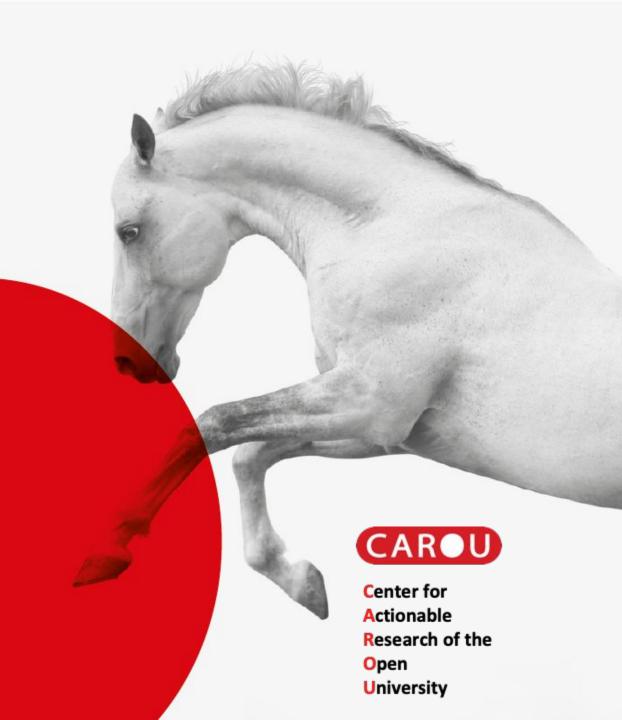


Online workshop Data Analytics

Lyana Curier, Daniele Di Mitri, Martine Hermans





Who we are?



Lyana Curier lyana.curier@gmail.com



Daniele DI MITRI daniele.dimitri@ou.nl



Martine HERMANS martine.hermans@ou.nl



Center for Actionable Research of the Open University



Agenda





Introduction

data rich, data driven world.

You've probably heard of kilobytes, megabytes, gigabytes, or even terabytes

By 2025, it's estimated that 463 exabytes of data will be created each day globally – that's the equivalent of 212,765,957 DVDs per day!



The science of analyzing raw data in order to draw conclusions.





Techniques and processes are automated into mechanical processes and algorithms to prepare raw data for human consumption.

descriptive
diagnostics
predictive
prescriptive

DATA ANALYTICS



Companies large and small, in traditional and non-traditional sectors, are using data to optimize their performance.

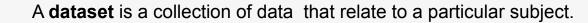






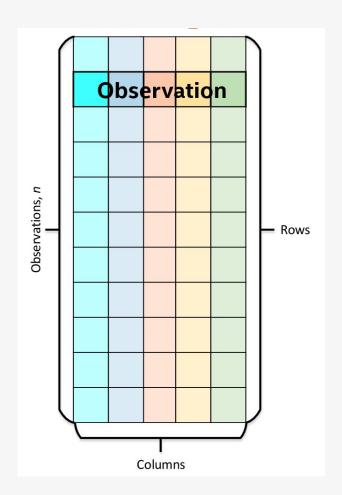






θ

tabular data



class	tal	per	sepal		
•00 00000000	width	th width length		length	
versicolor	1.3	4.4	2.3	6.3	
virginica	2.3	5.4	3.4	6.2	
setosa	0.2	1.4	3.4	5.2	
virginica	2.1	5.4	3.1	6.9	
setosa	0.4	1.5	4.4	5.7	
setosa	0.2	1.5	3.7	5.4	
setosa	0.2	1.4	3.3	5	
virginica	2.1	5.6	2.8	6.4	
virginica	1.8	4.8	3	6	
versicolor	1.3	4	2.5	5.5	

Iris Data set

THE CONCEPT OF DATA SET





tabular data

collection of images

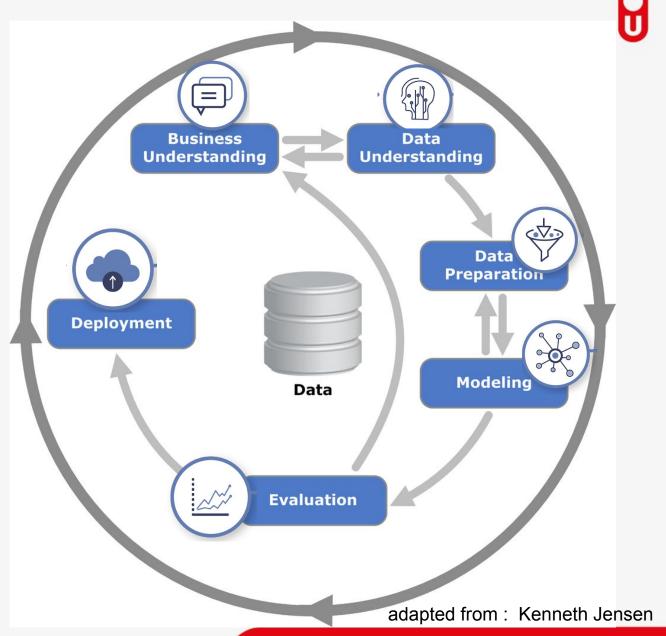


THE CONCEPT OF DATA SET

CRISP DM[†]:

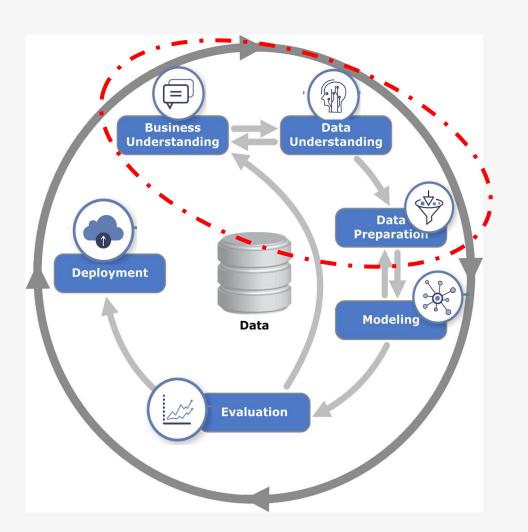
CRISP-DM breaks the process of data mining into six major steps.

sequence of the phases is not strict.





CRISP DM: Phase I (1/2)





thorough understanding of the business problem. Meet with stakeholders and domain/ subject experts to explicitly define "success criteria" for the project.



understand what data is collected in the business and for what purpose. Explore the data, present the data, assess quality and granularity. Formulate hypothesis on what can be found in the data.



given the hypothesis formulated in previous step access, transform, and condition available data into a format suitable for modeling and scoring i.e. for a machine learning prediction task



CRISP DM: Phase I (2/2)

Analyzing datasets to summarize their main characteristics Visual Exploration

- uncover underlying structure;
- extract important variables/features;

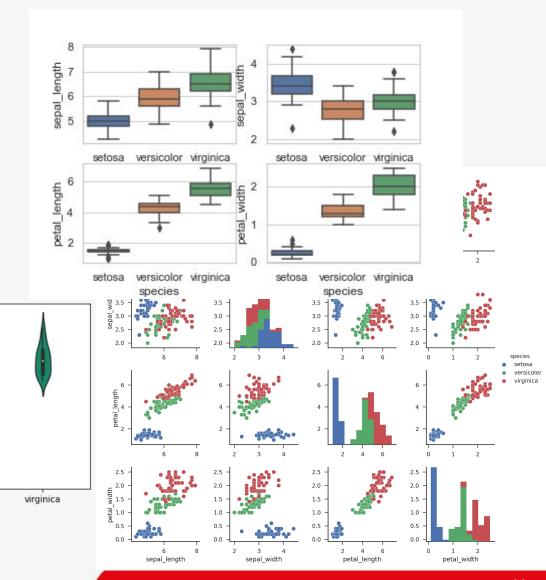
petal_length

setosa

versicolor

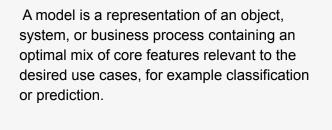
species

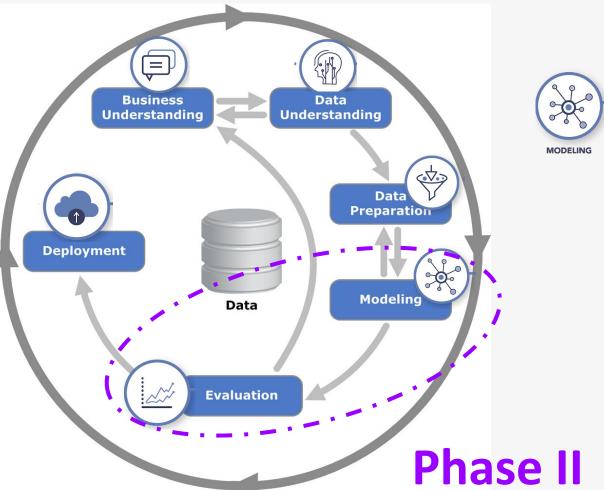
- detect outliers and anomalies;
- test underlying assumptions;





CRISP DM: Phase II (1/2)







Definition of a Parsimonious model determine optimal factor settings to represent the data/observations. for example regression, classification, clustering, text mining, computer vision

"Learning a Model"

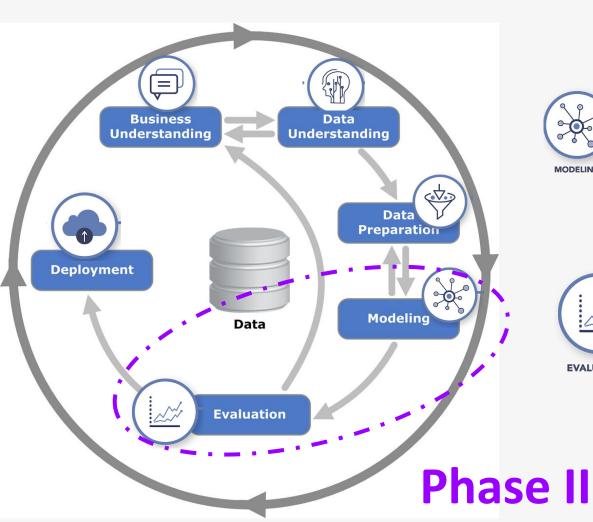
It is bottom up approach: maybe nothing or very little is known on the process generating the data.

It could be:

- a statistical model
- a set of rules
- a set of vectors
- a "description"



CRISP DM:Phase II (2/2)





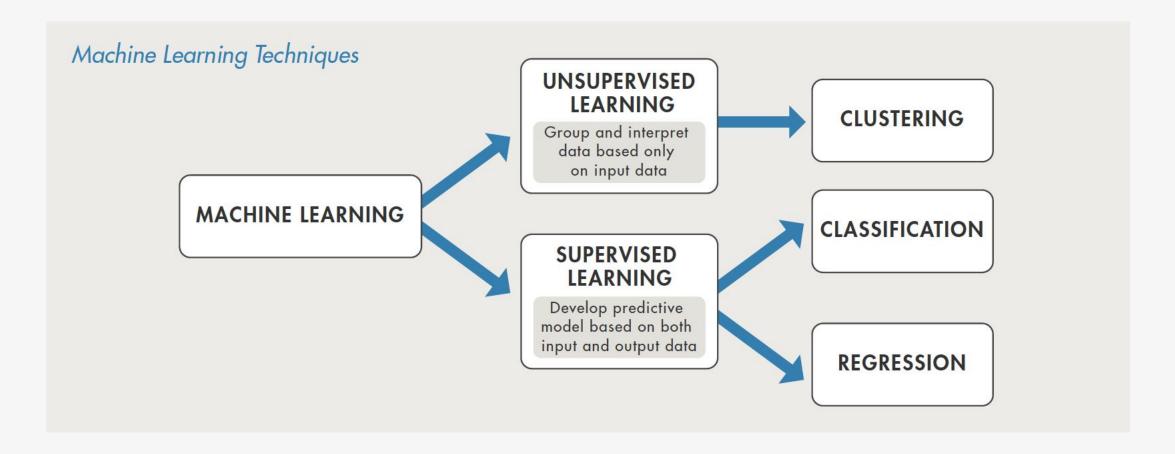
Definition of a Parsimonious model determine optimal factor settings to represent the data/observations. for example regression, classification, clustering, text mining, computer vision



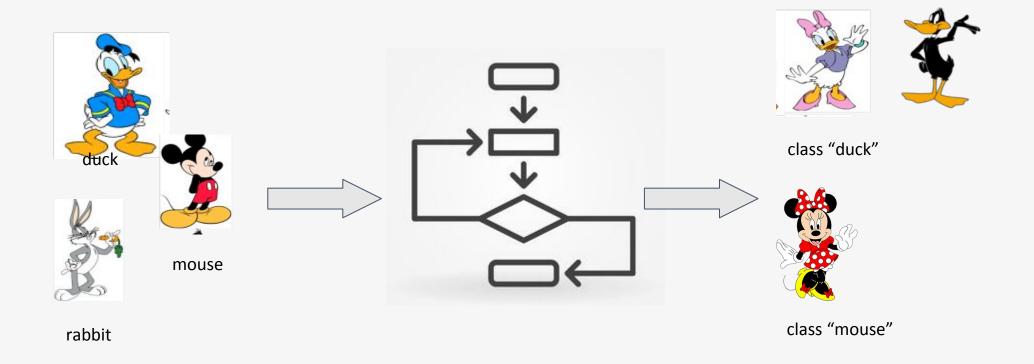
competing models are evaluated to determine which model (or model ensemble) best addresses the business objectives



CRISP-DM: Modelling





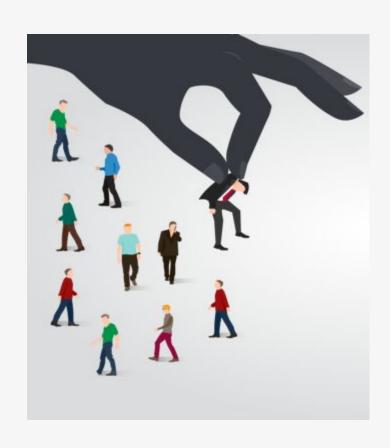


When you have a target

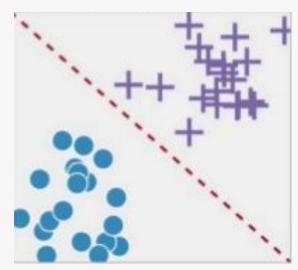
SUPERVISED LEARNING



Classification



attempts to predict, for each individual in a population, which of a (small) set of classes that individual belongs to

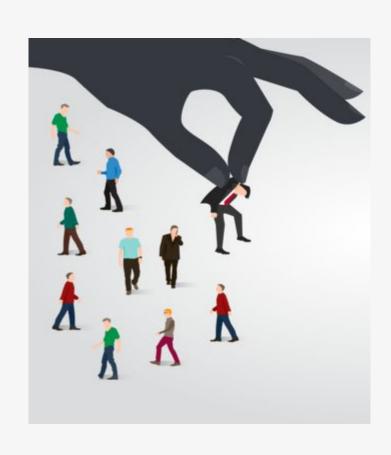


spam detection: binary classification spam and not spam.

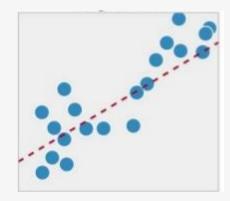
many more example: credit approval, medical diagnosis, target marketing, land use, etc...



Regression



attempts to estimate or predict, for each individual in a population, the numerical value of some variable.



housing market: Predict sale prices of a house given the features of house.

A regression problem where input variables are ordered by time is called a time series forecasting problem

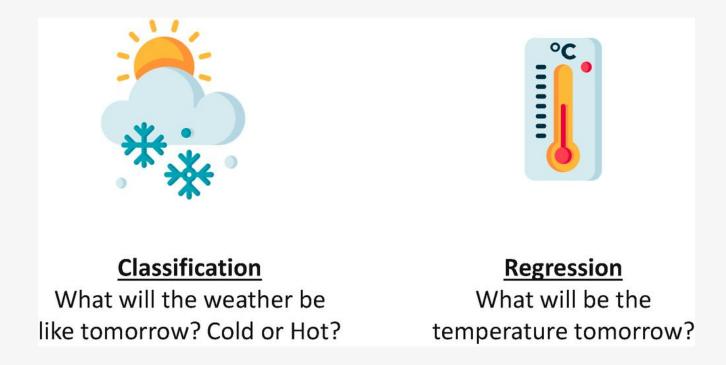


Examples: Classification vs. Regression

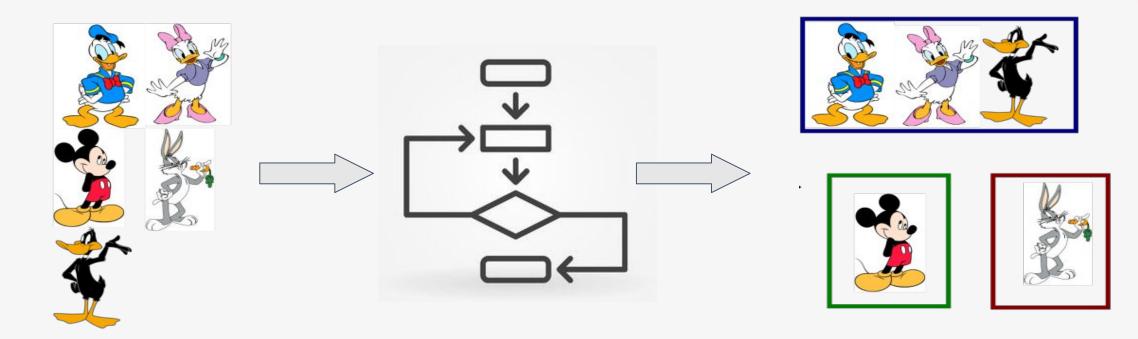
class	tal	pet	oal	ser	
-01 00000000	width	width length width		length	
versicolor	1.3	4.4	2.3	6.3	
virginica	2.3	5.4	3.4	6.2	
setosa	0.2	1.4	3.4	5.2	
virginica	2.1	5.4	3.1	6.9	
setosa	0.4	1.5	4.4	5.7	
setosa	0.2	1.5	3.7	5.4	
setosa	0.2	1.4	3.3	5	
virginica	2.1	5.6	2.8	6.4	
virginica	1.8	4.8	3	6	
versicolor	1.3	4	2.5	5.5	

Classify the iris species (virginica, setosa, or versicolor) based on the petal length, petal width

Predict petal length based on the iris species, petal width, sepal length, sepal width.







When you Don't have a target

UNSUPERVISED LEARNING

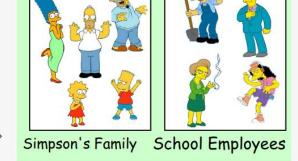


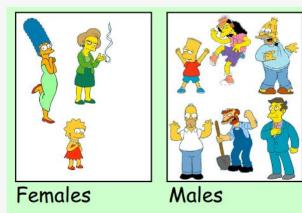


Aims to find homogeneous subgroups such that objects in the same group (clusters) are more similar to each other than the others.









(c) Eamonn Keogh, eamonn@cs.ucr.edu



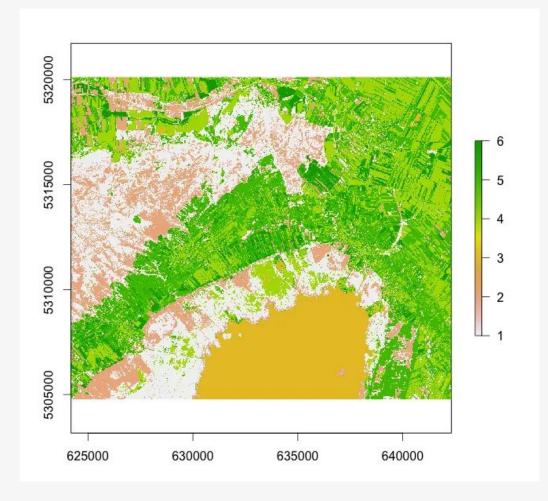
Examples: Clustering



true color composite (432)

a false color near infrared composite (843)

Sentinel-2 image





Learning Associations

Rule based **machine learning** and data mining technique that finds important relations between variables or features in a data set.

Basket analysis:

 n:
Given
MIVGII

- database of transaction
- each transaction is a list of items purchased by a customer in a visit

□ Find

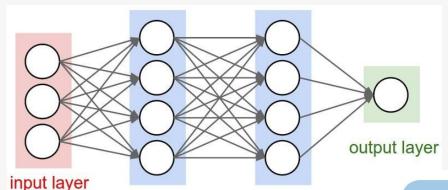
	all	rules tha	t correlate	e the	presence of	of one set	t of	items	with '	that o	fanot	her se	t of	it	ems
--	-----	-----------	-------------	-------	-------------	------------	------	-------	--------	--------	-------	--------	------	----	-----

Example:
$$P$$
 (chips | beer) = 0.7

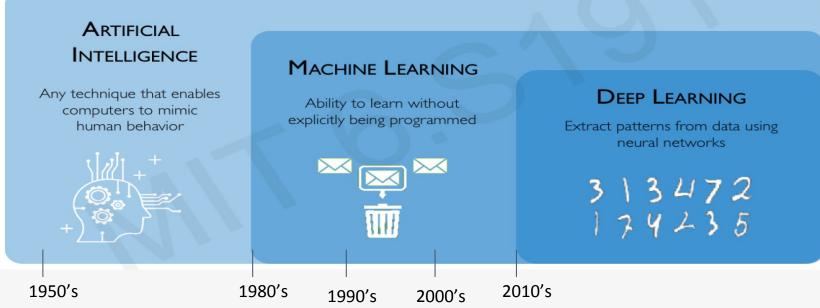
TID	Items
1	Bread, Milk
2	Bread, Diaper, Beer, Eggs
3	Milk, Diaper, Beer, Coke
4	Bread, Milk, Diaper, Beer
5	Bread, Milk, Diaper, Coke



Neural Networks & Deep Learning



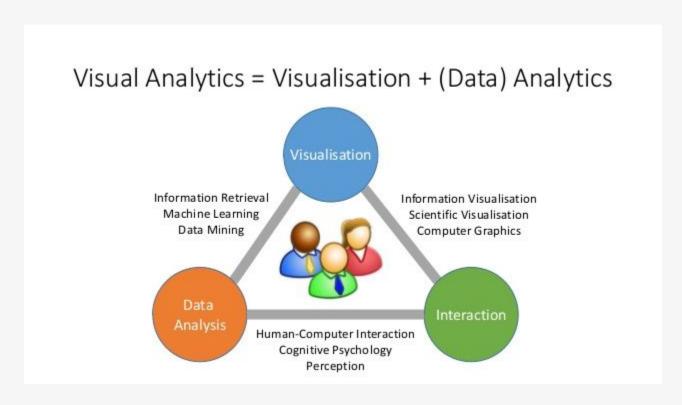
hidden layer 1 hidden layer 2



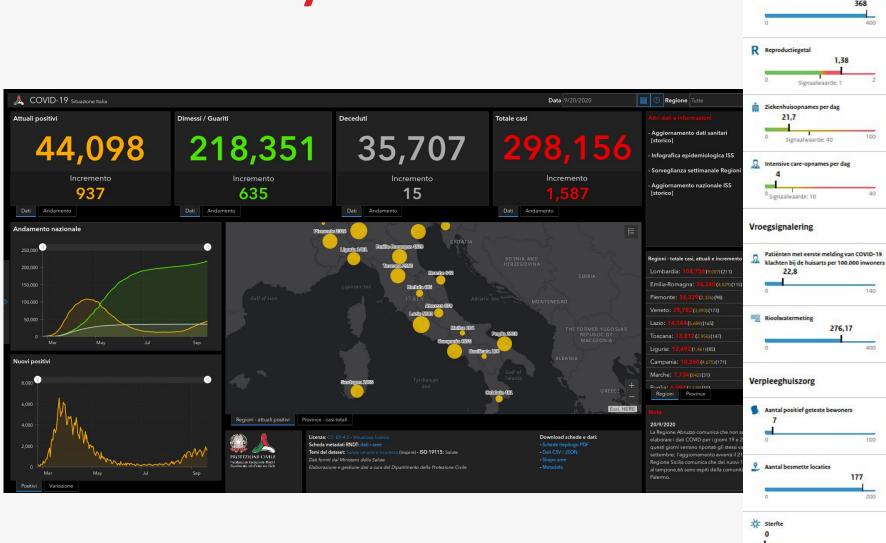


Visual Analytics

multidisciplinary field in which interactive visual interfaces are used to support analytical reasoning



Visual Analytics





Laatste ontwikkelingen Risiconiveaus per regio bekend

Medische indicatoren

Positief geteste mensen per 100.000 inwoners per dag

Aantal besmettelijke mensen per 100.000

Signaalwaarde: 7

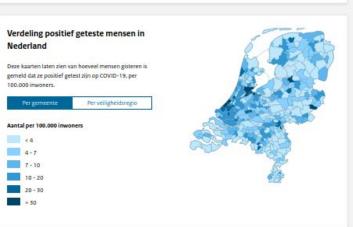
Taatste ontwikkelingen

Risiconiveaus per regio bekend

De afgelopen week (9 t/m 15 september) is het aantal nieuwe positief geteste personen wederom sterk toegenomen. De sterkste stijgingen waren weer te zien in de provincies Zuid-Holland en Noord-Holland. De nieuwe besmettingen nemen toe in alle leeftijdsgroepen, maar de meeste nieuwe besmettingen zijn gemeld in de leeftijdsgroep van 20-24 jaar. Ook het reproductiegetal is verder gestegen tot 1,38. De aantallen patiënten die vanwege COVID-19 zijn opgenomen in het ziekenhuis of op de intensive care zijn ook gestegen. Het risiconiveau per regio kan 'waakzaam' (niveau 1), 'zorgelijk' (niveau 2) of 'ernstig' (niveau 3) zijn. In 6 van de 25 regio's is de situatie

Bekijk de artuele informatie van het RIVM







Open Universiteit





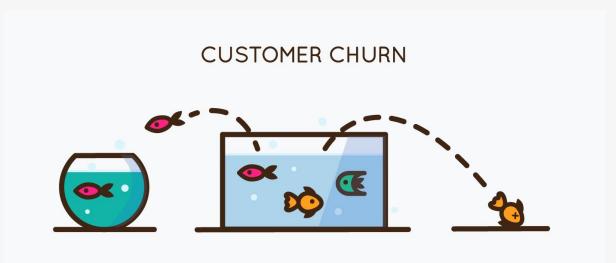
Break-out groups session



Customer Churn

THE CHALLENGE

A telecommunication company because of the strenuous competitions given by the concurrence competitors has to switch from growth strategy to protective approach of their existing customers. Their standard 'carpet' marketing campaigns had been delivering poor ROI and so they asked your team for help.





Solution design: the churn use case

SOLUTION DESIGN Design your data-driven solution for 'churn' prediction.

- 1. How do you define 'churn'? Provide an operational definition (e.g. no purchases in the last N months)
- 2. What kind of input data you need? (demographic data, customer service calls, transactional data, etc.)
- 3. How many data samples are you considering?
- 4. What machine learning approach is suitable for this task: classification, regression, clustering?
- 5. How is the TLC company supposed to use the outcomes of your solution?
- 6. What would be the benefits for TLC company in using your solution in the longer term?



Live tutorial: Data analytics in Python