

Data Science

Session 2 - Preparing data



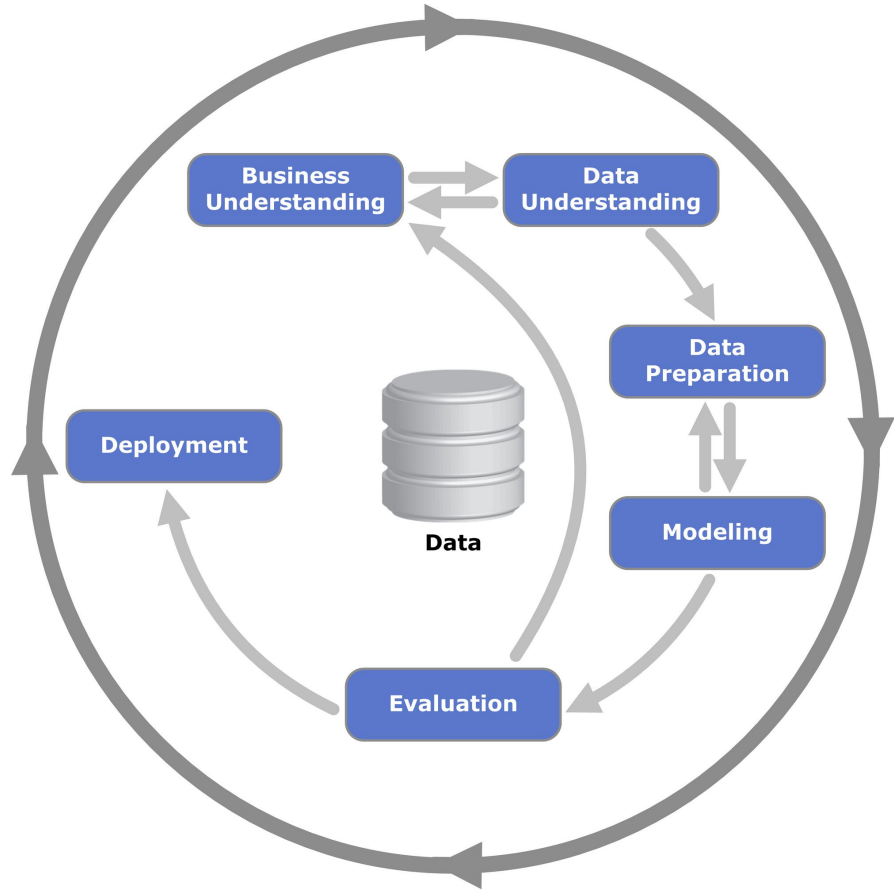
hadrien.salem@centralelille.fr



[introduction-to-data-science](#)

Introduction

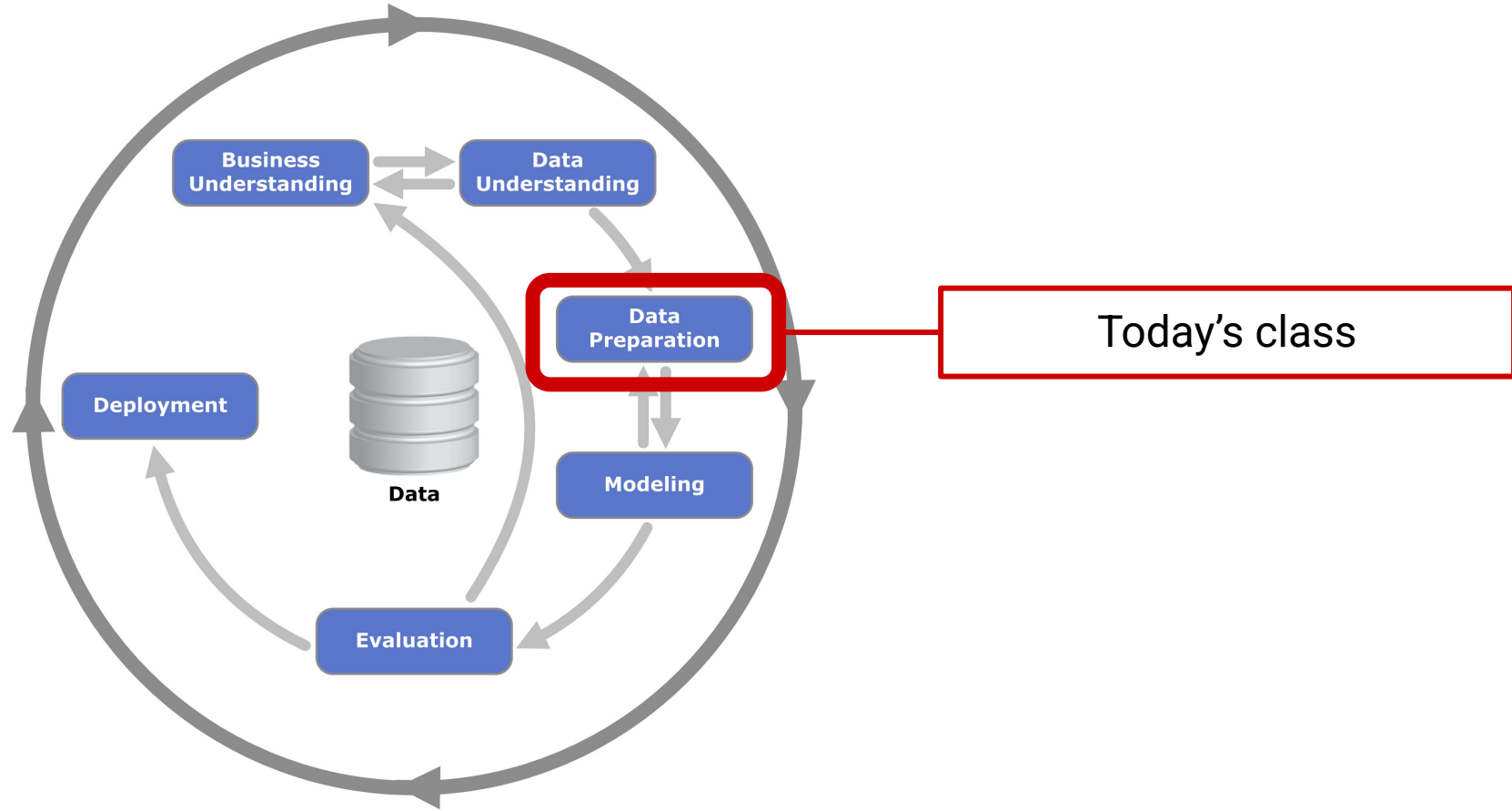
What did we do last time?



The CRISP-DM method

Cross-Industry Standard Process for Data Mining

- Published in 1999
- Common in the industry
- Still relevant today



Course outline

Data science course

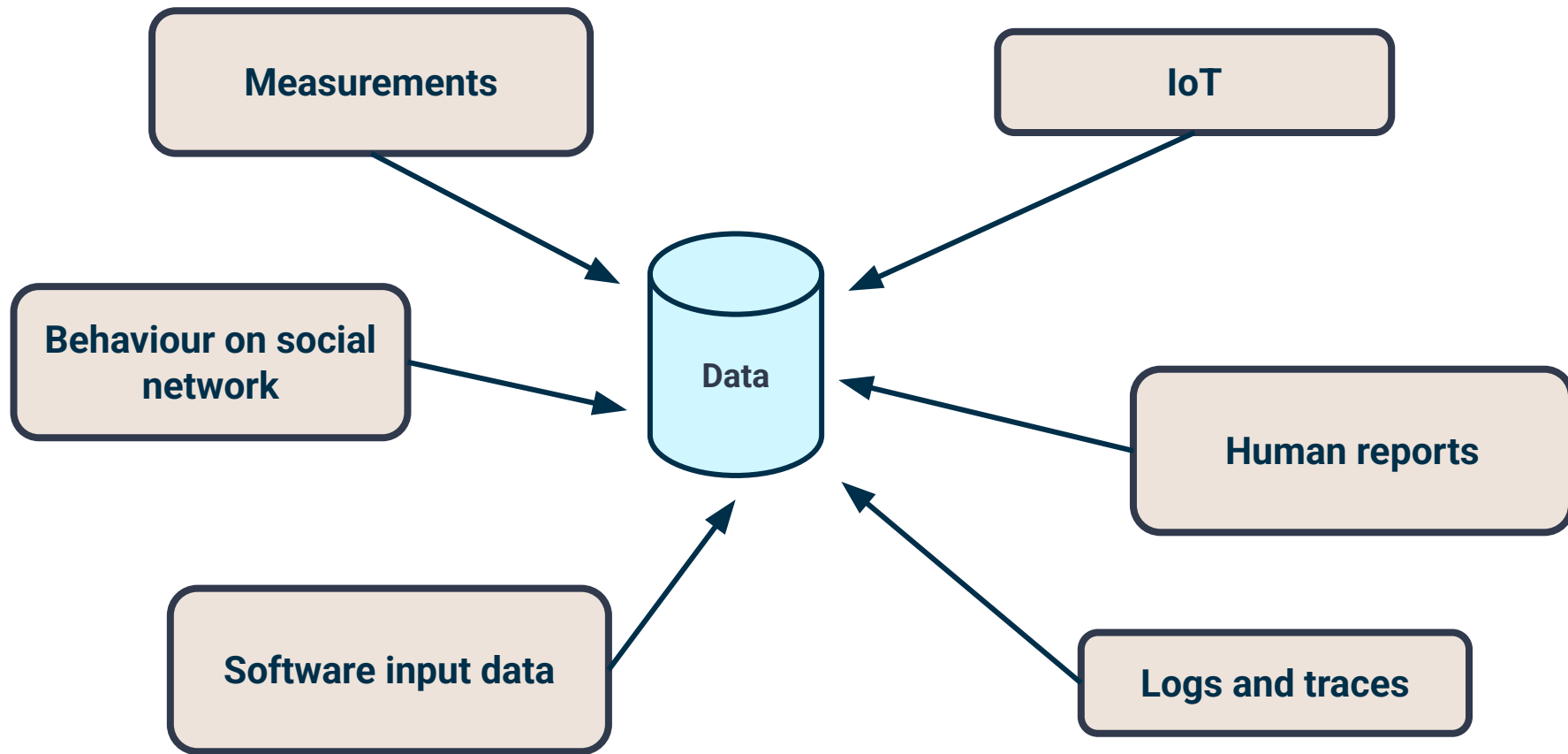
Session 1: Understanding data

Session 2: Preparing data

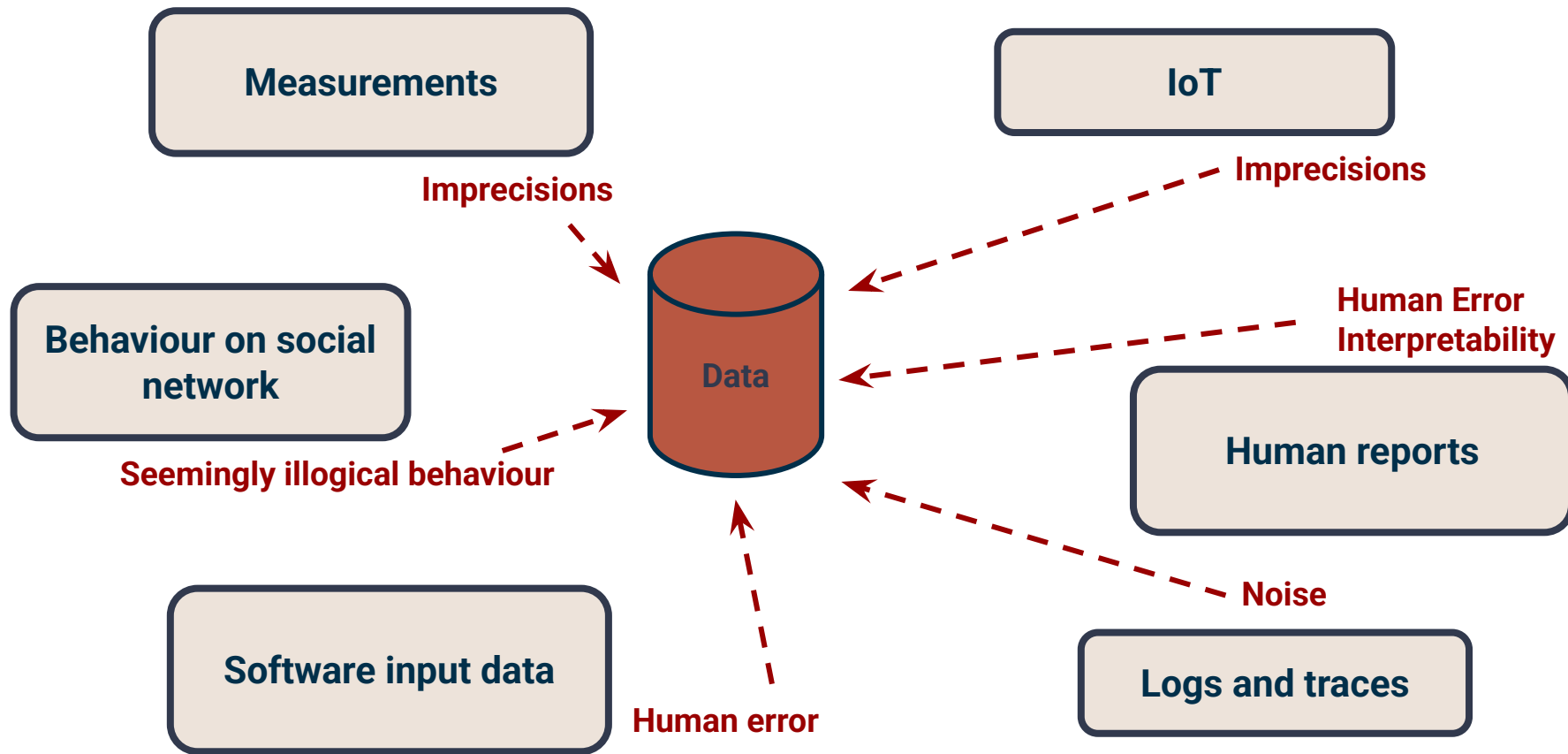


Machine learning course

What does it mean to prepare data?



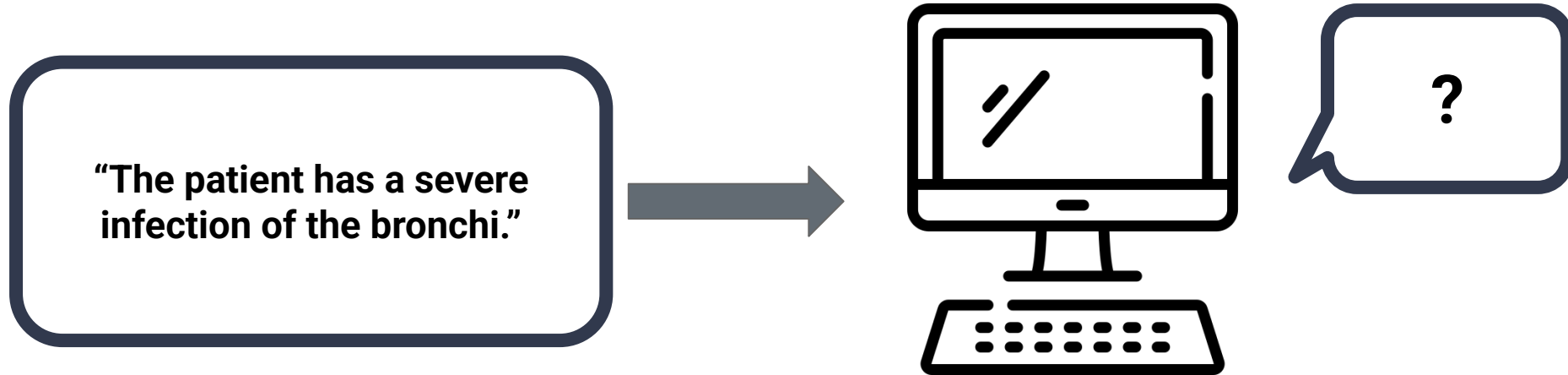
There are many sources of data...



...Which are all subject to uncleanliness

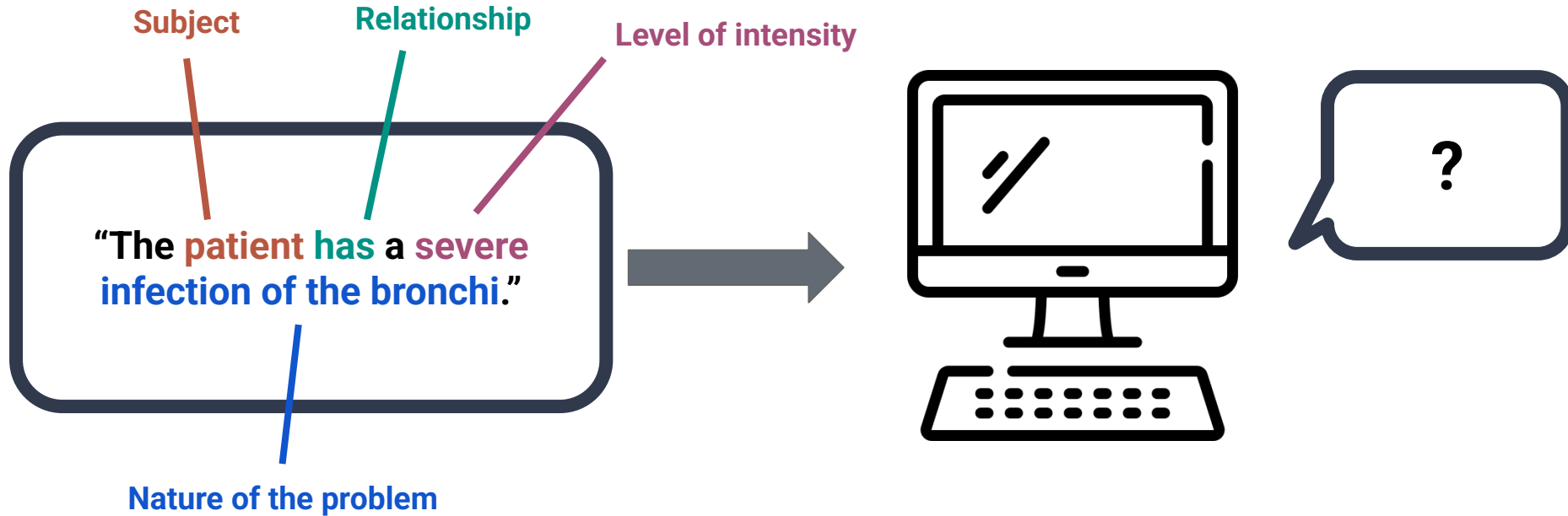
The risks of using unclean data

Example #1 : Inability to process data



The risks of using unclean data

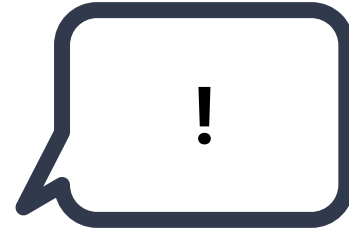
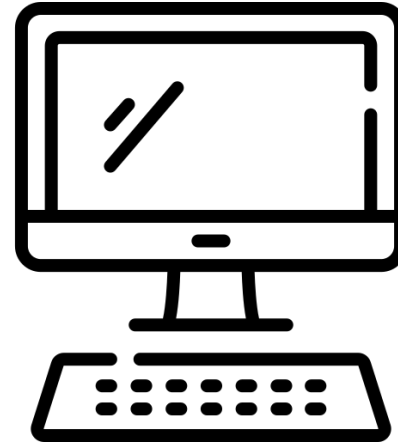
Example #1 : Inability to process data



The risks of using unclean data

Example #1 : Inability to process data

TARGET_TYPE: Patient
PROBLEM: Infection
LOCATION: 24
LOCATION_LABEL: Bronchi
INTENSITY: 3



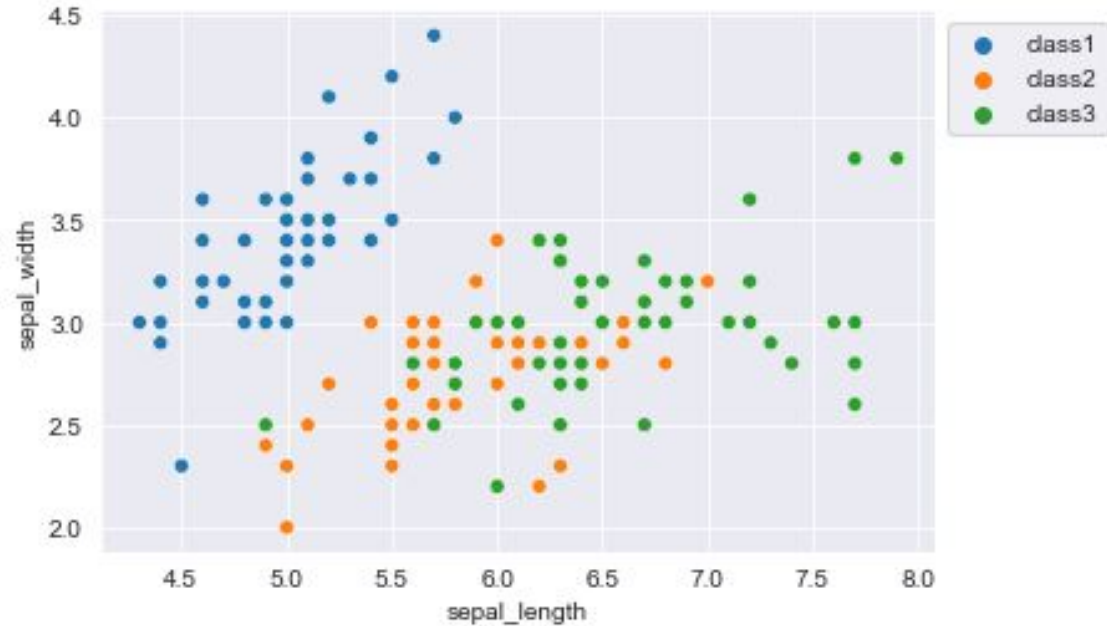
Data must be processed to be usable by a machine

The risks of using unclean data

Example #2 : Difficulty to model the data

An illustration from the first practical

Using only sepal information naively, the classification task is very difficult.



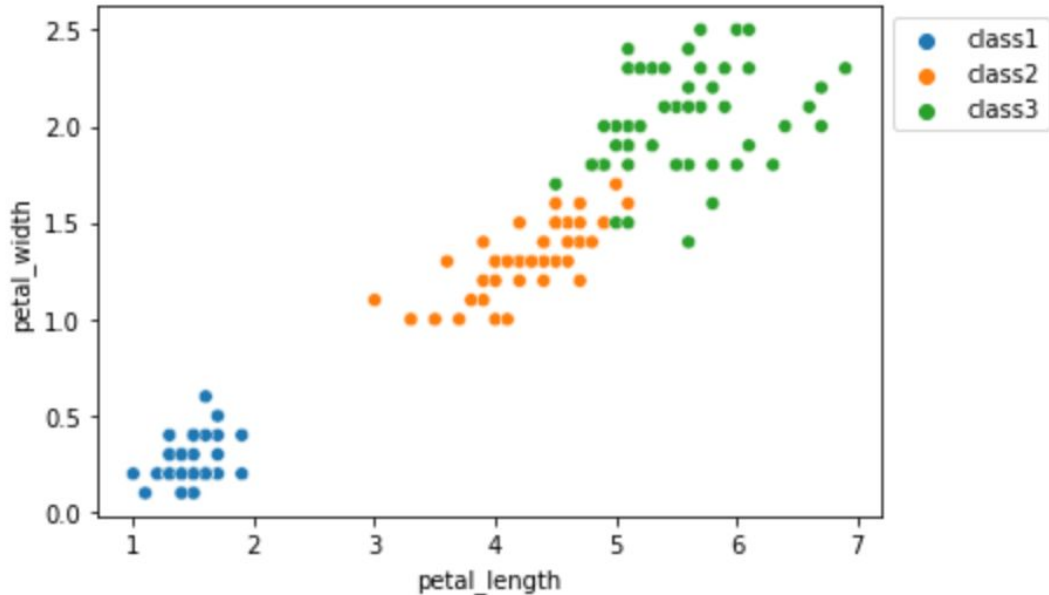
The risks of using unclean data

Example #2 : Difficulty to model the data

An illustration from the first practical

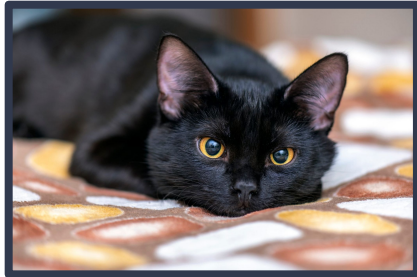
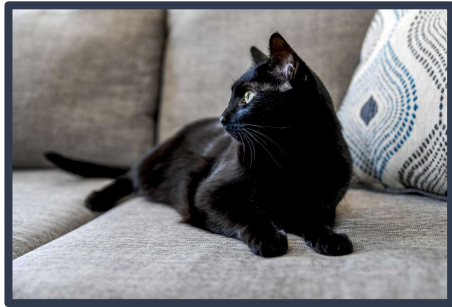
However, using petal information, it is much easier to choose a relevant model for classification.

Using relevant features is essential in machine learning.

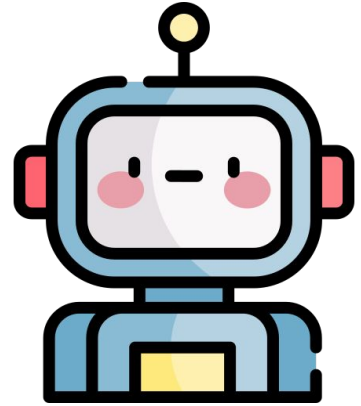


The risks of using unclean data

Example #3 : The introduction of bias

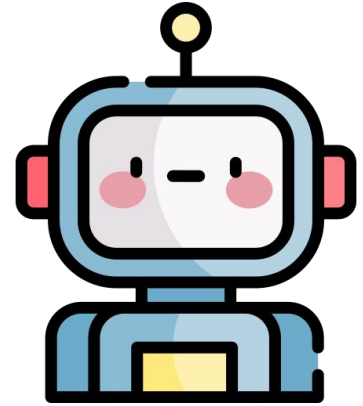


Learning



The risks of using unclean data

Example #3 : The introduction of bias



This is a white bear

The risks of using unclean data

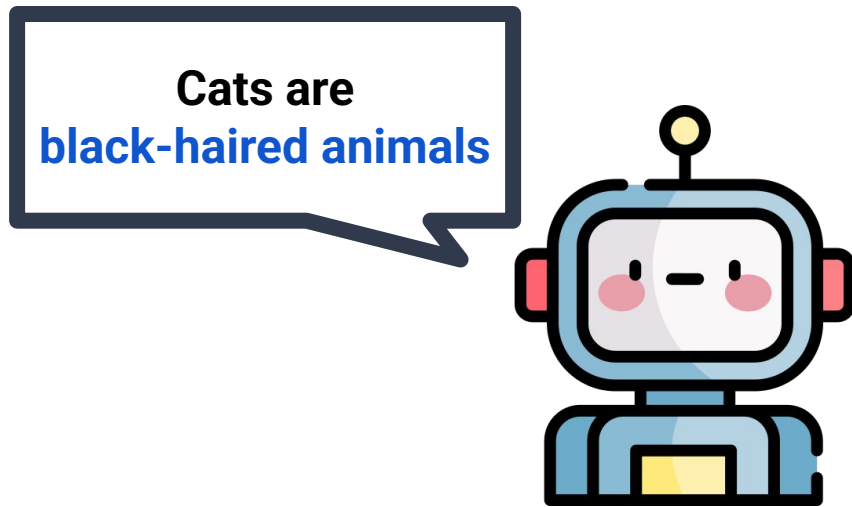
Example #3 : The introduction of bias

In the previous example, the training set is strongly biased.

Bias can have more severe consequences:

- ❖ Unusability in different regions
- ❖ Discrimination
- ❖ Sexism
- ❖ Maintaining human bias
- ❖ etc.

For the algorithms to **generalize** properly, bias is better avoided in a dataset.



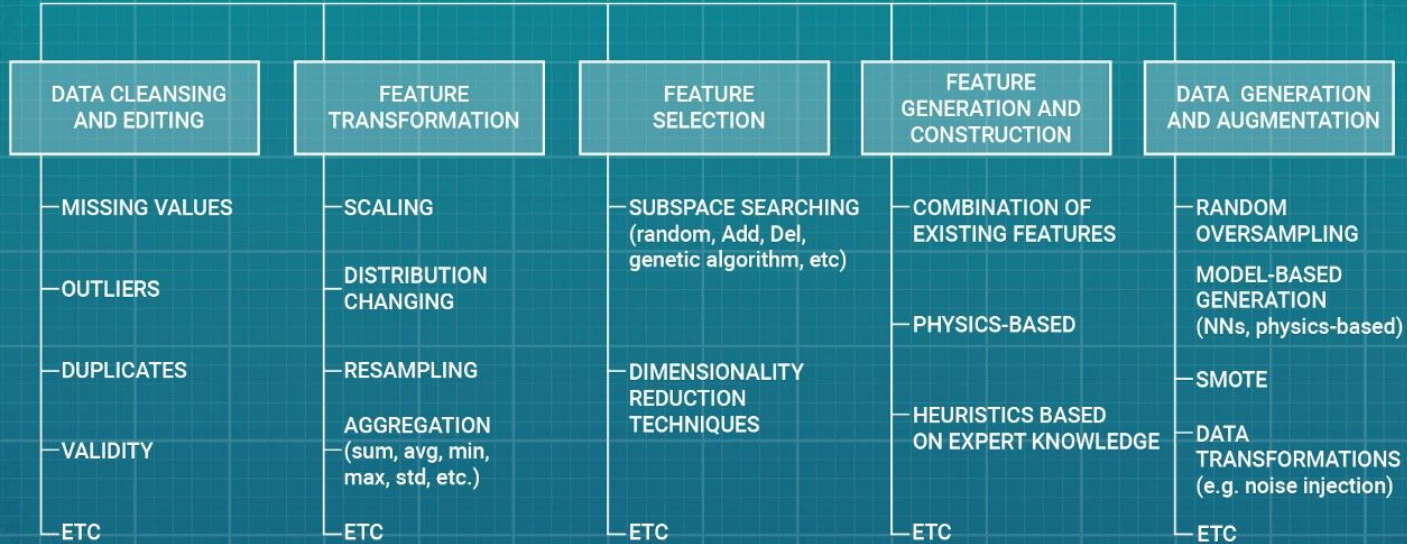
Preparing data is making it exploitable

Raw data is almost always **noisy** and **impractical**

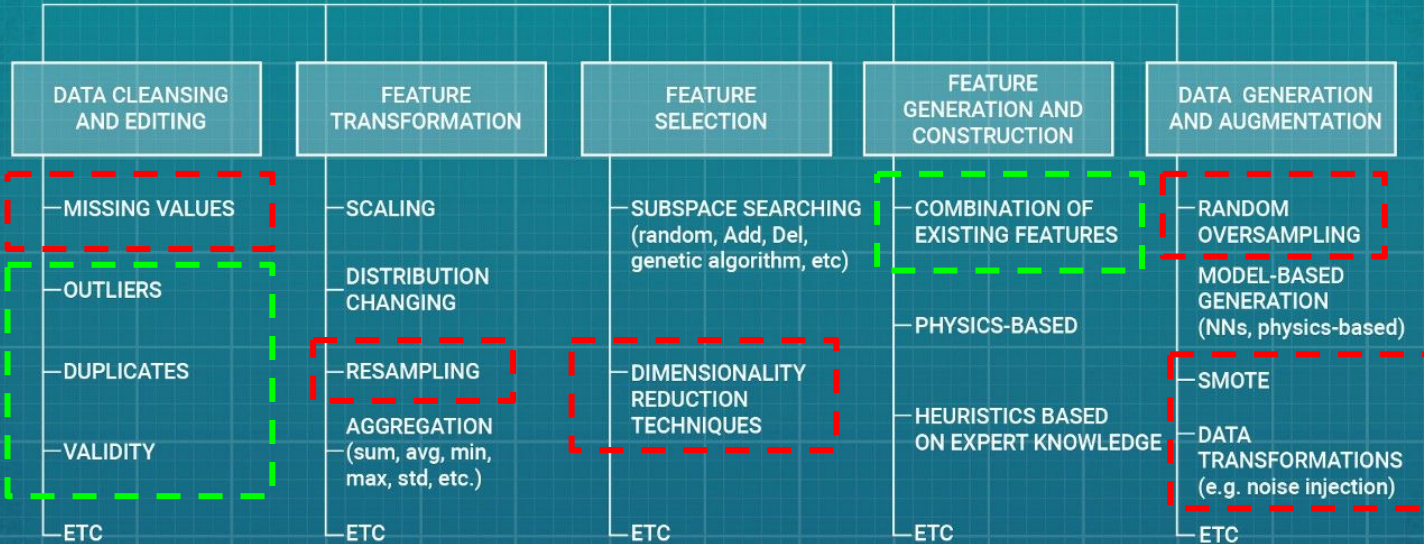
Preparatory work is systematically necessary for machine learning

How does one prepare data?

Data Pre-processing for ML



Data Pre-processing for ML



We have already seen some of those techniques, and we will see more.

In today's session,
we will see...

Missing data

Data imbalance

Dimensionality reduction

Managing missing data

Why can there be
missing data?



Why can there be missing data?

There can be several reasons...

- **Technical**
 - Faulty machines
 - Error during encoding
- **Human**
 - Typing errors
 - Deliberate choices (e.g. surveys)
- **Methodological**
 - Data collection not carried out for a certain part of the population (e.g. PSA for women)
 - Lack of measurement techniques

Why can there be missing data?

There are three main categories of missing data

- ***Missing Completely at Random***
 - There is no explainable pattern
 - Example : Human omitting to input data
- ***Missing at Random***
 - Patterns explainable from other columns
 - Example: A survey where men tend to reply less than women
- ***Missing Not at Random***
 - Explainable patterns, but not by observing the other columns
 - Example : People with lower incomes tend not to respond to questions about their salaries

Why can there be missing data?

There are three main categories of missing data

- ***Missing Completely at Random***
 - There is no explainable pattern
 - Example : Human omitting to input data
 - **Data missing for the training process**
- ***Missing at Random***
 - Patterns explainable from other columns
 - Example: A survey where men tend to reply less than women
 - **Generation of bias: The algorithm will generalize better for women**
- ***Missing Not at Random***
 - Explainable patterns, but not by observing the other columns
 - Example : People with lower incomes tend not to respond to questions about their salaries
 - **Generation of bias: The mean salary in the dataset will be inflated**

How can we deal
with missing data?



How can we deal with missing data?

Delete lines

Impute values

How can we deal with missing data?

Delete lines

Simple

Can drastically reduce the amount of data

Can introduce bias

Impute values

More robust with more missing data

We keep the “full” dataset

You have to experiment to find the best method

Can introduce bias or inconsistencies

Mean imputation

By definition, the mean is a value that makes some kind of sense.

It is however computed from observable data and can be influenced by existing bias.

Possible benefits

- Very simple to implement
- Gives a baseline with little effort

Possible pitfalls

- The mean is sensitive to outliers, especially if they are concentrated on one side of the distribution
- It reinforces the weight of the “mean individual”

Median imputation

In balanced datasets, the median tends to be close to the mean.

It is less sensitive to outliers.

NB: Outliers could also be managed specifically (removed or adjusted).

Possible benefits

- Very simple to implement
- Gives a baseline with little effort
- Less sensitive to outliers than the mean

Possible pitfalls

- Ignoring extreme values can be problematic in a dataset with high variance
- It reinforces the weight of the “median individual”

Random value imputation

Using random values can give surprising good results in machine learning.

Studying the distribution of features can help choose a probability distribution to draw from.

NB: This shows the importance of data visualization (cf. `kdeplot`) !

Possible benefits

- Not too difficult to implement
- The weight of existing values is not excessively increased

Possible pitfalls

- Finding a relevant probability distribution can be difficult
- The observable distribution could be biased
- Inconsistencies can be introduced in the data

Frequent value imputation

This method is mostly used for non-numerical data.

Similar to numerical data, more intelligent imputation methods can be implemented by studying the distribution of this data.

Possible benefits

- Extremely simple to implement

Possible pitfalls

- Giving too much weight to the most frequent value
- Introducing or maintaining bias

Interpolated value imputation

Interpolation is very useful when the value of a featured is determined by a known function.

Linear and polynomial interpolations are the most common.

NB: Here again, visualization can help.

Possible benefits

- If the actual distribution is close to the function we choose for interpolation, results can be very good

Inconvénients (possibles)

- Not always applicable in practice

Advanced imputation

Scikit-learn offers several imputers, such as `SimpleImputer`, `IterativeImputer`, or `KNNImputer`.

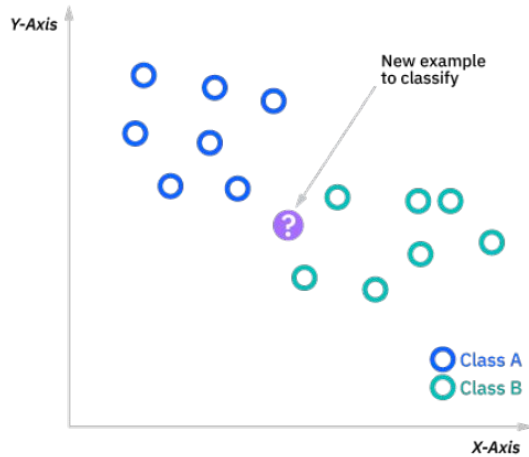
The `SimpleImputer` lets you do what we presented before, whereas the other two are based on machine learning.

Possible benefits

- These methods can help prevent the pitfalls listed before
- They are susceptible to find values that are close to the real ones

Inconvénients (possibles)

- Choosing the imputer is difficult
- The use of machine learning requires data to have been processed to some extent



Introduce a new example



Compute distances



Majority vote

K-nearest-neighbours

KNN is a very simple classification algorithm that can provide good results in some cases.

It can be used for data imputation.

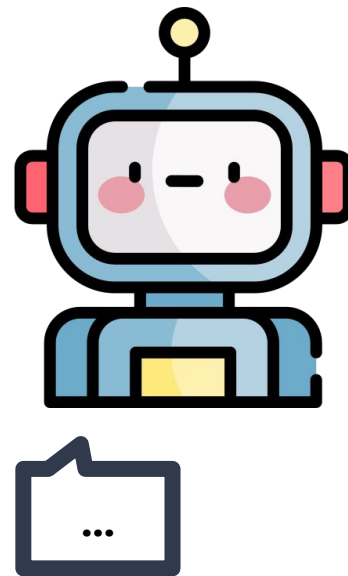
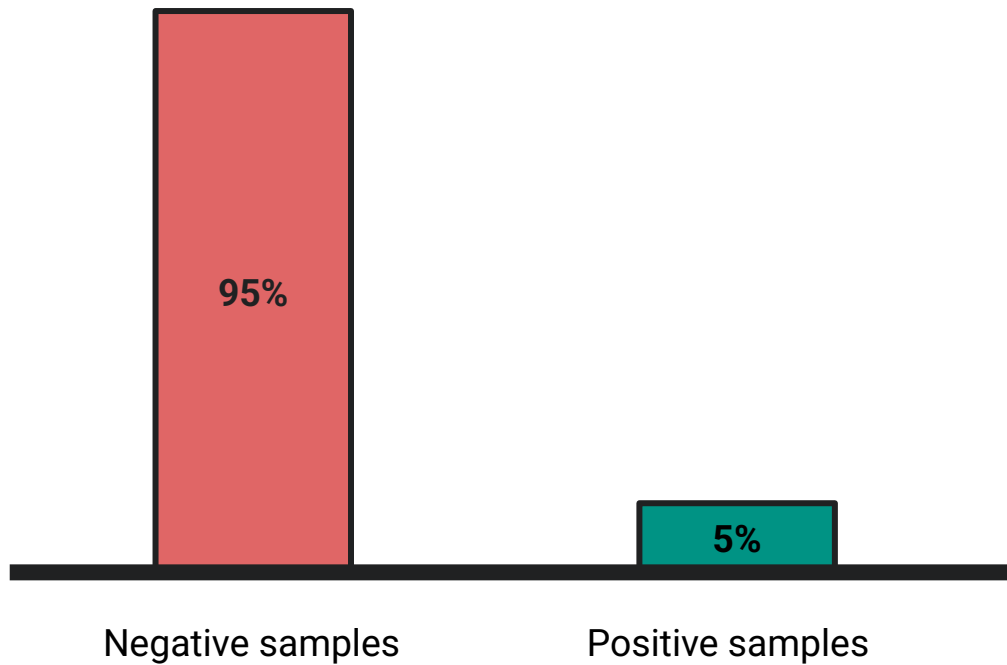
Possible benefits

- Easy to implement
- Few hyperparameters

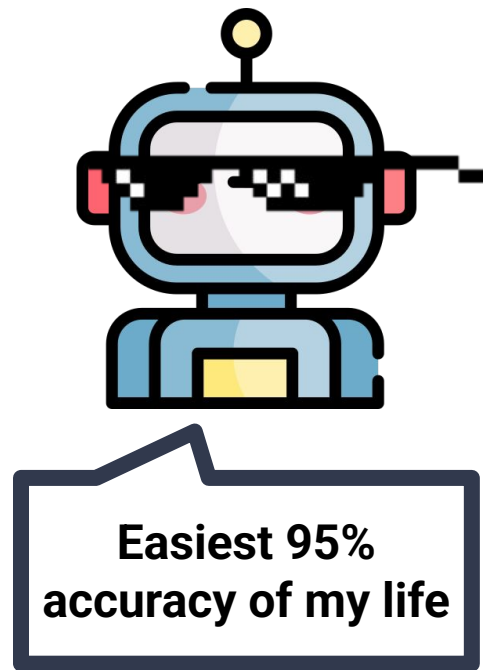
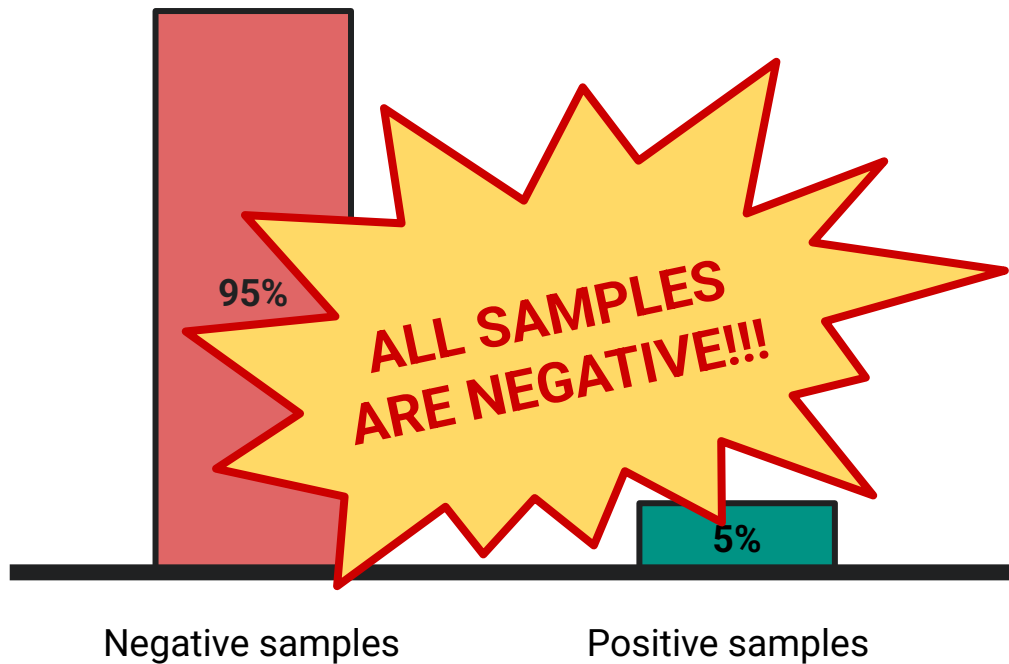
Possible pitfalls

- Choosing a distance is not always easy
- Can become computationally expensive
- Sensitive to the curse of dimensionality
- Sensitive to overfitting

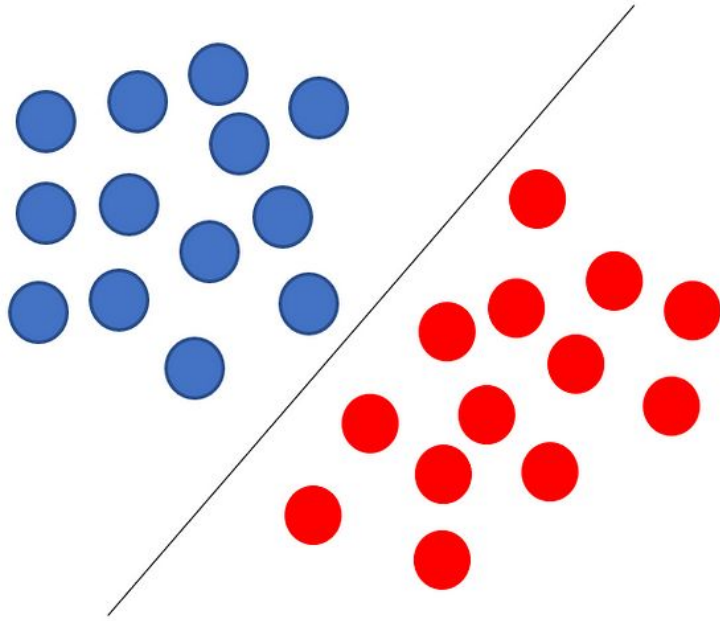
Data imbalance



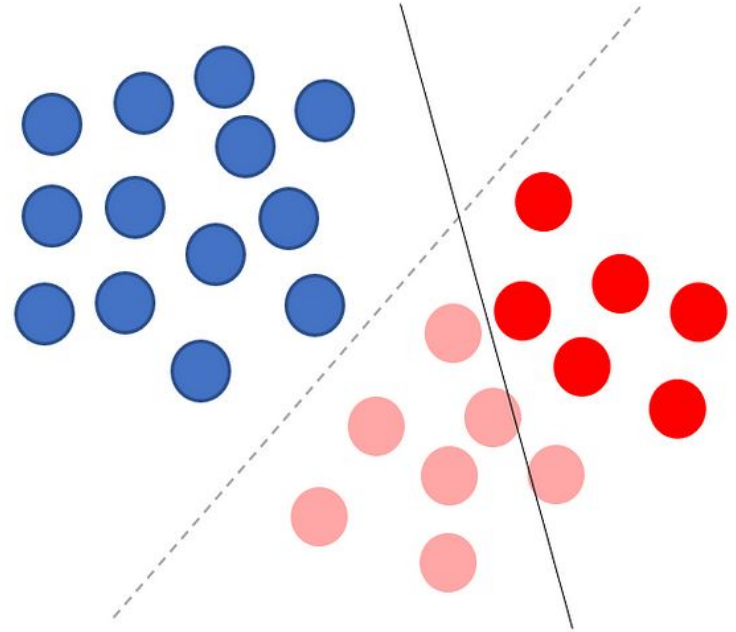
Class imbalance happens when one class has many more instances than the other(s)



The danger of class imbalance: unwarranted high accuracy



Classifier with balanced class



Classifier with imbalanced class

Class imbalance tends to skew the decision boundary of algorithms

How to deal with class imbalance

How can you deal with
class imbalance?

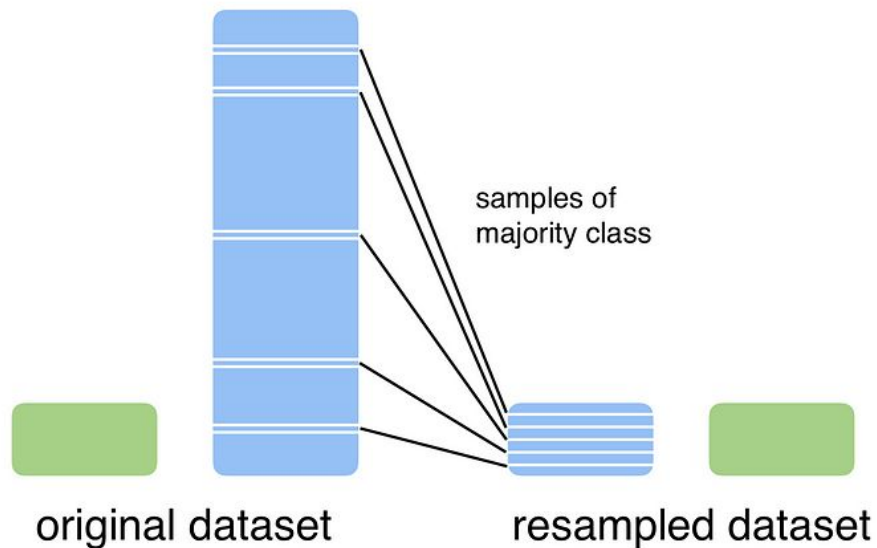


How can you deal with class imbalance?

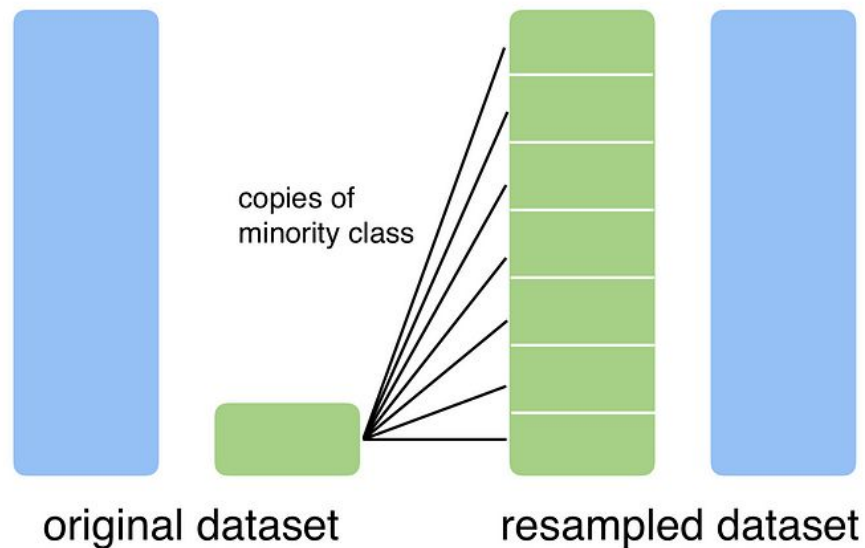
There are many methods to deal with class imbalance

- Undersampling your data
- Oversampling your data
- Generating artificial data
- Using imbalance-aware machine learning algorithms
 - ⇒ More on that in the ML course

Undersampling



Oversampling



Undersampling and oversampling

Undersampling

⇒ Removing data from the majority class

Addresses class imbalance

Reduces computational charge

Loss of information due to removing instances

Can introduce bias

Risk of underfitting when the imbalance is severe

Oversampling

⇒ Duplicating data from the minority class

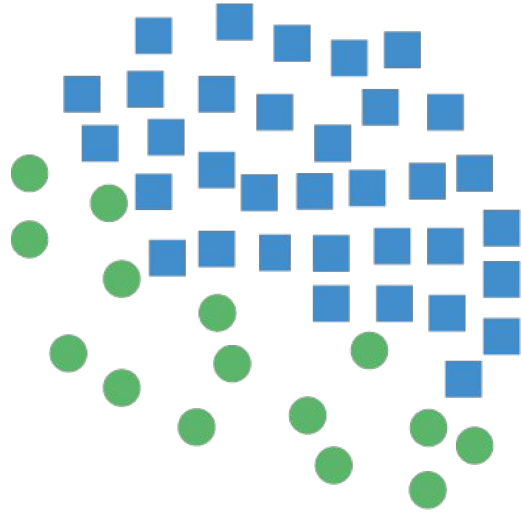
Addresses class imbalance

No loss of information

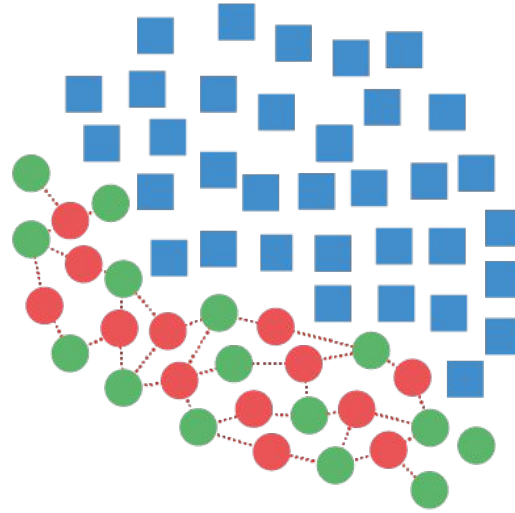
Risk of overfitting

May introduce noise from the minority class

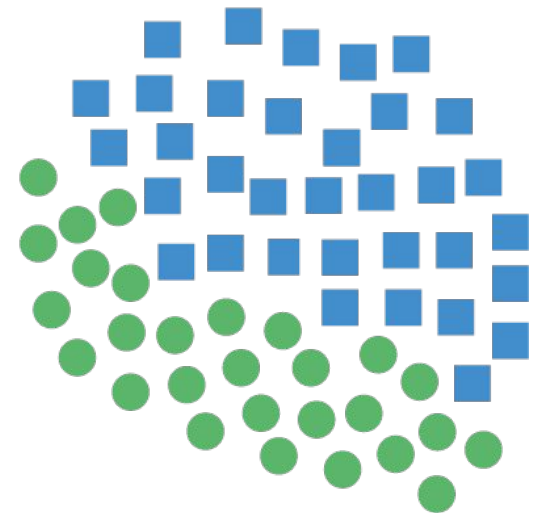
Synthetic Minority Oversampling Technique



Original Dataset



Generating Samples



Resampled Dataset

Synthetic Minority Oversampling Technique

Principle

- Choose a value for k
- For each instance in the minority class, identify the k nearest neighbours
- Interpolate new values linearly

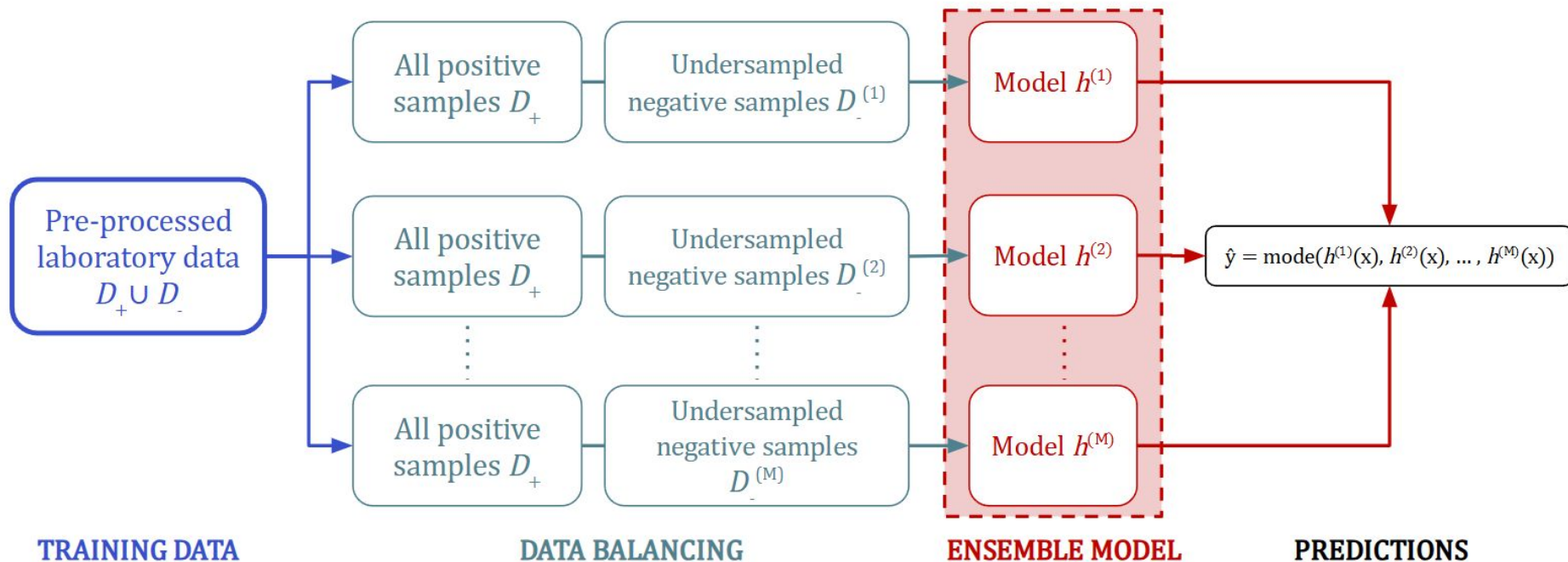
Variations

- ADASYN: Focuses on examples in low-density areas
- SMOTE-Tomek: Removes borderline noisy instances
- Borderline-SMOTE: Focuses on borderline instances



Generating artificial data with Generative Adversarial Networks (GAN)

[Image source](#)





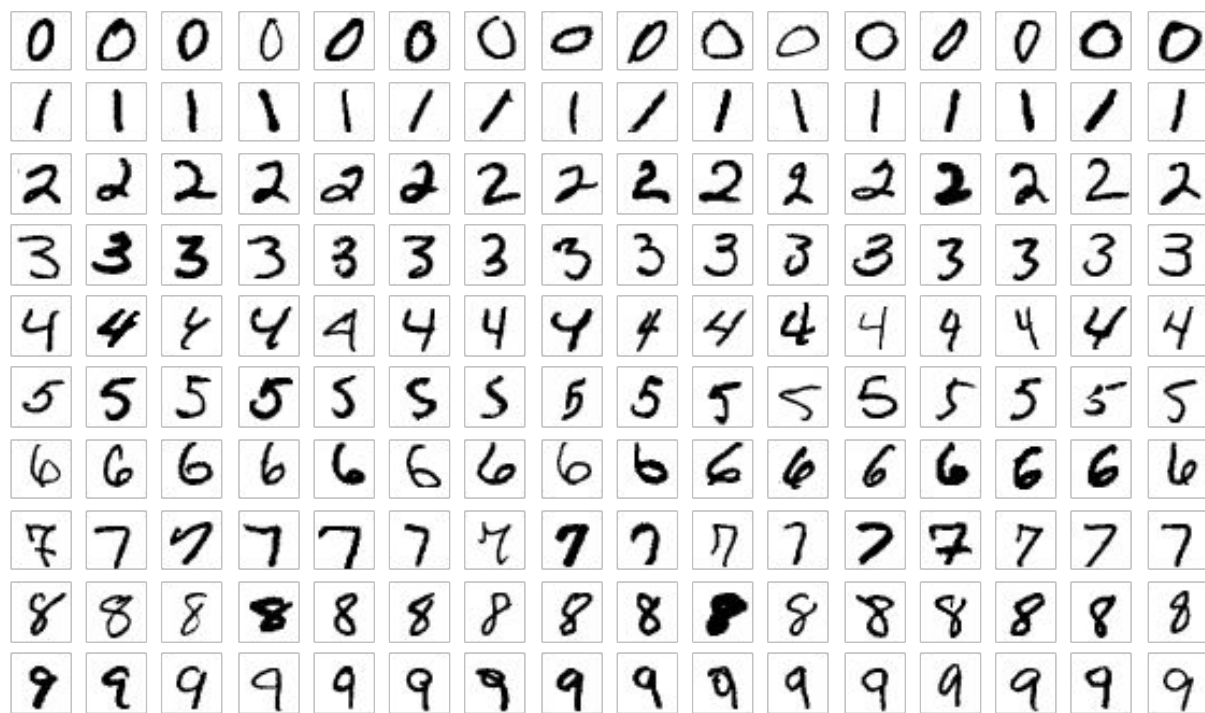
PERFORM RESAMPLING AFTER THE TRAIN-TEST SPLIT

Data leakage will artificially inflate your results

This is true for most pre-processing treatments: keep it in mind



Dimensionality reduction

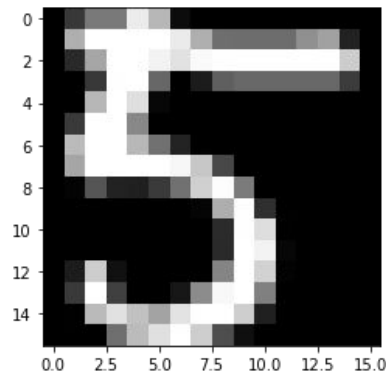


Example: the MNIST dataset

Example: the MNIST dataset

Modified National Institute of Standards and Technology database

- ❖ Database of hand-written digits
- ❖ Each digit is represented by the greyscale value of each pixel (between 0 and 255)
- ❖ Originally, the images are 28x28, but they will be 16x16 in today's practical

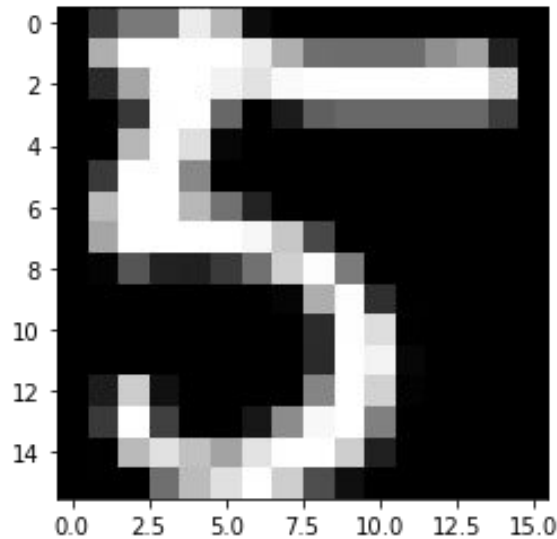


Example: the MNIST dataset

This example shows how the number of features can easily become very large.

Some datasets contain hundreds, or even thousands of features!

Imagine trying to train an algorithm with 4k photographs!



16 x 16 = 256 pixels

⇒ 256 dimensions

⇒ 256 features (columns)

Why would we want
to perform
dimensionality
reduction?



Why would we want to perform dimensionality reduction?

In many cases, having too many features is a disadvantage.

Dimensionality reduction allows for:

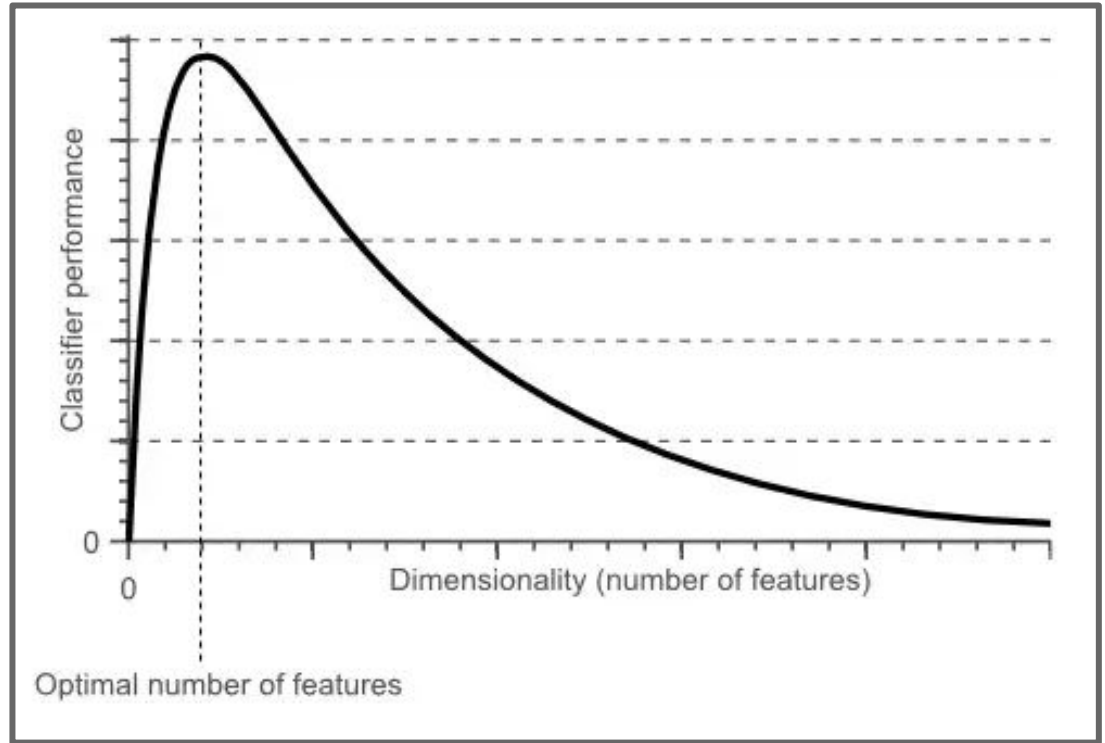
- The reduction of computational charge (i.e. reduction of computing time)
- The reduction of noise in the dataset
- The visualization of data in 2D / 3D
- The alleviation of the curse of dimensionality

⇒ Dimensionality reduction can help improve the performance of machine learning algorithms

Increasing the number of features can help up to a certain point.

However, when features are too numerous, it becomes difficult for algorithms to discernate patterns.

The effects of “high dimensionality” can appear with only 5 features!



How to perform dimensionality reduction

Method #1 : Feature Selection

How do you select
the most important
features?



How do you select the most important features?

There are many methods for feature selection

- Selection from expert knowledge (although it can be counter-productive)
- Deletion of low-variance variable
- Determining feature importance with a baseline machine learning model (e.g. Random Forests)
- Iterative choice using a machine learning model
 - Forward method: adding variables
 - Backward method: removing variables
 - Mixt method: doing both
 - Choosing randomly
- [Scikit-learn offers many methods for feature selection](#)

How to perform dimensionality reduction

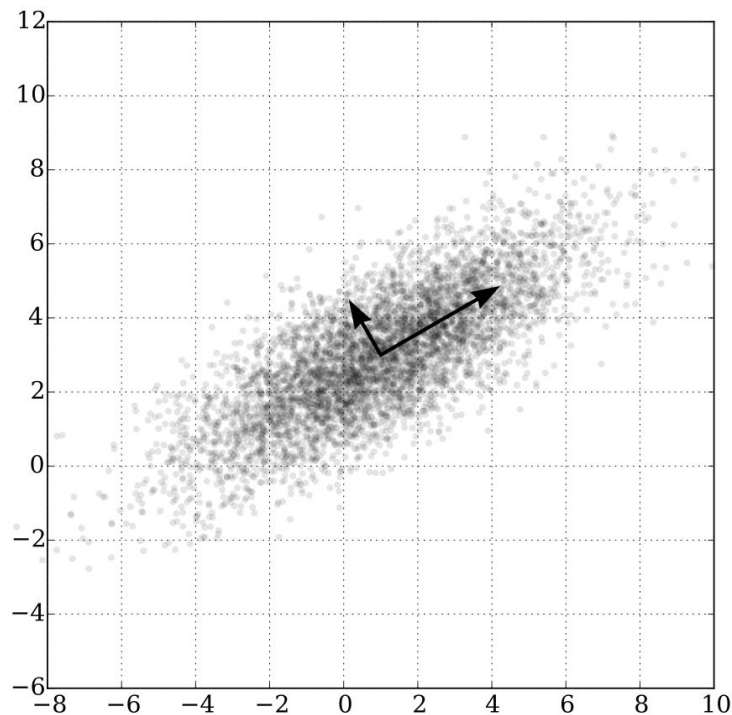
Method #2 : Feature extraction

Feature extraction

Using existing features to create new ones

The columns created in this way should be more significant
than the ones initially in the dataset

It can be simple linear combinations, or more complex functions



Intuitive principle

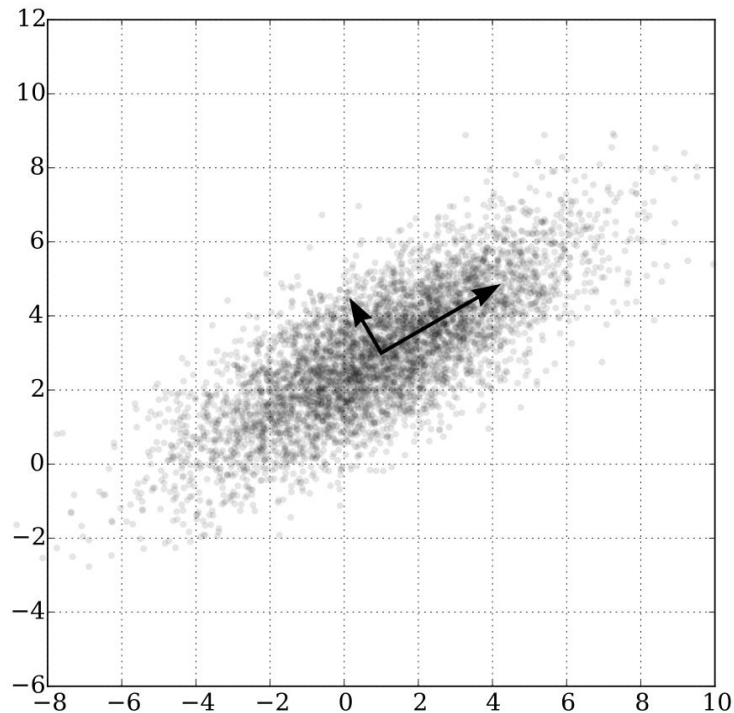
Finding the directions with the highest variance

Principle

We aim to find a **new basis** such that the variance of each projected component is maximized.

In other words, **PCA is not a dimension reduction method per se**. However, in constructing the new basis, we ensure that **the first vectors are the directions of the highest variance**.

To use PCA as a dimension reduction method, it is sufficient to **select the first N vectors** of the new basis.



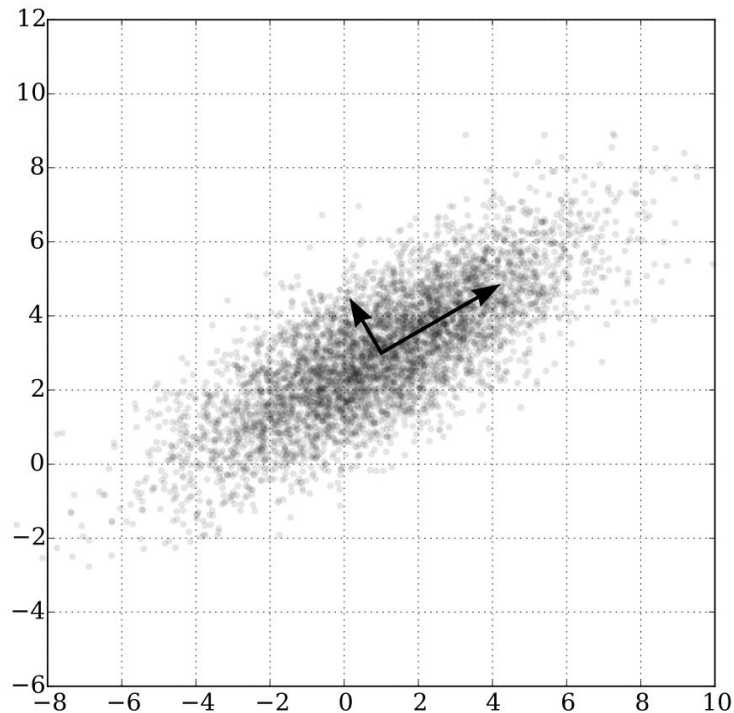
Formally

We are looking for a change of basis matrix that minimizes the approximation error.

$$U = \arg \min_{U^T U = 1} \sum_{n=1}^N \|x_n - \underbrace{UU^T x_n}_{x_{\text{approx}}}\|^2$$

It can be shown that this is equivalent to finding the **eigenvectors** of the covariance matrix.

$$\text{Cov}(X, Y) \equiv \text{E}[(X - \text{E}[X]) (Y - \text{E}[Y])]$$



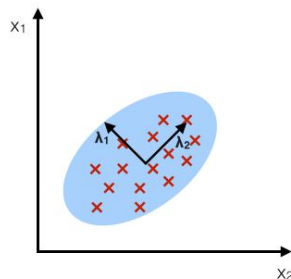
Using PCA

- **Standardization of data is necessary**
- Some information is lost despite keeping the most “important” dimensions
- Features cannot be interpreted anymore
- In a classification problem, PCA does not take interclass variance into account. Discriminant analysis takes both intra- and interclass variance into account.

⇒ The practical visually shows how PCA works with an example on MNIST

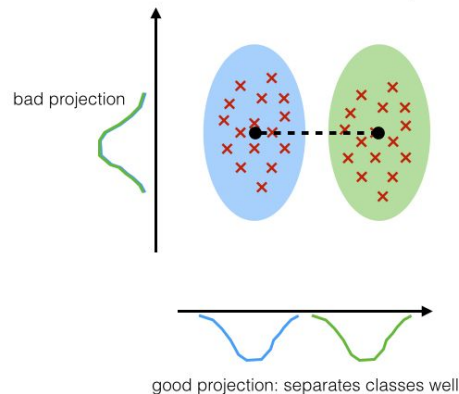
PCA:

component axes that maximize the variance



LDA:

maximizing the component axes for class-separation



[Image source](#)

Practical work

Get the latest version of the notebook from [GitHub](#)

Debrief

Debrief

What did we learn today?

What could we have done better?

What are we doing next time?

Data Science

Session 2 - Preparing data



hadrien.salem@centralelille.fr



[introduction-to-data-science](#)