

## Machine Learning 101

Jeanne Goossens
The Space Leiden

21-01-2021



Wat is Machine Learning?

Leren van data

Aan de slag met de Titanic



"Learning can be understood as a learning way to automatically find patterns and structure in data by optimizing the parameters of the model" (source).

"Learning can be understood as a learning way to automatically find patterns and structure in data by optimizing the parameters of the model" (source).

"Learning can be understood as a learning way to automatically find patterns and structure in data by optimizing the parameters of the model" (source).

Relational Databases

Relational Databases

Relational Algebra

Computer Science

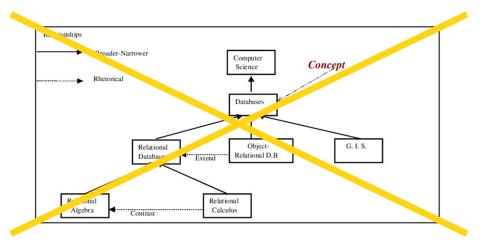
Concept

Concept

Relational Databases

Relational Calculus

"Learning can be understood as a learning way to automatically find patterns and structure in data by optimizing the parameters of the model" (source).



"Learning can be understood as a learning way to automatically find patterns and structure in data by optimizing the parameters of the model" (source).

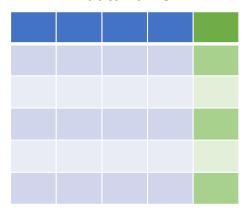
80% cleaning, 20% modelling

Onafhankelijke variabelen

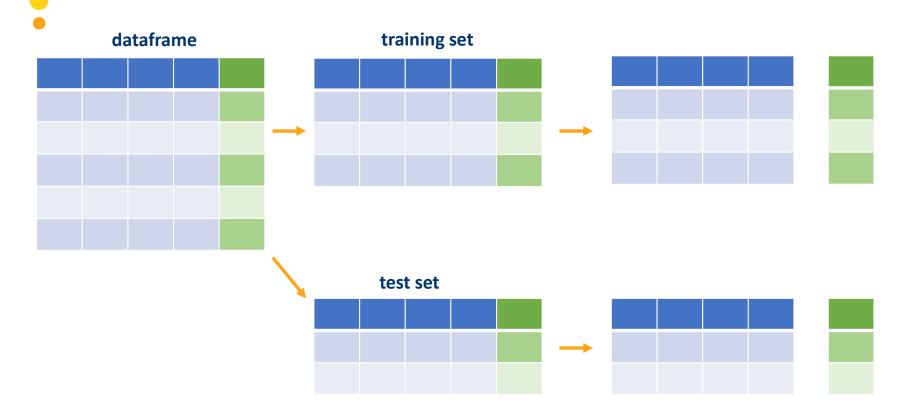
Afhankelijke variabele(n)

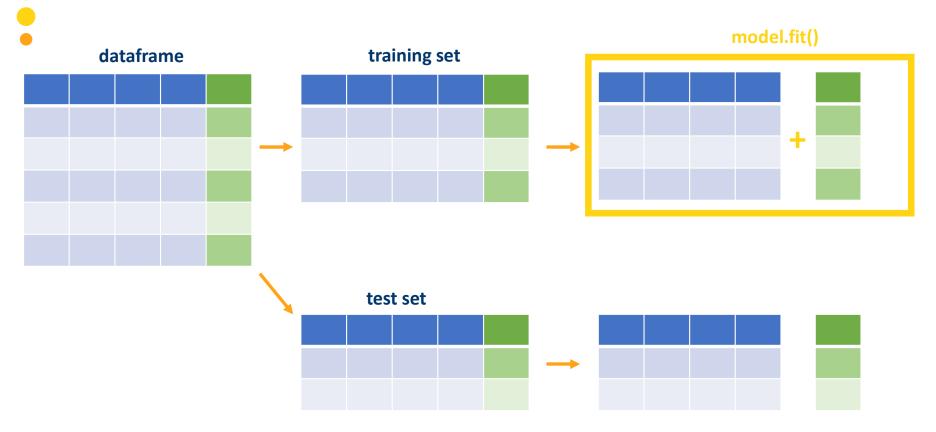
Passenger id	Class	Age	Gender	Survival

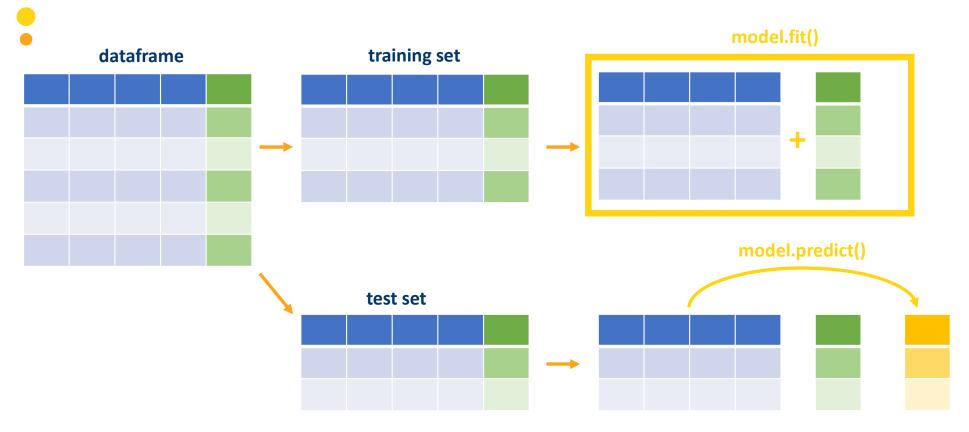


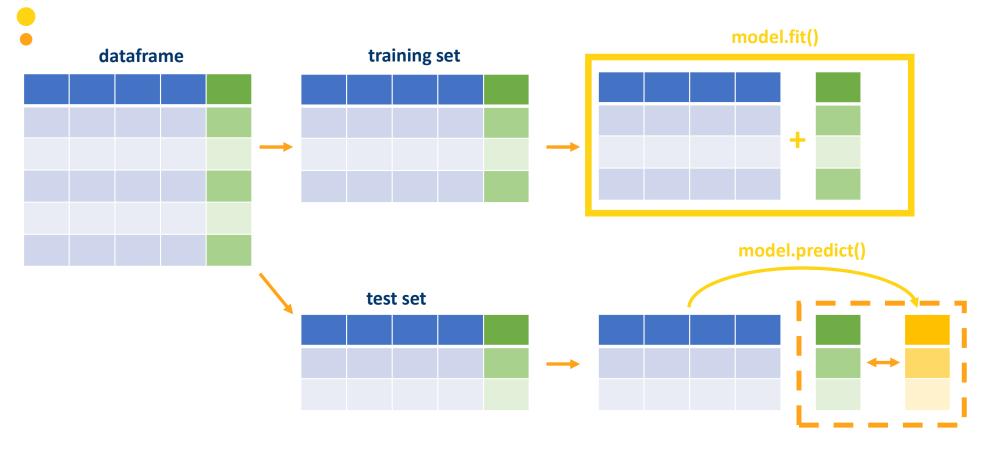














Als je lokaal wil draaien:

www.github.com/jeannegoossens/machine-learning-101

Als je in een virtual notebook wilt werken:

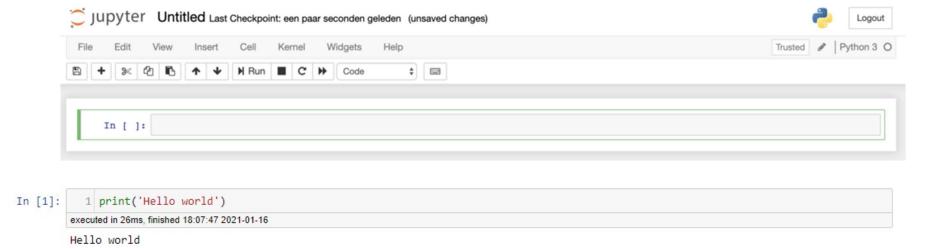
https://colab.research.google.com/drive/1\_de4HeQVsBSuXKEt5c-

gPTafY0k1b7R7?usp=sharing

## • Jupyter notebooks









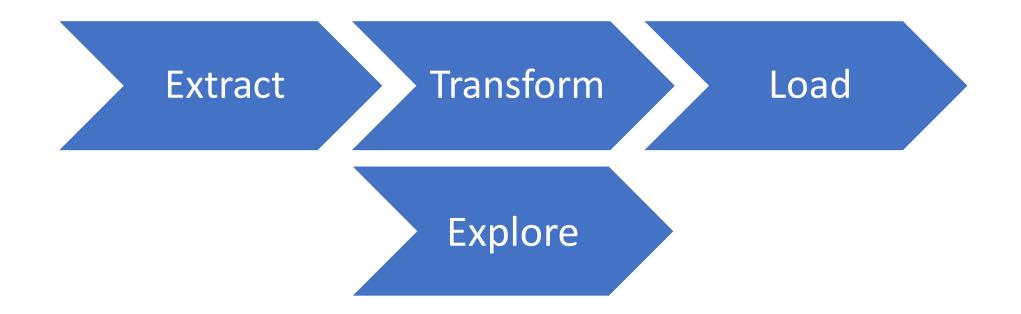
#### dataframe

class	name	sex	age	sibsp	parch	ticket	fare	cabin	emba rked	boat	body	home .dest	survi ved
	Allen, Miss. Elisabeth 1 Walton	female	29	0	0	24160	211.3375	B5	S	2	?	St Louis, MO	1
	Allison, Master. Hudson 1 Trevor	male	0.9167	1	2	113781	151.55	C22 C26	S	11	?	Montreal, PQ / Chesterville, ON	1
	Allison, Miss. Helen 1 Loraine	female	2	1	2			C22 C26	S	?	?	Montreal, PQ / Chesterville, ON	0
	Allison, Mr. Hudson Joshua 1 Creighton	male	30	1	2			C22 C26	S	?	135	Montreal, PQ / Chesterville, S ON	0
	Allison, Mrs. Hudson J C (Bessie Waldo 1 Daniels)	female	25	1	2	113781	151 55	C22 C26	S	2	2	Montreal, PQ / Chesterville, ON	0



Extract Transform Load







#### Out[3]:

	Passengerld	Survived	Pclass	Name	Sex	Age	SibSp	Parch	Ticket	Fare	Cabin	Embarked
0	431	Yes	1	Bjornstrom-Steffansson, Mr. Mauritz Hakan	male	28.0	0	0	110564	26.5500	C52	S
1	664	No	3	Coleff, Mr. Peju	male	36.0	0	0	349210	7.4958	NaN	S
2	44	Yes	2	Laroche, Miss. Simonne Marie Anne Andree	female	3.0	1	2	SC/Paris 2123	41.5792	NaN	С
3	347	Yes	2	Smith, Miss. Marion Elsie	female	40.0	0	0	31418	13.0000	NaN	S
4	891	No	3	Dooley, Mr. Patrick	male	32.0	0	0	370376	7.7500	NaN	Q

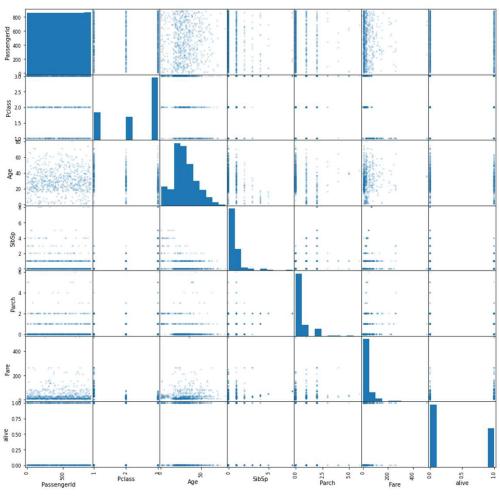
## Hypotheses?

Wat verwachten we dat de grootste voorspellende karakteristieken zijn?

- Leeftijd?
- Geslacht?
- Klasse?

Dit kunnen we testen op basis van onze data. Daarnaast kunnen we straks zien of het model deze kenmerken ook als relevante parameters ziet!

## • Explore



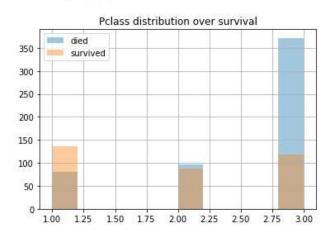
## • Hypotheses?

#### Klasse?

alive

#### Pclass

1 0.629630 2 0.472826 3 0.242363



## • Hypotheses?

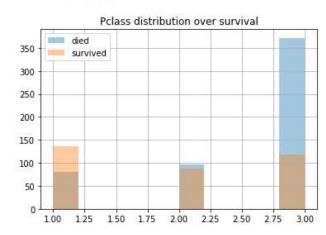
#### Klasse?

#### Geslacht?

alive

#### Pclass

1 0.629630 2 0.472826 3 0.242363



#### alive

Sex

female 0.742038

male 0.188908

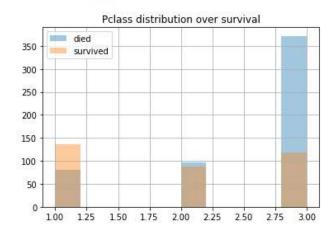
## • Hypotheses?

#### Klasse?

#### alive

#### Pclass

1 0.629630 2 0.472826 3 0.242363



#### Geslacht?

#### alive

Sex

female 0.742038

male 0.188908

#### Leeftijd?

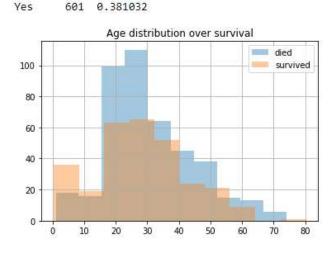
alive

count

mean

adult

lo 113 0.539823 'es 601 0.381032





Missende data opvullen?

- Alsnog achterhalen bij bron of Internet
- Achterhalen uit andere kolommen
- Opvullen met gemiddelde, modus, mediaan?
- Rij verwijderen of kolom verwijderen

Extra data toevoegen?

Embarked	
S	
С	
Q	
С	

Embarked
S
С
Q
С

Embarked_S	Embarked_C	Embarked_Q
1	0	0
0	1	0
0	0	1
0	1	0

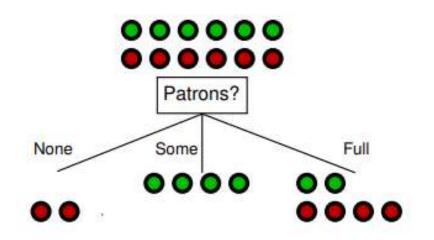


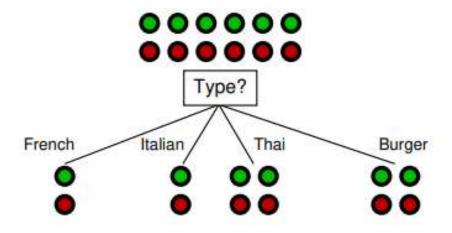
Embarked_S	Embarked_C	Embarked_Q
1	0	0
0	1	0
0	0	1
0	1	0

Em.	arked	
S		
С		
Q		
С		

Embarked_S	Embarked_C	Embarked_Q
1	0	0
0	1	0
0	0	1
0	1	0

## Model trainen: decision tree





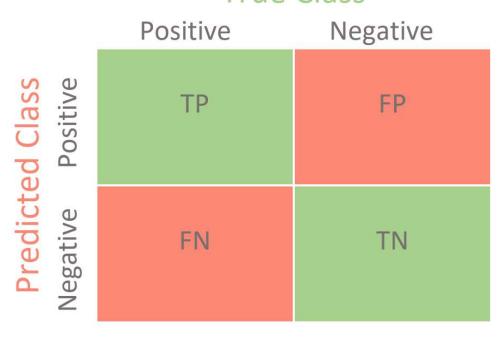
### Model trainen: decision tree

# Well Fit Decision Tree Ist split Srd split Srd split

Feature 1

## Model evalueren

#### True Class



## Model tweaken

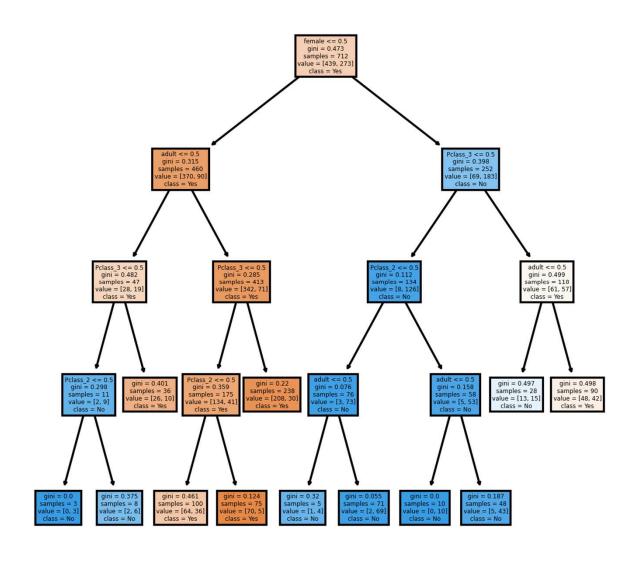
- Extra databewerking
  - Leeftijd categoriseren?
  - Persoonlijke titels gebruiken?
  - Fare categoriseren?

Probeer dingen uit, en test het effect op het model!

## Onze eigen hypotheses?

Wat zien we in het model terug van onze hypotheses? Wat als we een heel simple model draaien, met alleen die kolommen?

Het model wordt er beter van! Waarom is dat?



## Zou jij overleven?

Nu we een getraind model hebben, kunnen we ook kijken of wij zelf overleefd zouden hebben. Welke klasse zou jij minimal moeten reizen om de ramp overleefd te hebben?



Volgende week: The Space - Capture The Flag door Niels

Daarna: Github CI door Wesley