### Introduction to Machine Learning

Andrea De Lorenzo

A.Y.2020

< □ > < @ > < ≣ > < ≣ > ■ ● ○ Q @ 1/122

# Section 1

# General information

◆□ ▶ ◆ □ ▶ ◆ □ ▶ ◆ □ ▶ ● □ ● ○ ○ ○ 2/122

#### Lecturers



Dipartimento di Ingegneria e Architettura (DIA)

↓ □ ▶ ↓ □ ▶ ↓ ■ ▶ ↓ ■ ▶ ↓ ■ ♡ Q ○ 3/122

http://delorenzo.inginf.units.it/

### Course materials

#### Lecturer's slides

- http://delorenzo.inginf.units.it/project/ introduction-to-machine-learning-2020
- Suggested textbooks (for further reading)
  - Gareth James et al. An introduction to statistical learning. Vol. 6. Springer, 2013
- Other material:
  - I might point you to some scientific papers for discussing examples of application or specific details—just a "chat"

◆□▶ ◆□▶ ◆ □▶ ◆ □▶ ○ □ ○ ○ ○ ○ 4/122

Everything you are required to know is in the lecturer's slides

# Section 2

# Introduction

< □ ▶ < @ ▶ < ≣ ▶ < ≣ ▶ ■ ⑦ Q @ 5/122

# What is Machine Learning?

#### Definition

Machine Learning is the science of getting computer to learn without being explicitly programmed.

Definition Data Mining/Analytics is the science of discovering patterns in data.

### In practice

A set of mathematical and statistical tools for:

 building a model which allows to predict an output, given an input (*supervised learning*)

◆□ ▶ ◆□ ▶ ◆ ■ ▶ ◆ ■ ▶ ● ■ ⑦ Q ○ 7/122

example (input, output) pairs are available

 learn relationships and structures in data (unsupervised learning)

### Machine Learning: a computer science perspective



▲□▶ ▲圖▶ ▲ 필▶ ▲ 필▶ ■ ⑦ Q @ 8/122

#### Example problem: spam

#### Discriminate between spam and non-spam emails.

Google	in:spam		
Gmail •	•	C More ~	
COMPOSE			Delete all spam
		CSC Conference Secretari.	Call for Papers : 1st Annual Intern
Inbox (3) Starred Important Chats Sent Mail Drafts		Alexander Horn	Recently posted academic job vac
		Regalo di Benvenuto	emedvet@units.it per te uno Smar
		Peugeot Italia	Peugeot supervaluta il tuo usato. I
		CAP petite enfance	votre profil nous intéresse - Vous r
Spam (526)	□ ☆ >	Rachat de crédits	Réduisez vos mensualités jusqu'à
Categories Social Promotions (1) Updates (1) Purchases Travel Finance		Zalando	Le sneakers che conquistano la st
	□ ☆ >	Sondage National	Pour ou contre passer à 90 km/h s
		Oroscopo	Messaggio Privato per - Stai riceve
		Secret Escapes	Sconti Imbattibili su Hotel e Vacan
		Erogazione credito appro.	Fino a 50.000 euro, anche protesta

Figure: Spam filtering in Gmail.

#### Example problem: flight trajectories

Do flights over the same pair  $\langle origin, \, destination \rangle$  follow the "same" trajectory? Why?



#### Figure: Clustering of flight trajectories.

Example problem: image understanding Recognize objects in images.



Figure: Object recognition in Google Photos.

 $\ensuremath{\mathsf{Q}}\xspace$  : what type of learning (supervised/unsupervised) is in the examples?

<□ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □

- spam
- image understanding
- flight trajectories

## Why ML/DM "today"?

- we collect more and more data (big data)
- we have more and more computational power



Figure: From http://www.mkomo.com/cost-per-gigabyte-update.

# ML/DM is popular!



Figure: Popular areas of interest, from the Skill Up 2016: Developer Skills Report<sup>2</sup>

<sup>1</sup>https://techcus.com/p/r1zSmbXut/ top-5-highest-paying-programming-languages-of-2016/. <sup>2</sup>https://techcus.com/p/r1zSmbXut/ top-5-highest-paying-programming-languages-of-2016/.

Be able to:

- 1. design
- 2. implement
- 3. assess experimentally

an end-to-end Machine Learning or Data Mining system.

<□ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □

Be able to:

- 1. design
- 2. implement
- 3. assess experimentally

an end-to-end Machine Learning or Data Mining system.

Which is the problem to be solved? Which are the input and output? Which are the most suitable techniques? How should data be prepared? Does computation time matter?

◆□ → ◆□ → ◆ ■ → ▲ ● → ■ の Q ○ 15/122

Be able to:

- 1. design
- 2. implement
- 3. assess experimentally

an end-to-end Machine Learning or Data Mining system.

Which is the problem to be solved? Which are the input and output? Which are the most suitable techniques? How should data be prepared? Does computation time matter?

◆□ → ◆□ → ◆ ■ → ▲ ● → ■ の Q ○ 15/122

Write some code!

Be able to:

- 1. design
- 2. implement
- 3. assess experimentally

an end-to-end Machine Learning or Data Mining system.

- Which is the problem to be solved? Which are the input and output? Which are the most suitable techniques? How should data be prepared? Does computation time matter?
- Write some code!
- How to measure solution quality? How to compare solutions? Is my solution general?

Be able to:

- 1. design
- 2. implement
- 3. assess experimentally

an end-to-end Machine Learning or Data Mining system.

- Which is the problem to be solved? Which are the input and output? Which are the most suitable techniques? How should data be prepared? Does computation time matter?
- Write some code!
- How to measure solution quality? How to compare solutions? Is my solution general?
  - Itself: design and implementation

# Aims of the course: communication

Be able to:

- 1. design
- 2. implement
- 3. assess experimentally

an end-to-end Machine Learning or Data Mining system. And be able to convince the "client" that it is:

◆□ ▶ ◆□ ▶ ◆ ■ ▶ ◆ ■ ▶ ● ■ の Q @ 16/122

- technically sound
- economically viable
- in its larger context

#### Subsection 1

Motivating example

< □ > < ■ > < ≧ > < ≧ > ≧ の < ? 17/122

### The amateur botanist friend

He likes to collect Iris plants. He "realized" that there are 3 species, in particular, that he likes: *Iris setosa, Iris virginica,* and *Iris versicolor*. He'd like to have a tool to automatically *classify* collected samples in one of the 3 species.



Figure: Iris versicolor.

How to help him?

Which is the problem to be solved?

#### Which is the problem to be solved?

Assign exactly one specie to a sample.

Which is the problem to be solved?
Assign exactly one specie to a sample.
Which are the input and output?

Which is the problem to be solved?

Assign exactly one specie to a sample.

Which are the input and output?

Output: one species among I. setosa, I. virginica, I. versicolor.

Which is the problem to be solved?

- Assign exactly one specie to a sample.
- Which are the input and output?
  - Output: one species among I. setosa, I. virginica, I. versicolor.
  - Input: the plant sample...

- Which is the problem to be solved?
  - Assign exactly one specie to a sample.
- Which are the input and output?
  - Output: one species among I. setosa, I. virginica, I. versicolor.

- Input: the plant sample...
  - a description in natural language?

- Which is the problem to be solved?
  - Assign exactly one specie to a sample.
- Which are the input and output?
  - Output: one species among I. setosa, I. virginica, I. versicolor.

- Input: the plant sample...
  - a description in natural language?
  - a digital photo?

- Which is the problem to be solved?
  - Assign exactly one specie to a sample.
- Which are the input and output?
  - Output: one species among I. setosa, I. virginica, I. versicolor.

- Input: the plant sample...
  - a description in natural language?
  - a digital photo?
  - DNA sequences?

- Which is the problem to be solved?
  - Assign exactly one specie to a sample.
- Which are the input and output?
  - Output: one species among I. setosa, I. virginica, I. versicolor.

<□ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □

- Input: the plant sample...
  - a description in natural language?
  - a digital photo?
  - DNA sequences?
  - some measurements of the sample!

## Iris: input and output



Figure: Sepal and petal.

Input: sepal length and width, petal length and width (in cm) Output: the class Example:  $(5.1, 3.5, 1.4, 0.2) \rightarrow I$ . setosa

## Other information

The botanist friend asked a senior botanist to inspect several samples and label them with the corresponding species.

Sepal length	Sepal width	Petal length	Petal width	Species
5.1	3.5	1.4	0.2	I. setosa
4.9	3.0	1.4	0.2	I. setosa
7.0	3.2	4.7	1.4	I. versicolor
6.0	2.2	5.0	1.5	I. virginica

◆□ ▶ ◆□ ▶ ◆ ■ ▶ ◆ ■ ▶ ● ■ の Q @ 21/122

# Notation and terminology

- Sepal length, sepal width, petal length, and petal width are input variables (or independent variables, or features, or attributes).
- Species is the output variable (or dependent variable, or response).

### Notation and terminology

$$\mathbf{X} = \begin{pmatrix} x_{1,1} & x_{1,2} & \cdots & x_{1,p} \\ x_{2,1} & x_{2,2} & \cdots & x_{2,p} \\ \vdots & \vdots & \ddots & \vdots \\ x_{n,1} & x_{n,2} & \cdots & x_{n,p} \end{pmatrix} \quad \mathbf{y} = \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{pmatrix}$$

►  $x_1^T = (x_{1,1}, x_{1,2}, ..., x_{1,p})$  is an observation (or instance, or data point), composed of *p* variable values;

### Notation and terminology

$$\mathbf{X} = \begin{pmatrix} x_{1,1} & x_{1,2} & \cdots & x_{1,p} \\ x_{2,1} & x_{2,2} & \cdots & x_{2,p} \\ \vdots & \vdots & \ddots & \vdots \\ x_{n,1} & x_{n,2} & \cdots & x_{n,p} \end{pmatrix} \quad \mathbf{y} = \begin{pmatrix} \mathbf{y}_1 \\ \mathbf{y}_2 \\ \vdots \\ \mathbf{y}_n \end{pmatrix}$$

x<sub>1</sub><sup>T</sup> = (x<sub>1,1</sub>, x<sub>1,2</sub>, ..., x<sub>1,p</sub>) is an observation (or instance, or data point), composed of p variable values; y<sub>1</sub> is the corresponding output variable value
#### Notation and terminology

$$\mathbf{X} = \begin{pmatrix} x_{1,1} & x_{1,2} & \cdots & x_{1,p} \\ x_{2,1} & \mathbf{x}_{2,2} & \cdots & \mathbf{x}_{2,p} \\ \vdots & \vdots & \ddots & \vdots \\ x_{n,1} & \mathbf{x}_{n,2} & \cdots & \mathbf{x}_{n,p} \end{pmatrix} \quad \mathbf{y} = \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{pmatrix}$$

- ► x<sub>1</sub><sup>T</sup> = (x<sub>1,1</sub>, x<sub>1,2</sub>, ..., x<sub>1,p</sub>) is an observation (or instance, or data point), composed of p variable values; y<sub>1</sub> is the corresponding output variable value
- ▶  $\mathbf{x}_2^T = (x_{1,2}, x_{2,2}, \dots, x_{n,2})$  is the vector of all the *n* values for the 2nd variable (X<sub>2</sub>).

## Notation and terminology

Different communities (e.g., statistical learning vs. machine learning vs. artificial intelligence) use different terms and notation:

<□ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □

• 
$$x_i^{(i)}$$
 instead of  $x_{i,j}$  (hence  $x^{(i)}$  instead of  $x_i$ )

m instead of n and n instead of p

Focus on the meaning!

**>** . . .



Simplification: forget petal and I. virginica  $\rightarrow$  2 variables, 2 species (binary classification problem).

 Problem: given any new observation, we want to automatically assign the species.



- Problem: given any new observation, we want to automatically assign the species.
- Sketch of a possible solution:



- Problem: given any new observation, we want to automatically assign the species.
- Sketch of a possible solution:
  - 1. learn a model (classifier)



- Problem: given any new observation, we want to automatically assign the species.
- Sketch of a possible solution:
  - 1. learn a model (classifier)
  - 2. "use" model on new observations



#### "A" model?

There could be many possible models:

- how to choose?
- how to compare?
- **Q:** a model of what?

The choice of the model/tool/technique to be used is determined by many factors:

<□ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □

- Problem size (n and p)
- Availability of an output variable (y)
- Computational effort (when learning or "using")
- Explicability of the model

We will see some options.

...

#### Experimentally: does the model work well on (new) data?

Experimentally: does the model work well on (new) data? Define "works well":

◆□ → ◆□ → ◆ ■ → ▲ ■ → ● ○ 28/122

a single performance index?

how to measure?

repeatability/reproducibility...

Q: what's the difference?

We will see/discuss some options.

It does not work well...

Why?

- the data is not informative
- the data is not representative
- the data has changed
- the data is too noisy

We will see/discuss these issues.

◆□ ▶ ◆□ ▶ ◆ ■ ▶ ◆ ■ ▶ ● ■ の Q @ 29/122

### ML is not magic

*Problem*: find birth town from height/weight.



When "solving" a problem, we usually need:

- explore/visualize data
- apply one or more ML technique
- assess learned models

"By hands?" No, with software!

## ML/DM software

Many options:

- libraries for general purpose languages:
  - Java: e.g., http://haifengl.github.io/smile/
  - Python: e.g., http://scikit-learn.org/stable/

specialized sw environments:

- Octave: https://en.wikipedia.org/wiki/GNU\_Octave
- R: https:

//en.wikipedia.org/wiki/R\_(programming\_language)

from scratch

▶ ...

## ML/DM software: which one?

- production/prototype
- platform constraints
- degree of (data) customization
- documentation availability/community size

- ► ...
- previous knowledge/skills

## ML/DM software: why?

In all cases, sw allows to be more productive and concise. E.g., learn and use a model for classification, in Java+Smile:

## Section 3

#### Plotting data: an overview

<□ ▶ < @ ▶ < ≧ ▶ < ≧ ▶ ≧ り Q <sup>®</sup> 35/122

#### Advanced plotting

- many packages (e.g., ggplot2)
- many options

Which is the most proper chart to support a thesis?





うくで 37/122

bringing changing keep delayed generally maximum really formes building, extent currency requirements R p figures building boost almost La construction de la construcción de la construcci 5 minur deurse zu saka deurse seine saka deurse seine sei § indextan upt m state learn "Walketyndestolide defamiliet offstelling without warm i concerning many analysis period of financial titrid offstelling of the state of the s asenough 85 often path product of a series of the series o account recently chief plans Ppart offer union american monday pressure argest continue a back press outside difficult cost callexpected detailslooking west january predicted without joint prediction upon the second sec 2017 torus/Biblioteness pockesman added "We march since ... & average \_btaining occurs house board and the constraints however current major house board and the constraints however current major house board and the constraint and the cons introduced volume want analyst companies fainge past e buying slow makes proposal days start gave longer washington expansion whether problems increases imports expansion whether start and the start start committerimmediate local levels support december subsidiary tomorrowactions sharp commarce needs agriculture paying maaaurea legve previously balaved exchanges reduction manager requirement competitiveness transactions chance factors > coming meanable

gear am 0 drat mpg 10.4 vs 000 0.0 qsec 14.5 W 1.51 472 disp 71.1 су 335 hp carb

Car Milage Data in PC2/PC1 Order



<□ > < □ > < □ > < Ξ > < Ξ > Ξ の Q ↔ 40/122

# Section 4

#### Tree-based methods

<□ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □

#### The carousel robot attendant

*Problem*: replace the carousel attendant with a robot which automatically decides who can ride the carousel.



Observed human attendant's decisions.



How can the robot take the decision?

< □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □



How can the robot take

▲□▶ ▲□▶ ▲三▶ ★三▶ 三三 のへで

43/122

 $\blacktriangleright$  if younger than 10  $\rightarrow$ 



< □ > < @ > < 注 > < 注 > 注 の Q @ 43/122

Observed human attendant's decisions.



How can the robot take the decision?

- if younger than  $10 \rightarrow can't!$
- otherwise:

◆□▶ ◆□▶ ◆ □▶ ◆ □▶ ○ □ ○ ○ ○ ○

43/122

- if shorter than 120  $\rightarrow$  can't!
- otherwise  $\rightarrow$  can!

Observed human attendant's decisions.



How can the robot take the decision?

- if younger than  $10 \rightarrow can't!$
- otherwise:
  - if shorter than 120  $\rightarrow$  can't!
  - otherwise  $\rightarrow$  can!





#### How to build a decision tree

Dividi-et-impera (recursively):

- find a cut variable and a cut value
- for left-branch, dividi-et-impera
- for right-branch, dividi-et-impera

How to build a decision tree: detail

Recursive binary splitting function BUILDDECISIONTREE(X, y) if SHOULDSTOP(y) then  $\hat{y} \leftarrow \text{most common class in } \mathbf{y}$ **return** new terminal node with  $\hat{y}$ else  $(i, t) \leftarrow \text{BestBranch}(\mathbf{X}, \mathbf{y})$  $n \leftarrow$  new branch node with (i, t)append child BUILDDECISIONTREE( $\mathbf{X}|_{\mathbf{x}_i < t}, \mathbf{y}|_{\mathbf{x}_i < t}$ ) to n append child BUILDDECISIONTREE( $\mathbf{X}|_{\mathbf{x}_i > t}, \mathbf{y}|_{\mathbf{x}_i > t}$ ) to n return n end if end function

- Recursive binary splitting
- Top down (start from the "big" problem)

#### Best branch

$$\begin{array}{l} \text{function BestBranch}(\mathbf{X}, \mathbf{y}) \\ (i^{\star}, t^{\star}) \leftarrow \arg\min_{i,t} E(\mathbf{y}|_{\mathbf{x}_i \geq t}) + E(\mathbf{y}|_{\mathbf{x}_i < t}) \\ \text{return } (i^{\star}, t^{\star}) \\ \text{end function} \end{array}$$

Classification error on subset:

$$egin{aligned} \mathsf{E}(\mathbf{y}) &= rac{|\{y \in \mathbf{y} : y 
eq \hat{y}\}|}{|\mathbf{y}|} \ \hat{y} &= ext{the most common class in } \mathbf{y} \end{aligned}$$

Greedy (choose split to minimize error now, not in later steps)

#### Best branch

$$(i^{\star}, t^{\star}) \leftarrow \operatorname*{arg\,min}_{i,t} E(\mathbf{y}|_{\mathbf{x}_i \geq t}) + E(\mathbf{y}|_{\mathbf{x}_i < t})$$

◆□ ▶ ◆□ ▶ ◆ ■ ▶ ◆ ■ ▶ ● ■ の Q @ 47/122

The formula say what is done, not how is done!

**Q:** "how" can different methods differ?

## Stopping criterion

function SHOULDSTOP(y) if y contains only one class then return true else if  $|y| < k_{min}$  then return true else return false end if end function

Other possible criterion:

tree depth larger than d<sub>max</sub>

#### Best branch criteria

Classification error E() works, but has been shown to be "not sufficiently sensitive for tree-growing".

$$E(\mathbf{y}) = \frac{|\{y \in \mathbf{y} : y \neq \hat{y}\}|}{|\mathbf{y}|} = 1 - \max_{c} \frac{|\{y \in \mathbf{y} : y = c\}|}{|\mathbf{y}|} = 1 - \max_{c} p_{\mathbf{y},c}$$

Other two option:

► Gini index

$$G(\mathbf{y}) = \sum_{c} p_{\mathbf{y},c}(1-p_{\mathbf{y},c})$$

Cross-entropy

$$D(\mathbf{y}) = -\sum_{c} p_{\mathbf{y},c} \log p_{\mathbf{y},c}$$

For all indexes, the lower the better (node impurity).
Best branch criteria: binary classification



◆□ → ◆□ → ◆ ■ → ▲ ■ → ● ○ ○ ○ 50/122

Cross-entropy is rescaled.

Q: what happens with multiclass problems?

### Categorical independent variables

- Trees can work with categorical variables
- ▶ Branch node is  $x_i = c$  or  $x_i \in C' \subset C$  (*c* is a class)
- Can mix categorical and numeric variables

# Stopping criterion: role of $k_{\min}$



When the tree is "too complex"

- less readable/understandable/explicable
- maybe there was noise into the data
- Q: what's noise in carousel data?

Tree complexity is not related (only) with  $k_{\min}$ , but also with data

Tree complexity: other interpretation

maybe there was noise into the data

The tree *fits* the learning data too much:

- it overfits (overfitting)
- does not generalize (high variance: model varies if learning data varies)

◆□ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ <

"model varies if learning data varies": what? why data varies?

learning data is about the system/phenomenon/nature S

- ▶ a collection of *observations* of *S*
- a point of view on S

"model varies if learning data varies": what? why data varies?

- learning data is about the system/phenomenon/nature S
  - ▶ a collection of *observations* of *S*
  - a point of view on S
- learning is about understanding/knowing/explaining S

"model varies if learning data varies": what? why data varies?

learning data is about the system/phenomenon/nature S

- a collection of observations of S
- a point of view on S
- learning is about understanding/knowing/explaining S
  - if I change the point of view on S, my knowledge about S should remain the same!

# Spotting overfitting



<□ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □

Test error: error on unseen data

# Spotting overfitting



<□ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □

Test error: error on unseen data

### k-fold cross-validation

Where can I find "unseen data"? Pretend to have it!

- 1. split learning data (**X** and **y**) in k equal slices (each of  $\frac{n}{k}$  observations/elements)
- 2. for each split (i.e., each  $i \in \{1,\ldots,k\}$  )
  - 2.1 learn on all but k-th slice
  - 2.2 compute classification error on unseen k-th slice
- 3. average the k classification errors

In essence:

- can the learner generalize beyond available data?
- how the learned artifact will behave on unseen data?

## k-fold cross-validation



Or with any other meaningful (effectiveness) measure

Q: how should data be split?

### Fighting overfitting with trees

- large k<sub>min</sub> (large w.r.t. what?)
- when building, limit depth
- when building, don't split if low overall impurity decrease
- after building, prune

### Pruning: high level idea

- 1. learn a full tree  $t_0$
- 2. build from  $t_0$  a sequence  $T = \{t_0, t_1, \dots, t_n\}$  of trees such that
  - $t_i$  is a root-subtree of  $t_{i-1}$   $(t_i \subset t_{i-1})$
  - $\blacktriangleright$   $t_i$  is always less complex than  $t_{i-1}$
- 3. choose the  $t \in T$  with minimum classification error with *k*-fold cross-validation

### k-fold cross-validation: data splitting

#### **Q:** how should data be split? Example: Android Malware detection

- Gerardo Canfora et al. "Effectiveness of opcode ngrams for detection of multi family android malware". In: Availability, Reliability and Security (ARES), 2015 10th International Conference on. IEEE. 2015, pp. 333–340
- Gerardo Canfora et al. "Detecting android malware using sequences of system calls". In: Proceedings of the 3rd International Workshop on Software Development Lifecycle for Mobile. ACM. 2015, pp. 13–20

Using cross-validation (CV) for assessment (I)

How the learned artifact will behave on unseen data?

More precisely: How an artifact learned with **this learning technique** will behave on unseen data?

# Using CV for assessment (II)

"This learning technique" = BUILDDECISIONTREE() with  $k_{min} = 10$ 

- 1. repeat k times
  - 1.1 BUILDDECISIONTREE() with  $k_{min} = 10$  on all but one slice

<□ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ↓ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ } < □ ∧ < □ ∧ < □ ∧ < □ ∧ < □ ∧ < □ ∧ < □ ∧ < □ ∧ < □ ∧ < □ ∧ < □ ∧ < □ ∧ < □ ∧ < □ ∧ < □ ∧ < □ ∧ < □ ∧ < □ ∧ < □ ∧ < □ ∧ < □ ∧ < □ ∧ < □ ∧ < □ ∧ < □ ∧ < □ ∧ < □ ∧ < □ ∧ < □ ∧ < □ ∧ < □ ∧ < □ ∧ < □

- k-1/k n observations in each X passed to BUILDDECISIONTREE()
- $1.2\,$  compute classification error on left out slice
- 2. average computed classification errors
- k invocations of BUILDDECISIONTREE()

Using CV for assessment (III)

"This learning technique" = BUILDDECISIONTREE() with  $k_{min}$  chosen automatically with a 10-fold CV

For assessing this technique, we do two nested CVs:

- 1. repeat k times
  - 1.1 choose  $k_{\min}$  among *m* values with 10-CV (repeat BUILDDECISIONTREE() 10*m* times) on all but one slice
    - $\frac{k-1}{k} \frac{9}{10}n$  observations in each **X** passed to BUILDDECISIONTREE()!
  - 1.2 compute classification error on left out slice
    - usually, a new tree is built on  $\frac{k-1}{k}n$  observations
- 2. average computed classification errors

(10m + 1)k invocations of BUILDDECISIONTREE()

"This learning technique" = BUILDDECISIONTREE() with  $k_{min}$  chosen automatically with a 10-fold CV

Using just one CV is cheating (cherry picking)!

- $k_{\min}$  is chosen exactly to minimize error on the full dataset
- conceptually, this way of "fitting" k<sub>min</sub> is similar to the way we build the tree

#### Subsection 1

Regression trees

< □ ▶ < ■ ▶ < ≧ ▶ < ≧ ▶ Ξ の Q @ 66/122

#### Regression with trees

Trees can be used for regression, instead of classification.

decision tree vs. regression tree

<□ ▶ < @ ▶ < ≧ ▶ < ≧ ▶ ≧ り Q <sup>®</sup> 67/122

Tree building: decision  $\rightarrow$  regression

function BUILDDECISIONTREE(X, y) if SHOULDSTOP(y) then  $\hat{y} \leftarrow \text{most common class in } \mathbf{y}$ **return** new terminal node with  $\hat{y}$ else  $(i, t) \leftarrow \text{BestBranch}(\mathbf{X}, \mathbf{y})$  $n \leftarrow$  new branch node with (i, t)append child BUILDDECISIONTREE( $\mathbf{X}|_{\mathbf{x}_i < t}, \mathbf{y}|_{\mathbf{x}_i < t}$ ) to *n* append child BUILDDECISIONTREE( $\mathbf{X}|_{\mathbf{x}_i > t}, \mathbf{y}|_{\mathbf{x}_i > t}$ ) to *n* return n end if end function

Q: what should we change?

Tree building: decision  $\rightarrow$  regression

function BUILDDECISIONTREE(X, y) if SHOULDSTOP(y) then  $\hat{v} \leftarrow \bar{v}$ ⊳ mean **y return** new terminal node with  $\hat{y}$ else  $(i, t) \leftarrow \text{BestBranch}(\mathbf{X}, \mathbf{y})$  $n \leftarrow$  new branch node with (i, t)append child BUILDDECISIONTREE( $\mathbf{X}|_{\mathbf{x}_i < t}, \mathbf{y}|_{\mathbf{x}_i < t}$ ) to *n* append child BUILDDECISIONTREE( $\mathbf{X}|_{\mathbf{x}_i > t}, \mathbf{y}|_{\mathbf{x}_i > t}$ ) to *n* return n end if end function **Q:** what should we change?

<□ → < □ → < Ξ → < Ξ → Ξ の Q ↔ 68/122</p>

#### Best branch

```
 \begin{array}{l} \text{function BESTBRANCH}(\mathbf{X}, \mathbf{y}) \\ (i^{\star}, t^{\star}) \leftarrow \arg\min_{i, t} E(\mathbf{y}|_{\mathbf{x}_i \geq t}) + E(\mathbf{y}|_{\mathbf{x}_i < t}) \\ \text{return } (i^{\star}, t^{\star}) \\ \text{end function} \end{array}
```

Q: what should we change?

### Best branch

function BESTBRANCH(**X**, **y**)  $(i^*, t^*) \leftarrow \arg \min_{i,t} \sum_{y_i \in \mathbf{y}|_{\mathbf{x}_i \geq t}} (y_i - \bar{y})^2 + \sum_{y_i \in \mathbf{y}|_{\mathbf{x}_i < t}} (y_i - \bar{y})^2$ return  $(i^*, t^*)$ end function

↓ □ ▶ ↓ □ ▶ ↓ ■ ▶ ↓ ■ ⑦ Q ○ 69/122

Q: what should we change?

Minimize sum of residual sum of squares (RSS) (the two  $\bar{y}$  are different)

Stopping criterion

function SHOULDSTOP(y) if y contains only one class then return true else if  $|y| < k_{min}$  then return true else return false end if end function

**Q**: what should we change?

Stopping criterion

function SHOULDSTOP(y)
 if RSS is 0 then
 return true
 else if |y| < k<sub>min</sub> then
 return true
 else
 return false
 end if
end function

**Q:** what should we change?

### Interpretation



◆□ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ <

## Regression and overfitting



Image from F. Daolio

## Trees in summary

Pros:

- easily interpretable/explicable
- ▲ learning and regression/classification easily understandable

▲ can handle both numeric and categorical values

Cons:

▼ not so accurate (Q: always?)

### Tree accuracy?



Image from An Introduction to Statistical Learning

く置≯

≣⇒

< 17 →

#### Subsection 2

Trees aggregation

< □ ▶ < ■ ▶ < ≧ ▶ < ≧ ▶ Ξ の Q @ 75/122

### Weakness of the tree



Small tree:

- Iow complexity
- will hardly fit the "curve" part
- high bias, low variance

Big tree:

- high complexity
- may overfit the noise on the right part
- Iow bias, high variance

### The trees view



#### Small tree:

"a car is something that moves"

#### Big tree:

"a car is a made-in-Germany blue object with 4 wheels, 2 doors, chromed fenders, curved rear enclosing engine"

### Big tree view

A big tree:

- has a detailed view of the learning data (high complexity)
- "trusts too much" the learning data (high variance)

What if we "combine" different big tree views and ignore details on which they disagree?

What if we "combine" different big tree views and ignore details on which they disagree?

- many views
- independent views
- aggregation of views

 $\approx$  the wisdom of the crowds: a collective opinion may be better than a single expert's opinion
many views

independent views

aggregation of views

many views
 just use many trees
 independent views

aggregation of views

many views
 just use many trees
 independent views

#### aggregation of views

 just average prediction (regression) or take most common prediction (classification)

#### many views

- just use many trees
- independent views
  - ► ??? learning is deterministic: same data ⇒ same tree ⇒ same view
- aggregation of views
  - just average prediction (regression) or take most common prediction (classification)

# Independent views $\equiv$ different points of view $\equiv$ different learning data

But we have only one learning data!

## Independent views: idea! (Bootstrap)

Like in cross-fold, consider only a part of the data, but:

- instead of a subset
- ▶ a sample with repetitions

#### Independent views: idea! (Bootstrap)

Like in cross-fold, consider only a part of the data, but:

- instead of a subset
- a sample with repetitions

$$\begin{split} \mathbf{X} &= (x_1^T x_2^T x_3^T x_4^T x_5^T) & \text{original learning data} \\ \mathbf{X}_1 &= (x_1^T x_5^T x_3^T x_2^T x_5^T) & \text{sample 1} \\ \mathbf{X}_2 &= (x_4^T x_2^T x_3^T x_1^T x_1^T) & \text{sample 2} \\ \mathbf{X}_i &= \dots & \text{sample } i \end{split}$$

- (y omitted for brevity)
- learning data size is not a limitation (differently than with subset)

# Tree **bagging**

When learning:

- 1. Repeat B times
  - $1.1\,$  take a sample of the learning data
  - 1.2 learn a tree (unpruned)

When predicting:

- 1. Repeat B times
  - 1.1 get a prediction from *i*th learned tree
- 2. predict the average (or most common) prediction

For classification, other aggregations can be done: majority voting (most common) is the simplest Using independent, possibly different classifiers together: *ensemble* of classifiers

#### How many trees?

B is a parameter:

when there is a parameter, there is the problem of finding a good value

remember k<sub>min</sub>, depth (Q: impact on?)

#### How many trees?

B is a parameter:

when there is a parameter, there is the problem of finding a good value

<□ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □

- remember k<sub>min</sub>, depth (Q: impact on?)
- it has been shown (experimentally) that
  - ▶ for "large" *B*, bagging is better than single tree
  - increasing B does not cause overfitting
  - ▶ (for us: default B is ok! "large" ≈ hundreds)
- Q: how better? at which cost?

```
Bagging: impact of B
```



<□▶ < □▶ < □▶ < □▶ < □▶ < □▶ = りへで 85/122

Despite being learned on different samples, bagging trees may be correlated, hence views are not very independent

 e.g., one variable is much more important than others for predicting (*strong predictor*)

Idea: force point of view differentiation by "hiding" variables

#### Random forest

When learning:

- 1. Repeat B times
  - $1.1\,$  take a sample of the learning data
  - 1.2 consider only m on p independent variables
  - 1.3 learn a tree (unpruned)

When predicting:

- 1. Repeat B times
  - 1.1 get a prediction from *i*th learned tree
- 2. predict the average (or most common) prediction
- (observations and) variables are randomly chosen...
- ... to learn a forest of trees
- Q: are missing variables a problem?

#### Random forest: parameter m

How to choose the value for m?

•  $m = p \rightarrow bagging$ 

- it has been shown (experimentally) that
  - m does not relate with overfitting
  - $m = \sqrt{p}$  is good for classification

• 
$$m = \frac{p}{3}$$
 is good for regression

#### Random forest

Experimentally shown: one of the "best" multi-purpose supervised classification methods

Manuel Fernández-Delgado et al. "Do we need hundreds of classifiers to solve real world classification problems". In: J. Mach. Learn. Res 15.1 (2014), pp. 3133–3181



but...

#### No free lunch!

"Any two optimization algorithms are equivalent when their performance is averaged across all possible problems"

David H Wolpert. "The lack of a priori distinctions between learning algorithms". In: Neural computation 8.7 (1996), pp. 1341–1390

#### Why free lunch?

- many restaurants, many items on menus, many possibly prices for each item: where to go to eat?
- no general answer
- but, if you are a vegan, or like pizza, then a best choice could exist
- **Q:** problem? algorithm?

## Observation sampling

When learning:

- 1. Repeat B times
  - 1.1 take a sample of the learning data
  - 1.2 consider only m on p independent variables (only for RF)
  - 1.3 learn a tree (unpruned)

Each learned tree uses only a portion of the observation in the learning data:

▶ for each observation,  $\approx \frac{B}{3}$  trees did not considere it when learned

## Observation sampling

When learning:

- 1. Repeat B times
  - 1.1 take a sample of the learning data
  - 1.2 consider only m on p independent variables (only for RF)
  - 1.3 learn a tree (unpruned)

Each learned tree uses only a portion of the observation in the learning data:

- ▶ for each observation,  $\approx \frac{B}{3}$  trees did not considere it when learned
- those observation were *unseen* for those trees, like in cross-validation (OOB = out-of-bag)

#### Bonus 1: OOB error

- for unseen each observation there are  $\frac{B}{3}$  predictions
- can "average" prediction among trees, observation and obtain an estimate of the testing error (OOB error)

<□ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □

- like with cross-fold validation
- for free!

#### OOB error



Image from An Introduction to Statistical Learning

Why estimating the test error?

Because the test data, in real world, is not available!will my ML solution work?

## Bagging/RF and explicability

- $\blacktriangleright \text{ Trees are easily understandable} \rightarrow \text{explicability}$
- Hundreds of trees are not!



х

Image from F. Daolio

## Bagging/RF and explicability: idea!

While learning:

- $1. \,$  for each tree, at each split
  - 1.1 keep note of the split variable
  - 1.2 keep note of RSS/Gini reduction
- 2. for each variable, sum reductions

The largest reduction, the more important the variable!

#### Bonus 2: variable importance

Instead of explicability based on tree shape:

importance of variables based on RSS/Gini reduction

◆□ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ <

### Nature of the prediction

Consider classification:

 $\blacktriangleright \ \mathsf{tree} \to \mathsf{the} \ \mathsf{class}$ 

• forest  $\rightarrow$  the class, as resulting from a voting

#### Nature of the prediction

Consider classification:

- ▶ tree  $\rightarrow$  the class
  - "virginica" is just "virginica"
- forest  $\rightarrow$  the class, as resulting from a voting
  - "241 virginica, 170 versicolor, 89 setosa" is different than "478 virginica, 10 versicolor, 2 setosa"

Different confidence in the prediction

## Bonus 3: confidence/tunability

Voting outcome:

- ▶ in classification, a measure of confidence of the decision
- in binary classification, voting threshold can be tuned to adjust bias towards one class (*sensitivity*)

<□ → < □ → < Ξ → < Ξ → Ξ の Q O 99/122</p>

**Q:** in regression?

#### Subsection 3

Binary classification

< □ ▶ < ■ ▶ < ≣ ▶ < ≣ ▶ ■ の へ <sup>∞</sup> 100/122

#### Binary classification

Binary classification:

one of the most common classes of problems

(comparative) evaluation is important!

#### Binary classification: evaluation

Consider the problem of classifying a person ('s data) as suffering or not suffering from a disease X.

Suppose we have "an accuracy of 99.99%". Q: is it good?

## Binary classification: positives/negatives

Consider the problem of classifying a person ('s data) as suffering or not suffering from a disease X.

- positive: an observation of "suffering" class
- negative: an observation of "not suffering" class

In other problems, positive may mean a different thing: define it!

#### Effectiveness indexes: FPR, FNR

Given some labeled data and a classifier for the disease X problem, we can measure:

- the number of negative observations wrongly classified as positives: False Positives (FP)
- the number of positive observations wrongly classified as negatives: False Negatives (FN)

To decouple FP, FN from data size:

$$FPR = \frac{FP}{N} = \frac{FP}{FP + TN}$$
$$FNR = \frac{FN}{P} = \frac{FN}{FN + TP}$$

<□ > < @ > < E > < E > E の Q ℃ 104/122

#### Relation of FPR, FNR with accuracy and error rate

$$\label{eq:Accuracy} \begin{aligned} &\mathsf{Accuracy} = 1 - \mathsf{Error} \ \mathsf{Rate} \\ &\mathsf{Error} \ \mathsf{Rate} = \frac{\mathsf{FN} + \mathsf{FP}}{\mathsf{P} + \mathsf{N}} \end{aligned}$$

↓ □ ▶ ↓ ● ▶ ↓ ■ ▶ ↓ ■ → ○ ○ 105/122

**Q:** Error Rate  $\stackrel{?}{=} \frac{\text{FPR} + \text{FNR}}{2}$ 

### FPR, FNR and sensitivity

- Suppose FPR = 0.06, FNR = 0.04 with threshold set to 0.5 (default for RF)
- $\blacktriangleright$  One could be interested in "limiting" the FNR  $\rightarrow$  change the threshold

Experimentally:



#### Comparing classifiers with FPR, FNR

Classifier A: FPR = 0.06, FNR = 0.04
Classifier B: FPR = 0.10, FNR = 0.01
Which one is the better?

We'd like to have one single index  $\rightarrow$  EER, AUC

#### Equal Error Rate (EER)



#### EER: the FPR at the value of t for which FPR = FNR
### AUC: Area Under the Curve



AUC: the area under the TPR vs. FPR curve, plotted for different values of threshold t

the curve is called the Receiver operating characteristic (ROC)

## ROC and comparison



・ロト・西ト・山田・山田・山口・

**Q:** what does the bisector represent?

#### Other issues: robustness w.r.t. the threshold



"Same" with other parameters

Other issues: robustness w.r.t. random components

Consider A vs. B, AUC measured with cross-fold validation:

- A: 0.85, 0.73, 0.91,  $\cdots \rightarrow \mu = 0.83, \sigma = 0.15$
- ▶ B: 0.81, 0.78, 0.79, · · ·  $\rightarrow \mu = 0.81, \sigma = 0.03$

Can we say that A is better than B? (for effectiveness only)

In general, other sources of performance variability:

- random seed
- subclass of problem class (e.g., image recognition of dogs, cats, ...)

#### Comparing techniques

Technique A, B; different index (e.g., AUC) values:

• 
$$A \rightarrow (x_a^1, x_a^2, ...) \rightarrow \text{random variable } X_a$$
  
•  $B \rightarrow (x_a^1, x_a^2, ...) \rightarrow \text{random variable } X_b$ 

▶ 
$$\mathsf{B} \to (x_b^1, x_b^2, \dots) \to \mathsf{random}$$
 variable  $X_b$ 

Do  $X_a, X_b$  follow different distributions?

- yes: A and B are different (concerning the AUC)
- ▶ no: difference in  $\mu_a, \mu_b$  might be due to randomness  $\rightarrow$  A, B are not significantly different

# Statistical significance in a nutshell

Just the way of thinking:

1. State a set of assumptions (the *null hypothesis*  $H_0$ ), e.g.:

► X<sub>a</sub>, X<sub>b</sub> are normally distributed and independent

• 
$$\bar{x}_a = \bar{x}_b \text{ (or } \bar{x}_a \geq \bar{x}_b \text{)}$$

any other assumption in the statistical model

# Statistical significance in a nutshell

Just the way of thinking:

- 1. State a set of assumptions (the *null hypothesis*  $H_0$ ), e.g.:
  - ► X<sub>a</sub>, X<sub>b</sub> are normally distributed and independent

• 
$$\bar{x}_a = \bar{x}_b \text{ (or } \bar{x}_a \geq \bar{x}_b \text{)}$$

- any other assumption in the statistical model
- 2. Perform a statistical test, appropriate choice depending on many factors, e.g.:

◆□ ▶ ◆□ ▶ ◆ ■ ▶ ◆ ■ ▶ ● ■ のへで 114/122

- Wilcoxon test (many versions)
- Friedman (many versions)
- ...

# Statistical significance in a nutshell

Just the way of thinking:

- 1. State a set of assumptions (the *null hypothesis*  $H_0$ ), e.g.:
  - ► X<sub>a</sub>, X<sub>b</sub> are normally distributed and independent

• 
$$\bar{x}_a = \bar{x}_b \text{ (or } \bar{x}_a \geq \bar{x}_b \text{)}$$

- any other assumption in the statistical model
- 2. Perform a statistical test, appropriate choice depending on many factors, e.g.:

◆□ ▶ ◆□ ▶ ◆ ■ ▶ ◆ ■ ▶ ● ■ のへで 114/122

- Wilcoxon test (many versions)
- Friedman (many versions)
- ▶ ...
- 3. . . . which outputs a  $\textit{p-value} \in [0,1]$ 
  - 0 is "good", 1 is "bad"

#### *p*-value: meaning

0 is "good", 1 is "bad"

The *p*-value is the degree to which the data conform to the pattern predicted by the null hypothesis

• 
$$p$$
-value =  $P(x_a^1, x_a^2, \dots, x_b^1, x_b^2, \dots | H_0)$ 

If *p*-value is low:

• we've been very (un)lucky in having observed  $x_a^1, x_a^2, \ldots, x_b^1, x_b^2, \ldots$ 

"maybe" because H<sub>0</sub> is not true

#### *p*-value: meaning

0 is "good", 1 is "bad"

The *p*-value is the degree to which the data conform to the pattern predicted by the null hypothesis

• *p*-value = 
$$P(x_a^1, x_a^2, ..., x_b^1, x_b^2, ... | H_0)$$

If *p*-value is low:

we've been very (un)lucky in having observed x<sup>1</sup><sub>a</sub>, x<sup>2</sup><sub>a</sub>,..., x<sup>1</sup><sub>b</sub>, x<sup>2</sup><sub>b</sub>,...

"maybe" because H<sub>0</sub> is not true

• Warning! Any part of  $H_0$ , not necessarily the  $\bar{x}_a = \bar{x}_b$  part!

#### Statistical significance

Things are much more complex than this...

#### Some interesting papers:

- Joaquín Derrac et al. "A practical tutorial on the use of nonparametric statistical tests as a methodology for comparing evolutionary and swarm intelligence algorithms". In: Swarm and Evolutionary Computation 1.1 (2011), pp. 3–18
- Cédric Colas, Olivier Sigaud, and Pierre-Yves Oudeyer. "How Many Random Seeds? Statistical Power Analysis in Deep Reinforcement Learning Experiments". In: arXiv preprint arXiv:1806.08295 (2018)
- Sander Greenland et al. "Statistical tests, P values, confidence intervals, and power: a guide to misinterpretations". In: European journal of epidemiology 31.4 (2016), pp. 337–350

#### Subsection 4

Boosting

< □ ▶ < ፼ ▶ < ≣ ▶ < ≣ ▶ ■ の < ? 117/122

# Many views and aggregation

In bagging/RF (regression):

- many views are different samples
- aggregation is average

Alternative:

many views are subsequent residuals

aggregation is the sum

# Boosting

When learning:

- 1. Current data is learning data
- 2. Repeat B times
  - 2.1 learn a tree on current data
  - 2.2 current data becomes residuals of learned tree  $(\textbf{y}-\hat{\textbf{y}})$

<□ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ↓ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪ < □ ♪

When predicting:

- 1. Repeat B times
  - 1.1 get a prediction from *i*th learned tree
- 2. sum prediction
- Q: implementation differences w.r.t. RF?

# Boosting (regression)

```
function BOOSTTREES(X, y)

t(X) \leftarrow 0

for i \in \{1, 2, ..., B\} do

t_i \leftarrow \text{BUILDREGRESSIONTREE}(X, y, d)

t(X) \leftarrow t(X) + \lambda t_i(X)

y \leftarrow y - \lambda t_i(X)

end for

return t

end function
```

Each learned tree should be simple (maximum splits d)

•  $\lambda$  slows down learning

Trickier with classification.

### Boosting parameters

- $\blacktriangleright$   $\lambda$  usually set to 0.01 or 0.001
- $\blacktriangleright$   $\lambda$  and *B* interact: for small  $\lambda$ , *B* should be large
- ▶ large *B* can lead to overfitting (unlike bagging/RF, **Q**: why)

◆□ ▶ ◆□ ▶ ◆ ■ ▶ ◆ ■ ▶ ● ■ のへで 121/122

Find a good value for B with cross-validation

(Both boosting and bagging general techniques)

# Bagging/RF/boosting in summary

interpretability numeric/categorical accuracy test error estimate variable importance confidence/tunability fast to learn (almost) non-parametric



▲□▶ ▲□▶ ▲ 三▶ ▲ 三▶ - 三 - のへぐ

\*: **Q:** how faster? when? does it matter?