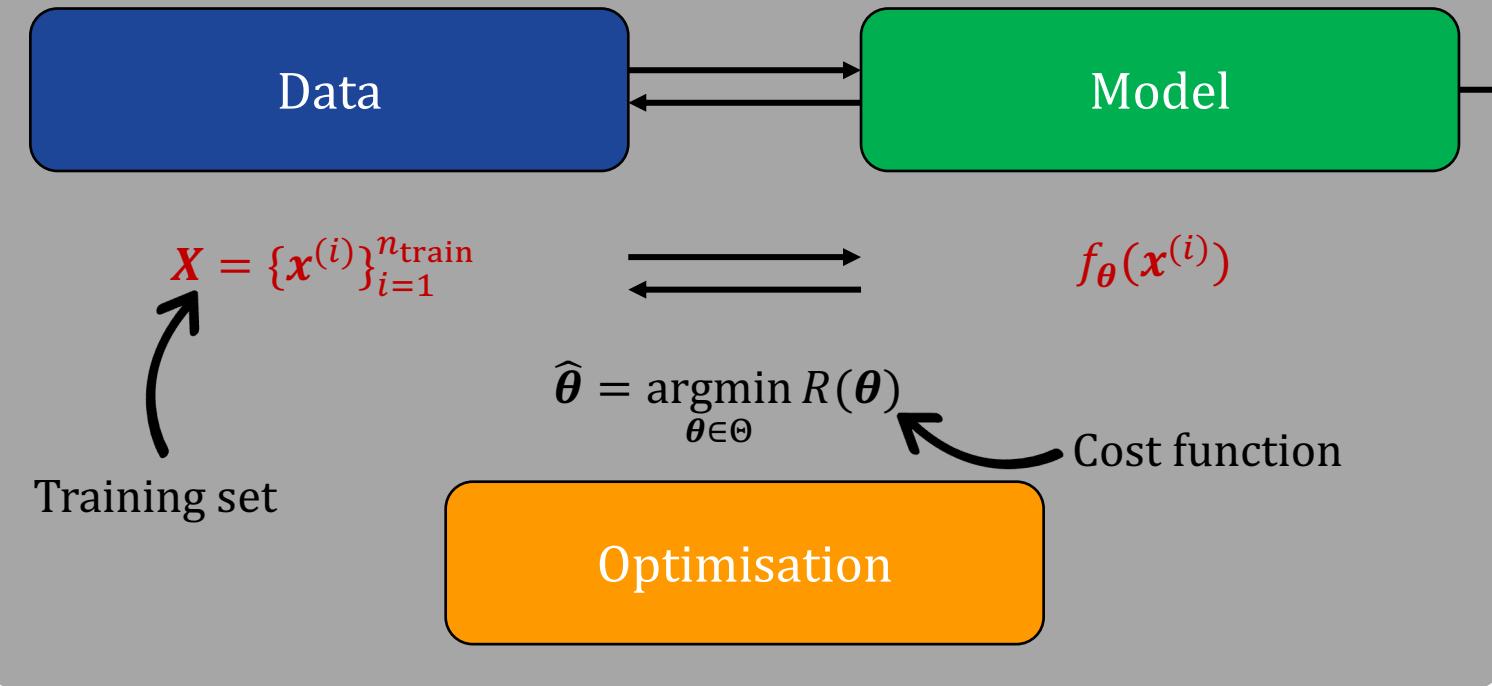
The gradient is propagated layer by layer from output to input

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

$$\mathbf{W}^{[1]} = \begin{bmatrix} b_1^{[1]} & w_{11}^{[1]} & w_{12}^{[1]} & w_{31}^{[1]} \\ b_2^{[1]} & w_{12}^{[1]} & w_{22}^{[1]} & w_{32}^{[1]} \end{bmatrix}$$

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

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**.

Example on decision trees

55

Generalities

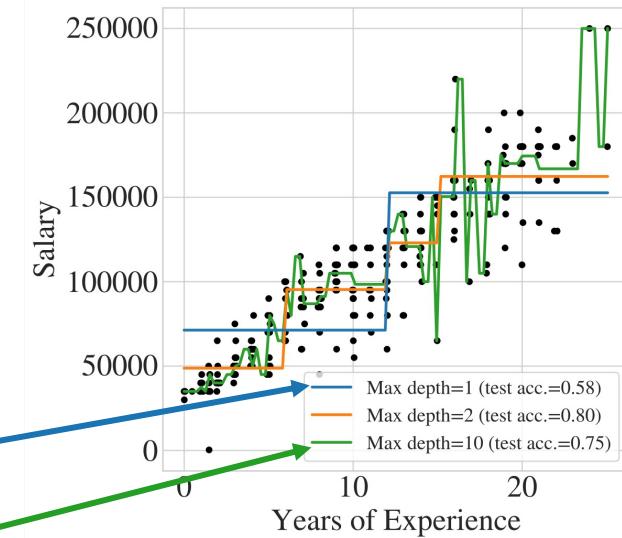
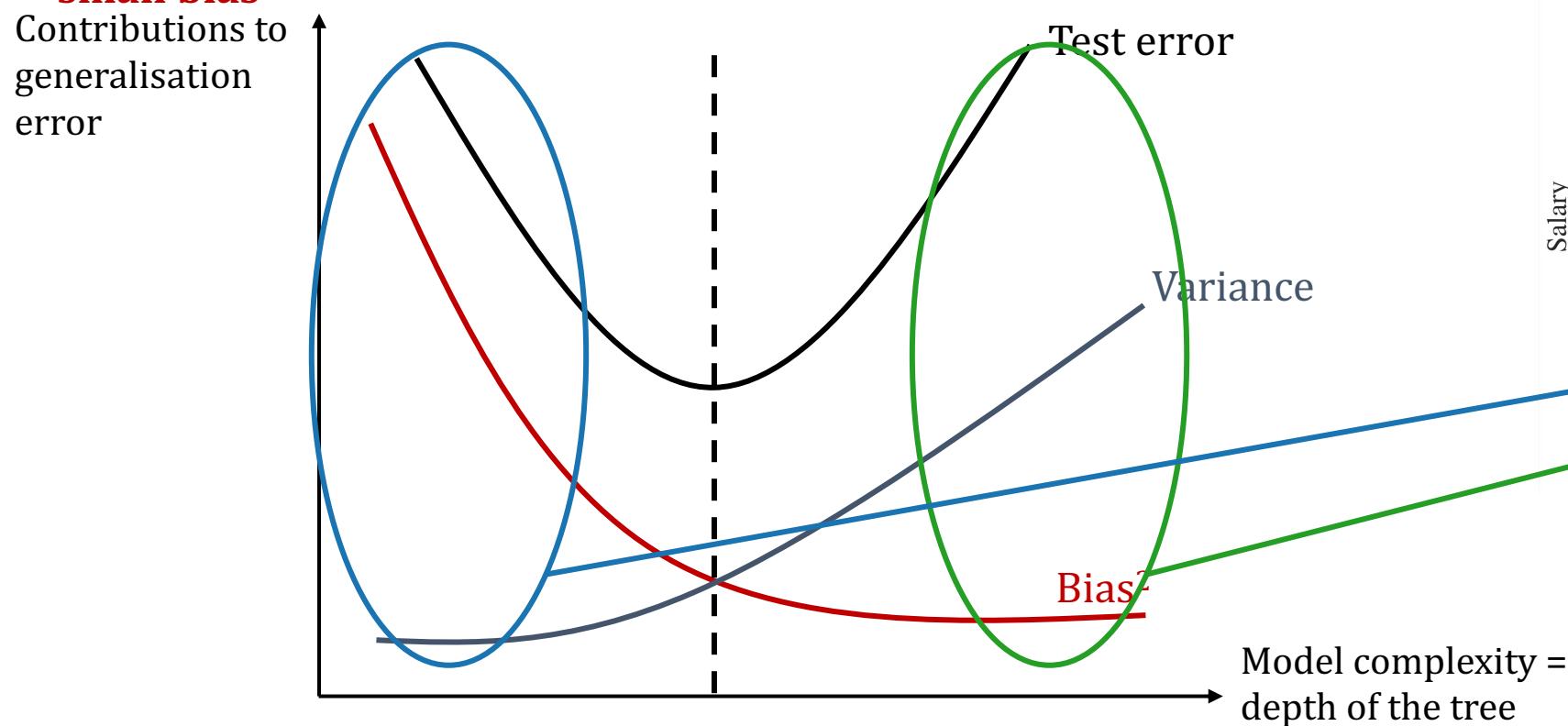
Introduction to ML

Basic supervised models

Optimisation of ML models

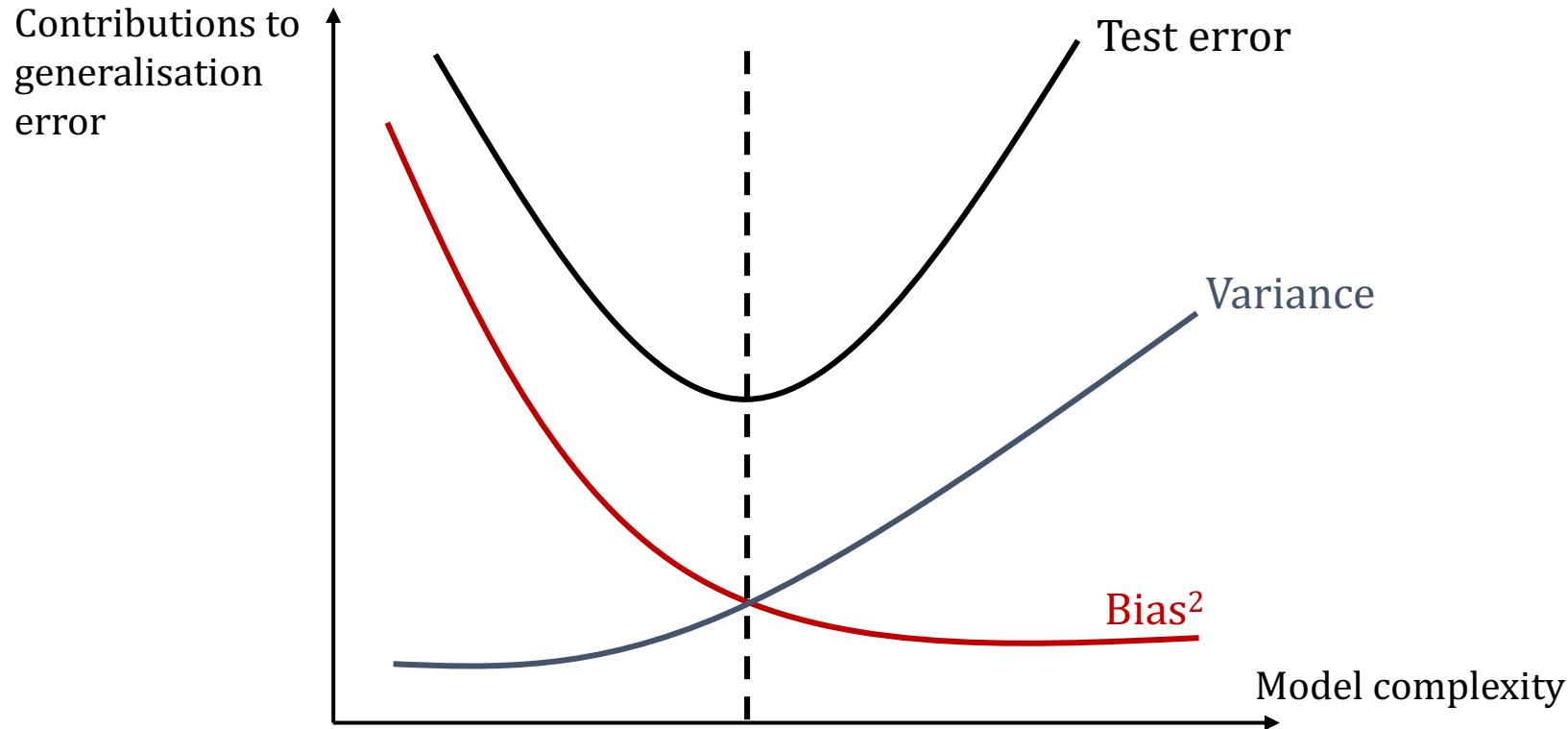
Generalization

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



Remember this?

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



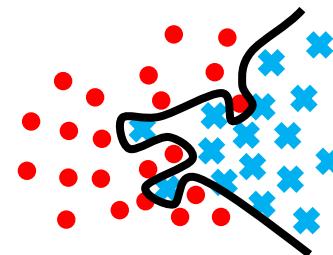
Models need to be built such that they are not too flexible to fit the noise in the data but also not too restrictive to avoid bias

* 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](#))

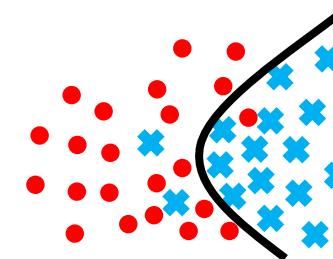
- To measure if we really learnt something useful in supervised learning in practice, we use a **test dataset** that the model has never seen but for which we know the labels and check that $R_{\text{train}}(\hat{\theta})$ and $R_{\text{test}}(\hat{\theta})$ are of the same order
- Based on the previous 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

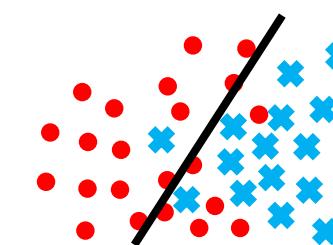
Overfitting



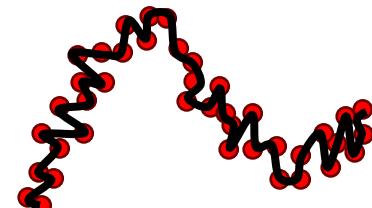
Good model



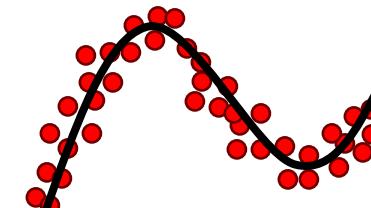
Underfitting



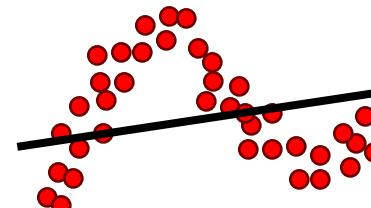
Regression



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



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



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

Some fixes to under and over fitting

58

Generalities

Introduction to ML

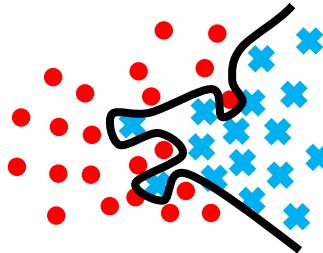
Basic supervised models

Optimisation of ML models

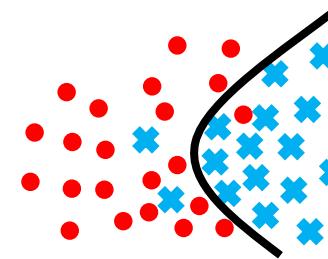
Generalization

Overfitting

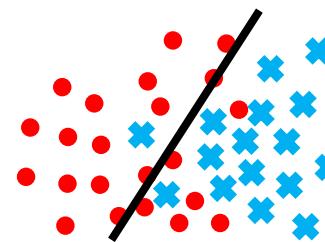
Classification



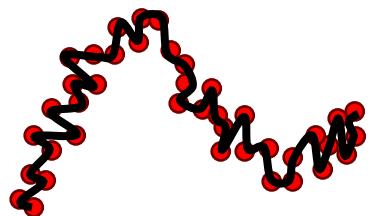
Good model



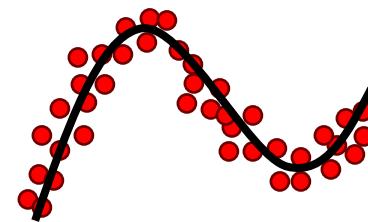
Underfitting



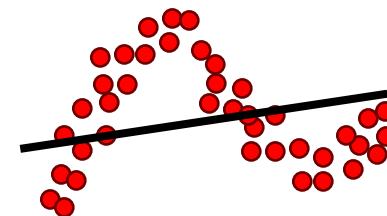
Regression



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



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



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

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

- 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

$$\hat{\theta} = \min_{\theta \in \Theta} R(\theta)$$

Unregularized
optimization

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

Regularized
optimization

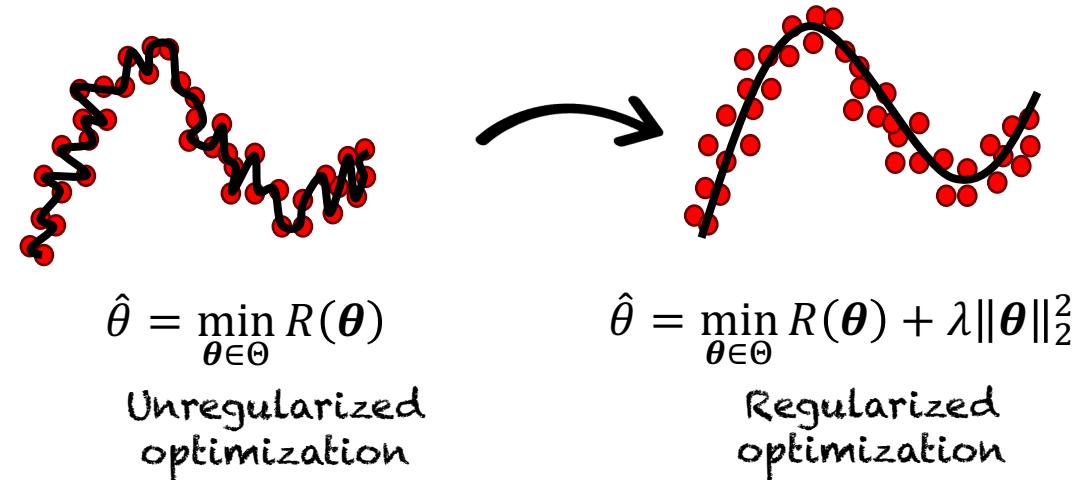
- The function $P(\theta)$ is called **penalty function** and the regularized problem can in fact be written equivalently as

$$\hat{\theta} = \min_{\theta \in \Theta} R(\theta) + \lambda P(\theta)$$

1. λ 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

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

- By solving instead the regularized optimization with an appropriate value of λ , we can reduce the overfitting



- **The model on the left lacks smoothness** which is introduced by the L_2 penalty
- The additional term here penalizes large weight values and **reduce the variance** of the estimator (even though increasing the bias)
- Regularization can also be viewed in terms of Bayesian statistics as **a prior distribution** you impose of the set of parameters