

Modulo 7: Clustering

*Unsupervised Machine Learning
(Lab)*

Ing. Marco Pirrone Ph.D.

Ing. Massimo Romano Ph.D.



DIPARTIMENTO DI METODI E MODELLI PER
L'ECONOMIA IL TERRITORIO E LA FINANZA
MEMOTEF



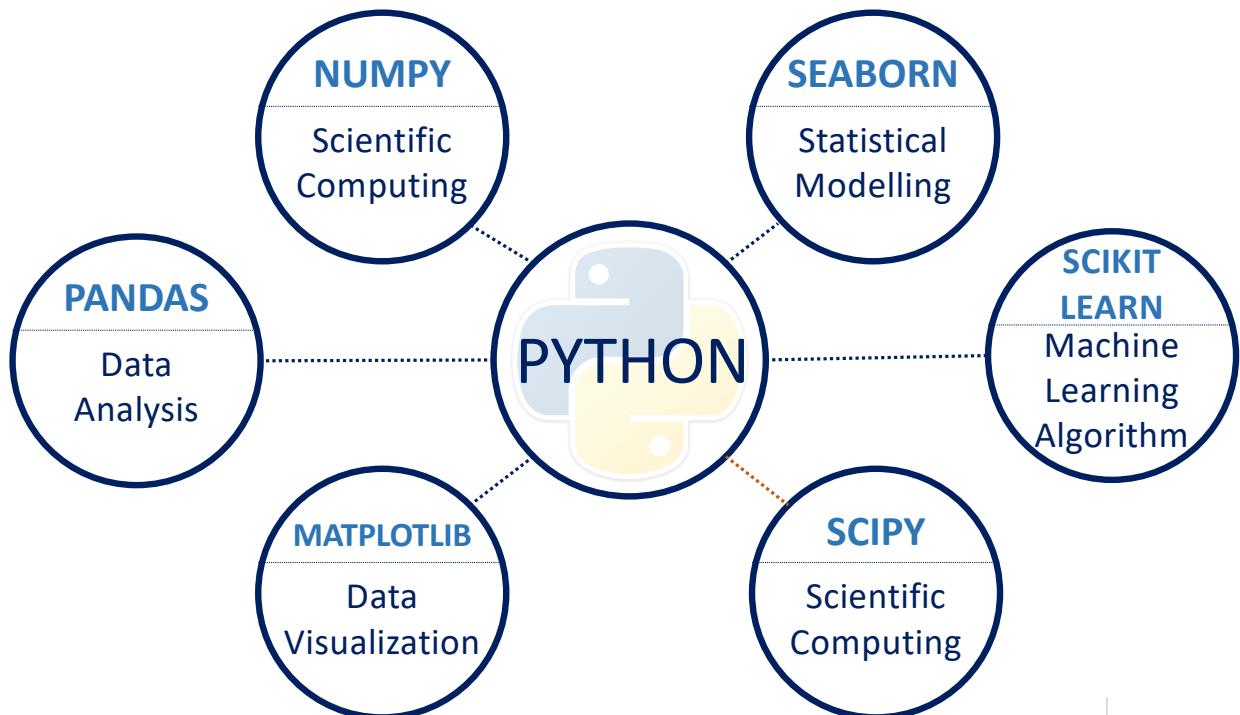
SAPIENZA
UNIVERSITÀ DI ROMA



A large, bright white rectangular area is positioned in the upper right quadrant of the slide, partially overlapping the abstract digital background. This area contains the word "Pillars" in a large, bold, black sans-serif font.

Pillars

Library



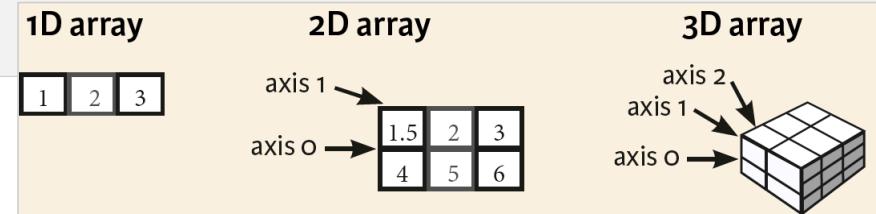
NumPy (1/2)

- **NumPy** (Numerical Python or Numeric Python), is an open source Python module which provides fast **mathematical computation** on **arrays** and **matrices**.
- Arrays and matrices are an essential part of the Machine Learning ecosystem

```
np.array(object, dtype=None, copy=True, order='K', subok=False, ndmin=0)
object:array_like - An array, any object exposing the array interface, an
object whose array method returns an array, or any (nested) sequence.

out:ndarray - An array object satisfying the specified requirements.
```

ndarray -> N-dimensional



DIPARTIMENTO DI METODI E MODELLI PER
L'ECONOMIA IL TERRITORIO E LA FINANZA
MEMOTEF



SAPIENZA
UNIVERSITÀ DI ROMA

NumPy (2/2)

Creating Arrays

```
>>> a = np.array([1,2,3])
>>> b = np.array([(1.5,2,3), (4,5,6)], dtype = float)
>>> c = np.array([(1.5,2,3), (4,5,6)], [(3,2,1), (4,5,6)]),
                           dtype = float)
```

Initial Placeholders

```
>>> np.zeros((3, 4))
>>> np.ones((2,3,4),dtype=np.int16)
>>> d = np.arange(10,25,5)

>>> np.linspace(0,2,9)

>>> e = np.full((2,2),7)
>>> f = np.eye(2)
>>> np.random.random((2,2))
>>> np.empty((3,2))
```

Create an array of zeros
Create an array of ones
Create an array of evenly spaced values (step value)
Create an array of evenly spaced values (number of samples)
Create a constant array
Create a 2X2 identity matrix
Create an array with random values
Create an empty array

Pandas (1/4)

- **Pandas** (PANel + DAta) is a library based on NumPy that allows easy and fast data analysis and **manipulation** by providing numerical **tables** and **time series** data structures called **DataFrame** and **Series**, respectively.
- **Series** -> One-dimensional **ndarray** with axis labels (including time series).

```
pd.Series(data=None,      index=None,      dtype=None,      name=None,      copy=False,  
          fastpath=False)  
  
  data: array-like, Iterable, dict, or scalar value  
  index: array-like or Index (1d) - Values must be hashable and have the same  
          length as data.  
  name: The name to give to the Series.
```

One-dimensional ndarray with axis labels (including time series).



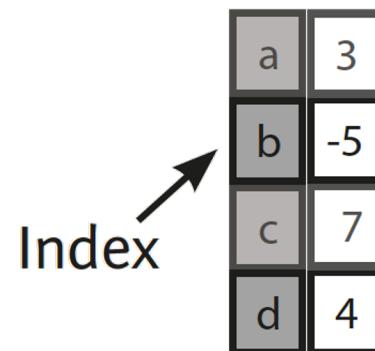
DIPARTIMENTO DI METODI E MODELLI PER
L'ECONOMIA IL TERRITORIO E LA FINANZA
MEMOTEF



Pandas (2/4)

Series

A one-dimensional labeled array
capable of holding any data type



a	3
b	-5
c	7
d	4

```
>>> s = pd.Series([3, -5, 7, 4], index=['a', 'b', 'c', 'd'])
```



Pandas (3/4)

- **DataFrame** -> Two-dimensional, size-mutable, potentially heterogeneous tabular data.

```
pd.DataFrame(data=None, index=None, columns=None, dtype=None, copy=False)
```

data: ndarray (structured or homogeneous), Iterable, dict, or DataFrame;
index: Index or array-like - Index to use for resulting frame;
columns: Index or array-like - Column labels to use for resulting frame;

Pandas (4/4)

DataFrame

	Country	Capital	Population
0	Belgium	Brussels	11190846
1	India	New Delhi	1303171035
2	Brazil	Brasília	207847528

A **two-dimensional** labeled data structure with columns of potentially different types

```
>>> data = {'Country': ['Belgium', 'India', 'Brazil'],
   'Capital': ['Brussels', 'New Delhi', 'Brasília'],
   'Population': [11190846, 1303171035, 207847528]}

>>> df = pd.DataFrame(data,
   columns=['Country', 'Capital', 'Population'])
```

Matplotlib (1/3)

- **Matplotlib** is a python library used to create 2D graphs and plots by using python scripts.
- It has a module named **pyplot** which makes things easy for plotting by providing feature to control line styles, font properties, formatting axes etc.
- Two type of interfaces for plotting
 - For simple plotting the **pyplot** module provides a MATLAB-like interface.
 - For the power user, you have full control of line styles, font properties, axes properties, etc, via an ***Object Oriented Interface***.

Matplotlib (2/3)

- MATLAB-like interface.

Import Pyplot

```
>>> import matplotlib.pyplot as plt
```

Plotting

```
>>> plt.plot(...)  
>>> plt.scatter(...)  
>>> plt.bar(...)
```

- Pandas has a built in **.plot()** function as part of the **DataFrame** class using matplotlib.



DIPARTIMENTO DI METODI E MODELLI PER
L'ECONOMIA IL TERRITORIO E LA FINANZA
MEMOTEF



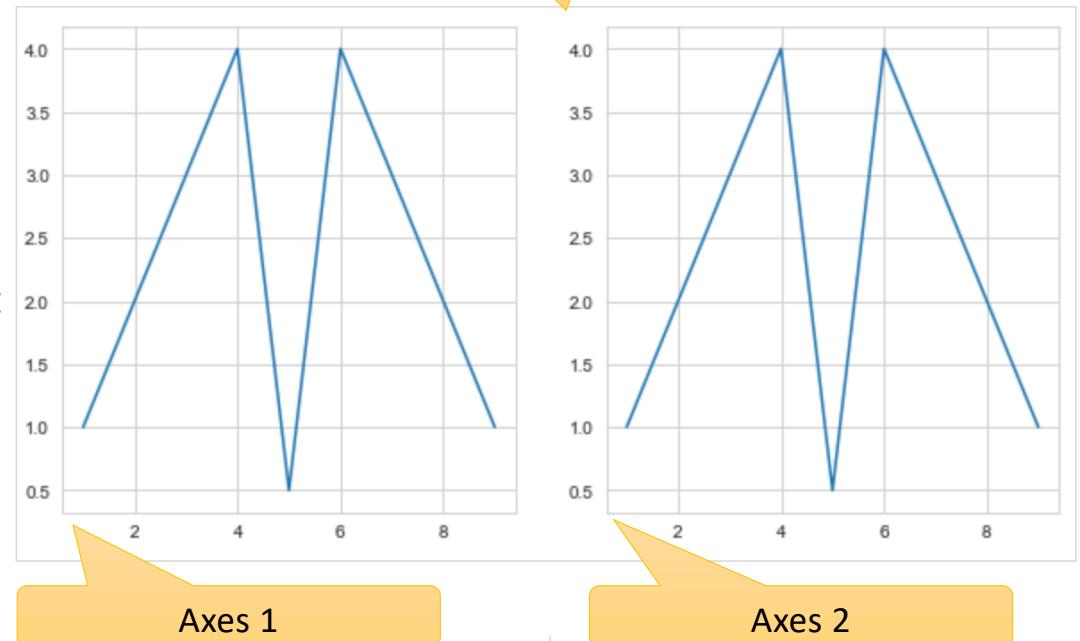
SAPIENZA
UNIVERSITÀ DI ROMA

Matplotlib (3/3)

- Object Oriented Interface.

- A Figure in matplotlib is divided into two different objects.
 - Figure object
 - Axes object
- A Figure object can contain one or more axes objects. One axes represents one plot inside figure.

Figure 1 row : 2 column



Axes 1

Axes 2

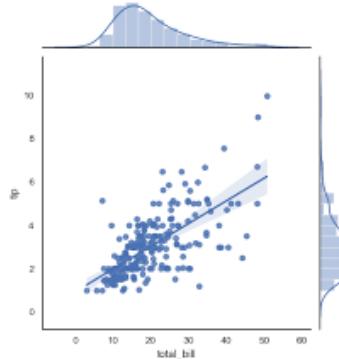
DIPARTIMENTO DI METODI E MODELLI PER
L'ECONOMIA IL TERRITORIO E LA FINANZA
MEMOTEF



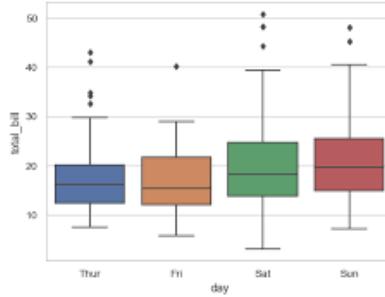
SAPIENZA
UNIVERSITÀ DI ROMA

Seaborn

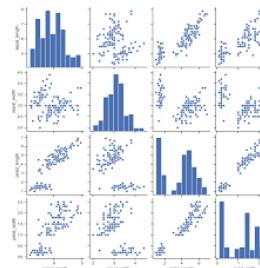
- **Seaborn** is a Python data visualization library based on Matplotlib. It provides a high-level interface for drawing attractive and informative statistical graphics.
- **Seaborn** is designed to work very well with Pandas DataFrame objects.



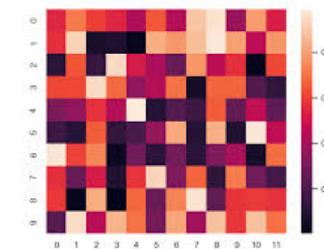
Jointplot



Boxplot



Pairplot

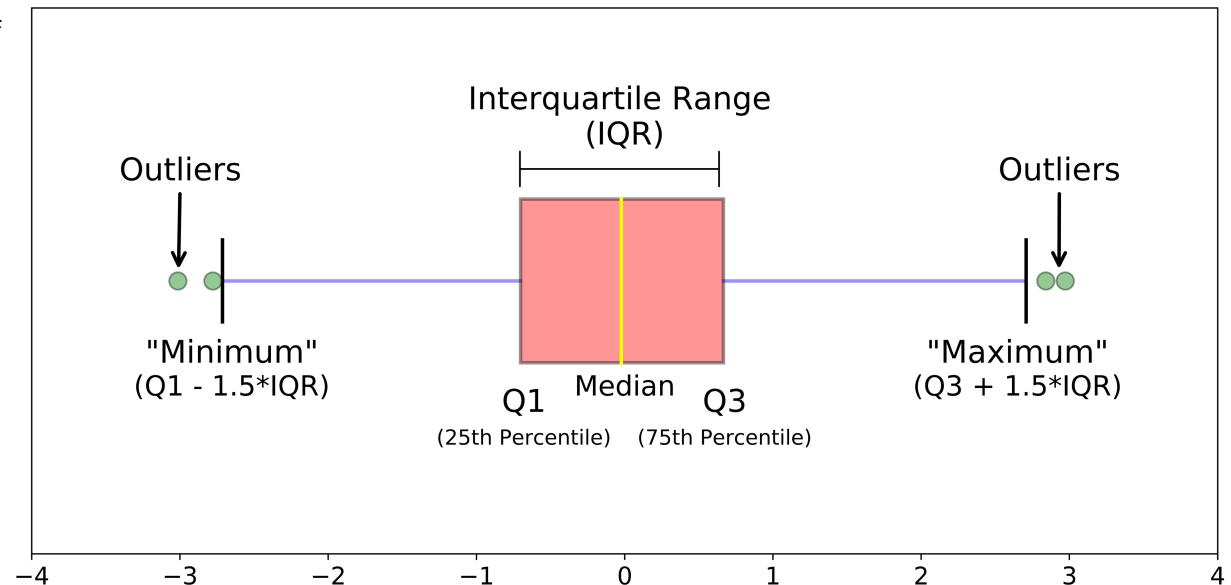


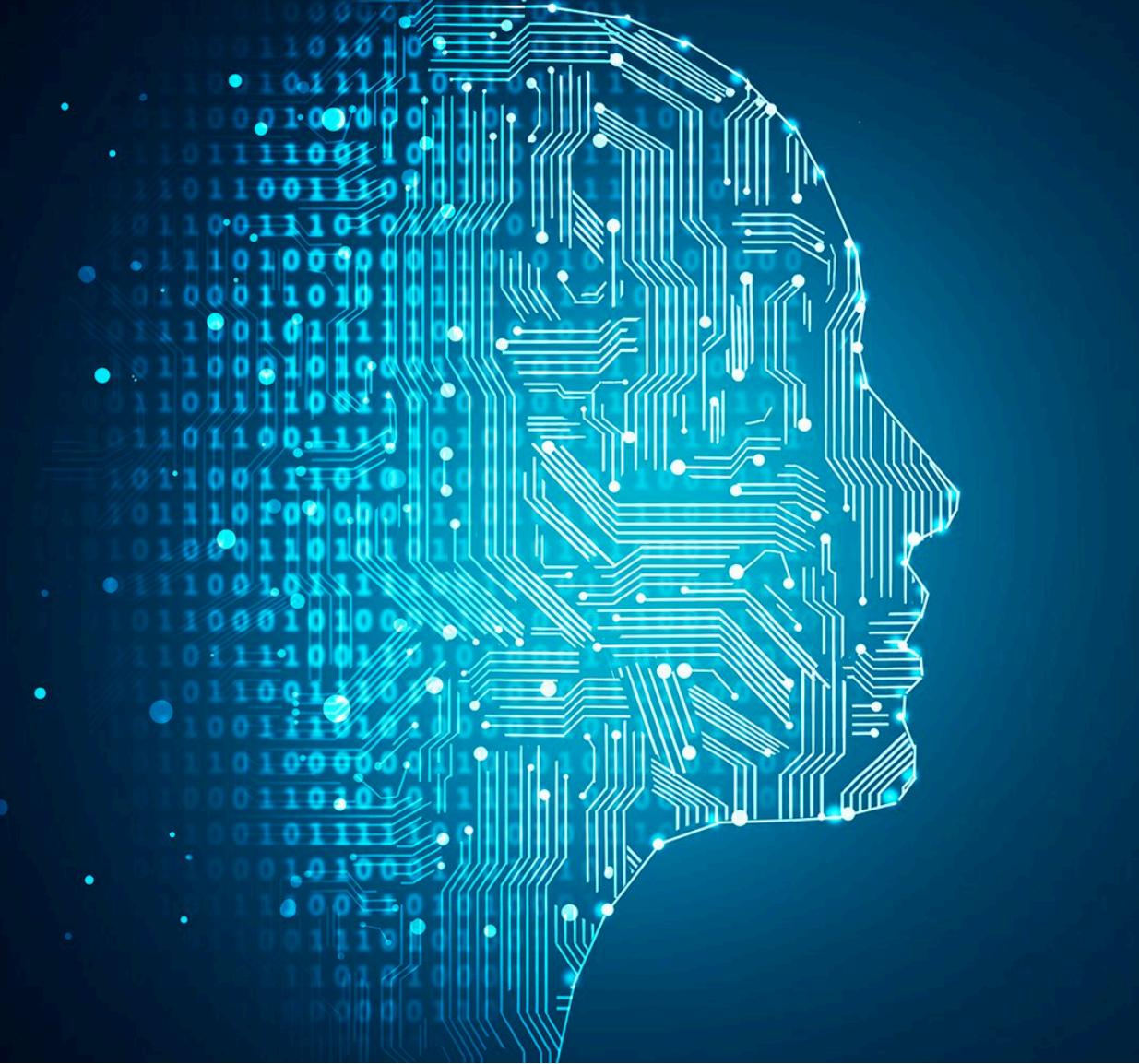
Heatmap

Boxplot

- A boxplot is a standardized way of displaying the distribution of data based on a five number summary (“minimum”, first quartile (Q1), median, third quartile (Q3), and “maximum”). It can tell you about your outliers and what their values are.

- median (Q2/50th Percentile)**: the middle value of the dataset.
- first quartile (Q1/25th Percentile)**: the middle number between the smallest number (not the “minimum”) and the median of the dataset.
- third quartile (Q3/75th Percentile)**: the middle value between the median and the highest value (not the “maximum”) of the dataset.
- interquartile range (IQR)**: 25th to the 75th percentile.
- whiskers (shown in blue)**
- outliers (shown as green circles)**
- “maximum”**: $Q3 + 1.5 \times IQR$
- “minimum”**: $Q1 - 1.5 \times IQR$





Scikit
Learn

Scikit-Learn



- **Scikit-Learn** is python's core machine learning package that has most of the necessary modules to support a basic machine learning project.
- It features various **classification**, **regression** and **clustering** algorithms including: support vector machines, random forests, gradient boosting, k-means and DBSCAN, etc.
- The package is written heavily in python, and it incorporates C++ libraries like LibSVM and LibLinear for support vector machines and generalized linear model implementation.
- The package depends on **Pandas** (mainly for the dataframe processes), **NumPy** (for the ndarray construct) and **SciPy** (for sparse matrices).



DIPARTIMENTO DI METODI E MODELLI PER
L'ECONOMIA IL TERRITORIO E LA FINANZA
MEMOTEF



<https://scikit-learn.org/>

Classification

Identifying to which category an object belongs to.

Applications: Spam detection, Image recognition.

Algorithms: SVM, nearest neighbors, random forest, ...

— Examples

Regression

Predicting a continuous-valued attribute associated with an object.

Applications: Drug response, Stock prices.

Algorithms: SVR, ridge regression, Lasso, ...

— Examples

Clustering

Automatic grouping of similar objects into sets.

Applications: Customer segmentation, Grouping experiment outcomes

Algorithms: k-Means, spectral clustering, mean-shift, ...

— Examples

Dimensionality reduction

Reducing the number of random variables to consider.

Applications: Visualization, Increased efficiency

Algorithms: PCA, feature selection, non-negative matrix factorization.

— Examples

Model selection

Comparing, validating and choosing parameters and models.

Goal: Improved accuracy via parameter tuning

Modules: grid search, cross validation, metrics.

— Examples

Preprocessing

Feature extraction and normalization.

Application: Transforming input data such as text for use with machine learning algorithms.

Modules: preprocessing, feature extraction.

— Examples



SAPIENZA
UNIVERSITÀ DI ROMA

DIPARTIMENTO DI METODI E MODELLI PER
L'ECONOMIA IL TERRITORIO E LA FINANZA
MEMOTEF



Preprocessing

Data Preparation

- Data **pre-processing** refers to the transformations applied to our data before feeding it to the algorithm.
 - Missing Value
 - Categorical Data
 - Feature Scaling



DIPARTIMENTO DI METODI E MODELLI PER
L'ECONOMIA IL TERRITORIO E LA FINANZA
MEMOTEF



SAPIENZA
UNIVERSITÀ DI ROMA

Missing Values (1/2)

- Very common to have data where values are missing for some of the attributes.
- Different ways of handling missing data have difference effects on the performance of our models.
- The approach to deal with missing values is heavily dependent on the nature of the dataset:
 - Ignore the data row.
 - Back-fill or forward-fill to propagate next or previous values respectively.
 - Replace with some constant value outside fixed value.
 - Replace with mean, median value.



Missing Values (2/2)

- Ignore the data row.
 - Preferred in cases where the percentage of missing values is relatively low (<5%).
 - You drop one whole observation just because one of the features had a missing value.
- Back-fill or forward-fill to propagate next or previous values respectively.
 - Note that the NaN value will remain even after forward filling or back filling if a next or previous value isn't available or it is also a NaN value.
- Replace with some constant value outside fixed value range.
 - This method is useful as it gives the possibility to group missing values as a separate category represented by a constant value.
 - The downside is that performance of linear models can suffer.
- Replace with mean, median value.
 - Treat every variable individually, ignoring any interrelationships with other variables.

Categorical Data (1/2)

- Categorical data are variables that contain label values rather than numeric values.
- Some algorithms can work with categorical data directly: **Decision Tree**.
- Many machine learning algorithms cannot operate on label data directly. They require all input variables and output variables to be numeric.
- This involves two steps:
 - Integer Encoding
 - One-Hot Encoding

Categorical Data (2/2)

INTEGER ENCODING

POSITION	ENCODE
First	1
Second	2
Third	3
Fourth	4
...	...

ONE-HOT ENCODING

COLOUR	IS_RED	IS_YELLOW	IS_GREEN
Red	1	0	0
Yellow	0	1	0
Green	0	0	1
...
...



Scaling (1/2)

- Data may contain attributes with a **mixtures of scales** for various quantities such as dollars, kilograms and sales volume.
- Many machine learning methods expect or are **more effective** if the data attributes have the same scale.
- Feature Scaling is a technique to **standardize** the independent features present in the data in a fixed range.
- If feature scaling is not done, then a machine learning algorithm **tends to weigh greater values**, higher and **consider smaller values as the lower values**, regardless of the unit of the values.

Scaling (2/2)

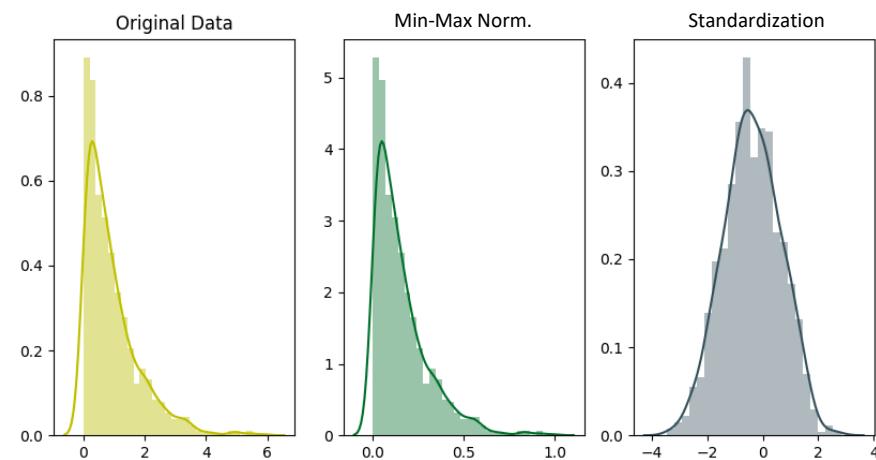
- Two popular data scaling methods are **normalization** and **standardization**.

Min-Max Normalization: rescales a feature or observation value with distribution value between 0 and 1.

$$X_{\text{new}} = \frac{X_i - \min(X)}{\max(x) - \min(X)}$$

Standardization: rescales a feature value so that it has distribution with 0 mean value and variance equals to 1.

$$X_{\text{new}} = \frac{X_i - X_{\text{mean}}}{\text{Standard Deviation}}$$



DIPARTIMENTO DI METODI E MODELLI PER
L'ECONOMIA IL TERRITORIO E LA FINANZA
MEMOTEF

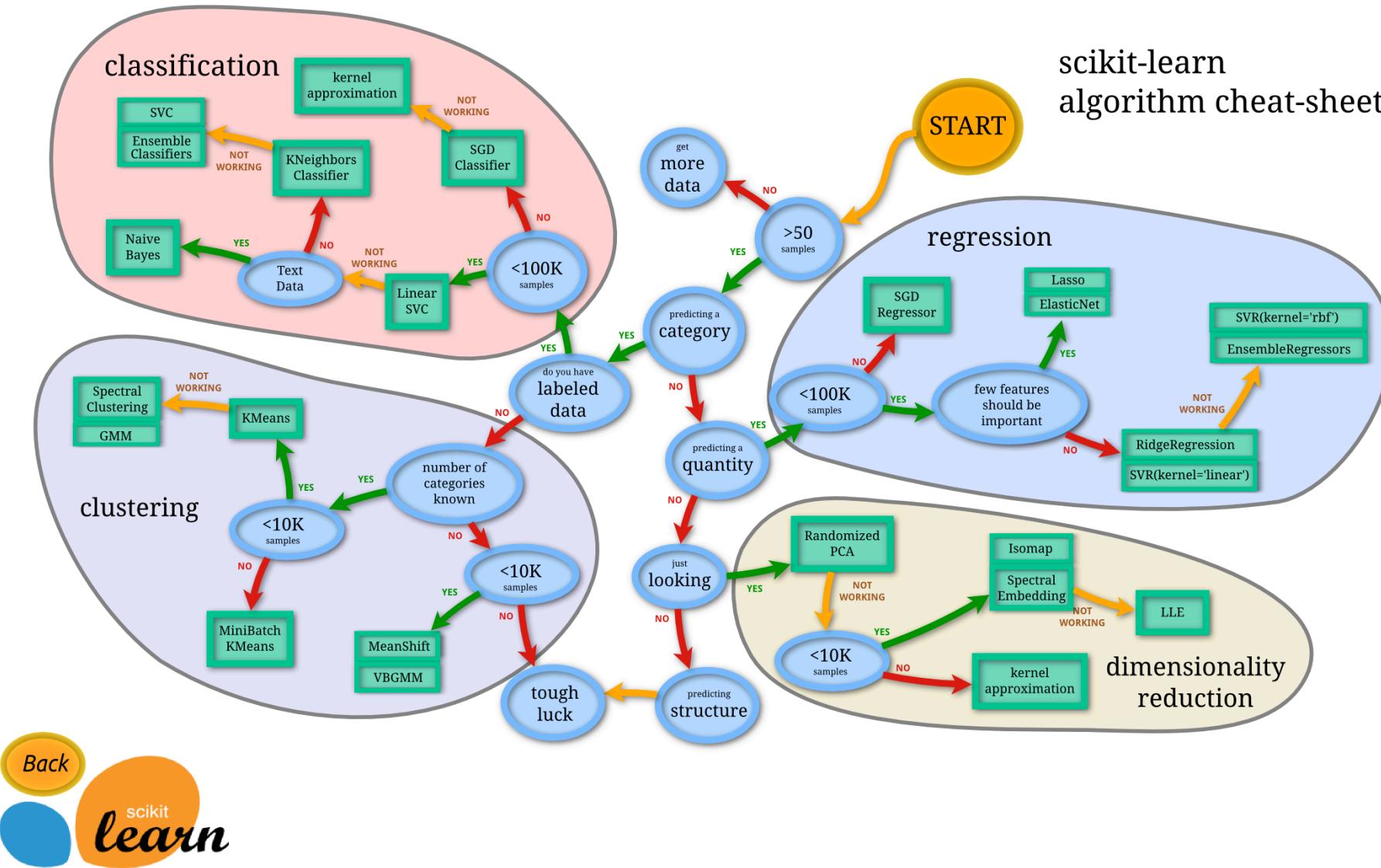


SAPIENZA
UNIVERSITÀ DI ROMA

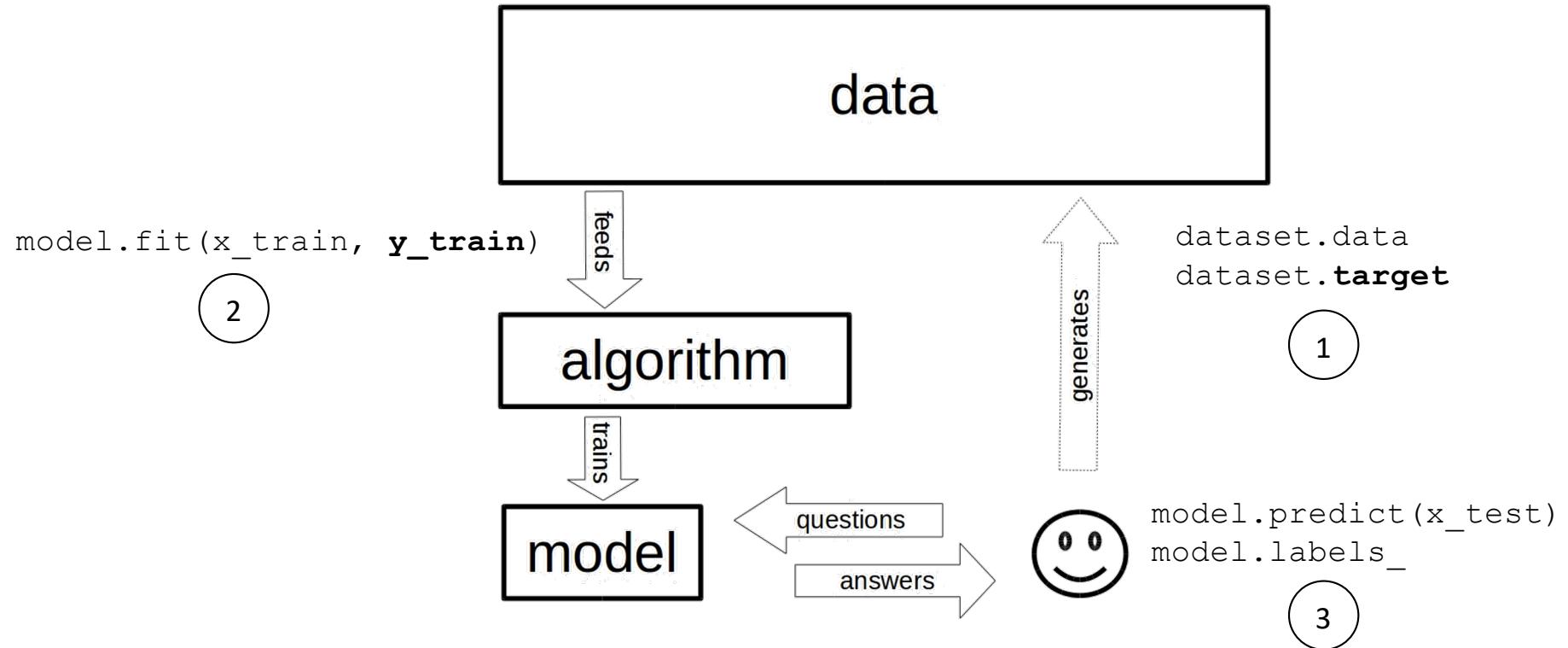


Clustering

scikit-learn algorithm cheat-sheet



Learning Process



Iris Dataset

- Number of Instances: 150 (50 in each of three classes)
- Number of Attributes: 4 numeric
 - sepal length in cm
 - sepal width in cm
 - petal length in cm
 - petal width in cm
- Class:
 - Iris Setosa
 - Iris Versicolour
 - Iris Virginica



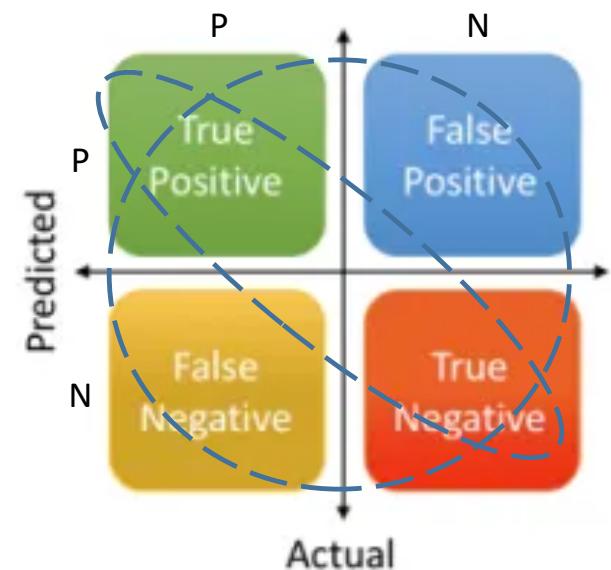
Evaluating Classification Algorithm (1/4)

Accuracy: is the most intuitive performance metric, it is simply the ratio of correctly predicted observations to total observations

$$\text{Accuracy} = \frac{\text{correctly predicted observation}}{\text{total observations}}$$

$$\text{Accuracy} = \frac{\text{TP} + \text{TN}}{\text{TP} + \text{FP} + \text{FN} + \text{TN}}$$

$$\text{Accuracy} = \frac{\text{True Positive} + \text{True Negative}}{\text{Total}}$$



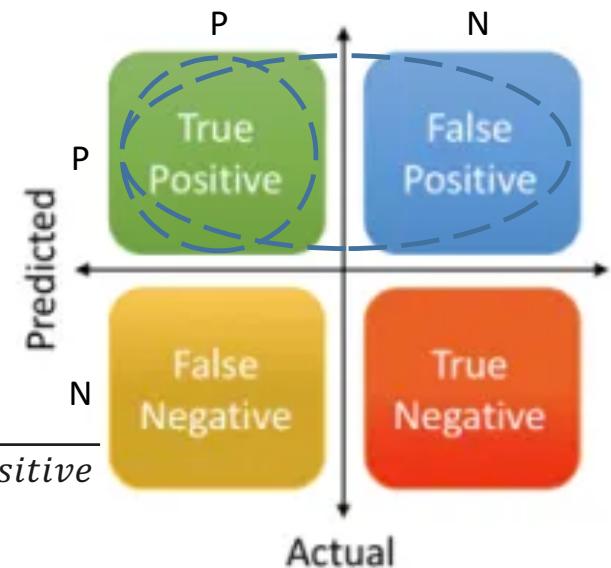
Evaluating Classification Algorithm (2/4)

Precision: is the ratio of correctly predicted ***positive*** observations of the total predicted ***positive*** observations.

$$\text{Precision} = \frac{\text{correctly predicted positive observation}}{\text{total predicted positive observations}}$$

$$\text{Precision} = \frac{\text{TP}}{\text{TP} + \text{FP}}$$

$$\text{Precision} = \frac{\text{True Positive}}{\text{Predicted Results}} \text{ or } \frac{\text{True Positive}}{\text{True Positive} + \text{False Positive}}$$



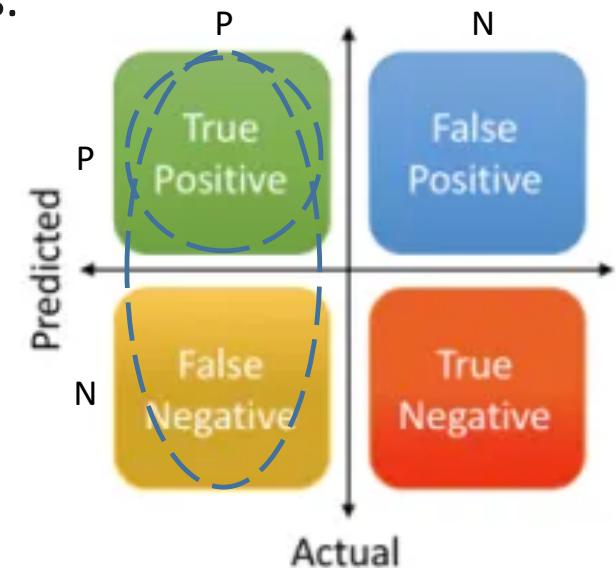
Evaluating Classification Algorithm (3/4)

Recall (Sensitivity): is the ratio of the **positive** observations correctly predicted to all observations in the actual class – yes.

$$Recal = \frac{\text{correctly predicted positive observation}}{\text{total actual positive observations}}$$

$$Recal = \frac{TP}{TP + FN}$$

$$Recal = \frac{\text{True Positive}}{\text{Actual Results}} \text{ or } \frac{\text{True Positive}}{\text{True Positive} + \text{False Negavite}}$$

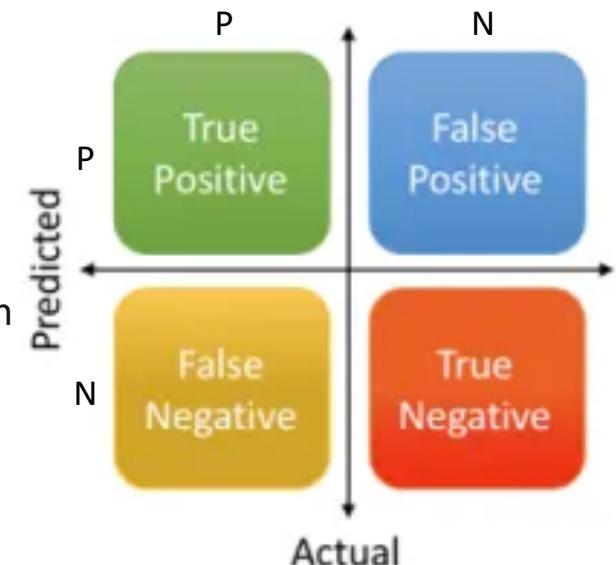


Evaluating Classification Algorithm (4/4)

F1 Score: is the harmonic mean of Precision and Recall.

$$F1\ score = 2 \frac{recall \times accuracy}{recall + accuracy}$$

- Accuracy is best if false positives and false negatives have similar costs.
- If the cost of false positives and false negatives is very different, then it's best to look at Precision and Recall at the same time.

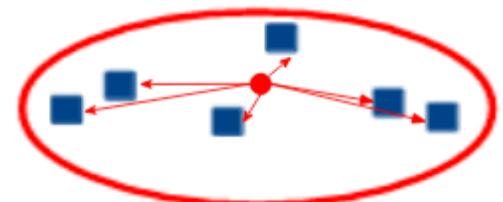


Evaluating Clustering Algorithm (1/5)

Inertia

- Within-cluster sum of squares of distances to the cluster center.
- It's a measure of how internally coherent clusters are.
- It suffers from various drawbacks:
 - Inertia makes the assumption that clusters are convex and isotropic.
 - Inertia is not a normalized metric, we just know that lower values are better and zero is optimal.

$$\sum_{i=0}^n \min_{\mu_j \in C} (\|x_i - \mu_j\|^2)$$

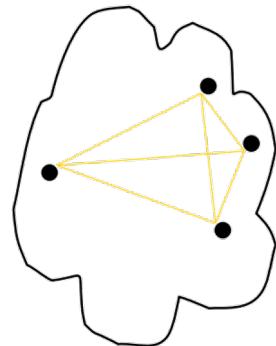


Evaluating Clustering Algorithm (2/5)

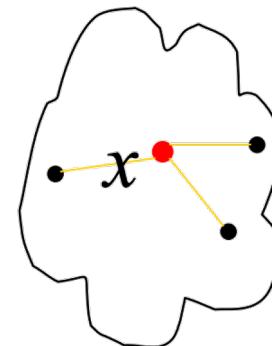
Cohesion

- Measures how closely related are points in a cluster.

$$Cohesion(C) = \sum_{x,y \in C} proximity(x,y)$$



$$Cohesion(C) = \sum_{x \in C} proximity(x, centroid)$$

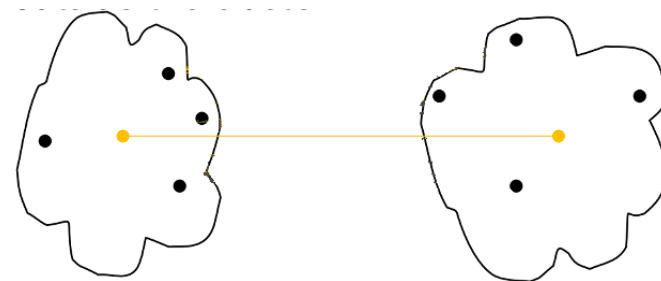
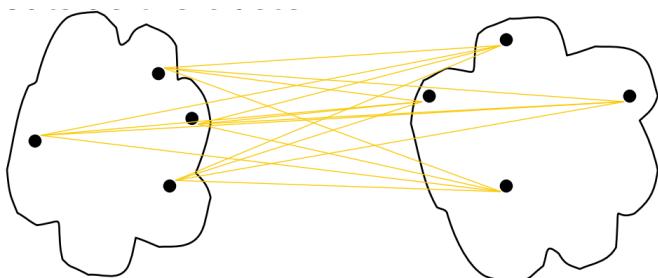


Evaluating Clustering Algorithm (3/5)

Separation

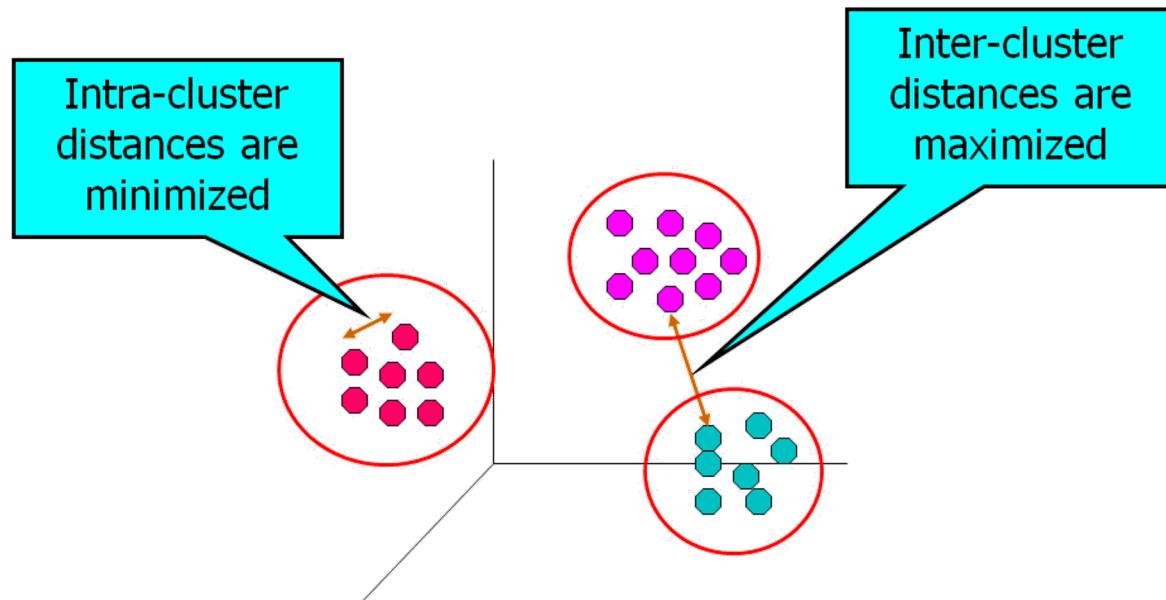
- Measure how distinct or well separated a cluster is from other clusters.

$$\text{Separation}(C_i, C_j) = \sum_{\substack{x \in C_i \\ y \in C_j}} \text{proximity}(x, y) \quad \text{Separation}(C_i, C_j) = \text{proximity}(\text{centroid}_i, \text{centroid}_j)$$



Evaluating Clustering Algorithm (4/5)

- Cohesion vs Separation



DIPARTIMENTO DI METODI E MODELLI PER
L'ECONOMIA IL TERRITORIO E LA FINANZA
MEMOTEF



SAPIENZA
UNIVERSITÀ DI ROMA

mands
masterandskills

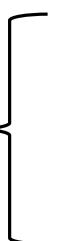
Evaluating Clustering Algorithm (3/5)

Silhouette Coefficient

Combination of the mean cohesion and separation.

$$\text{Silhouette} = \frac{b - a}{\max(a, b)} \quad a = \text{cohesion} \quad b = \text{separation}$$

Normalization factor

Range from [-1, 1] 

- 1: poor clustering
- 0: indifferent
- +1: excellent clustering

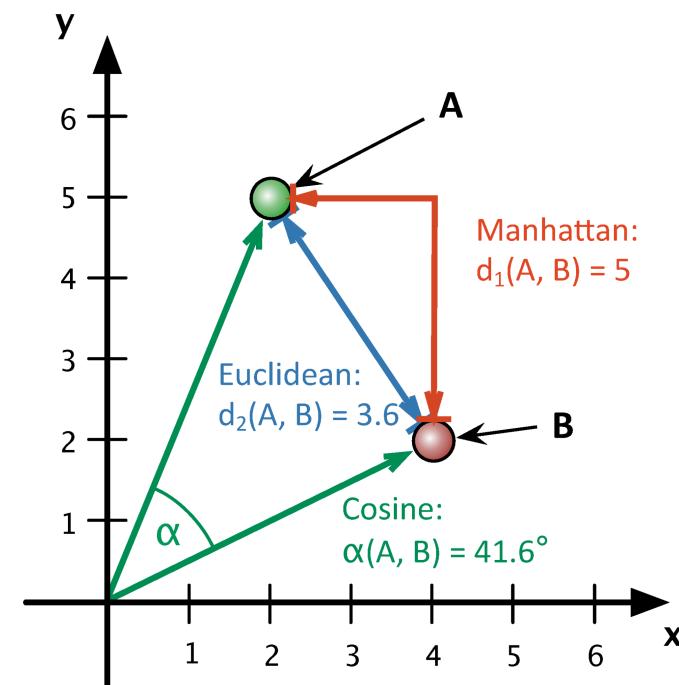
Agglomerative Clustering (1/2)

Affinity

$$\text{Manhattan} \Rightarrow d_1(A, B) = \sum_{i=1}^n |a_i - b_i|$$

$$\text{Euclidean} \Rightarrow d_2(A, B) = \sqrt{\sum_{i=1}^n (a_i - b_i)^2}$$

$$\text{Cosine} \Rightarrow \cos(\alpha) = \frac{A \cdot B}{\|A\| \|B\|} = \frac{\sum_{i=1}^n a_i b_i}{\sum_{i=1}^n a_i^2 \sum_{i=1}^n b_i^2}$$



DIPARTIMENTO DI METODI E MODELLI PER
L'ECONOMIA IL TERRITORIO E LA FINANZA
MEMOTEF



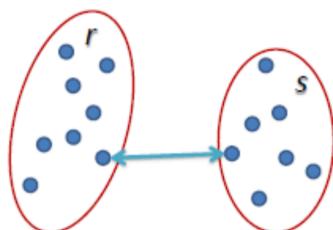
SAPIENZA
UNIVERSITÀ DI ROMA

Agglomerative Clustering (2/2)

Linkage

Single

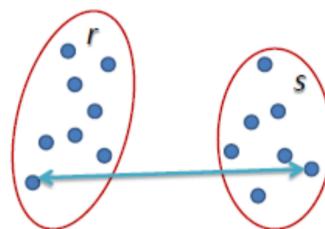
The distance between two clusters is the shortest distance between two points in each cluster.



$$L(r,s) = \min(D(x_{ri}, x_{sj}))$$

Complete

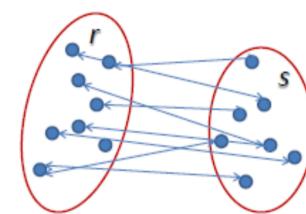
The distance between two clusters is the longest distance between two points in each cluster.



$$L(r,s) = \max(D(x_{ri}, x_{sj}))$$

Average

The distance between clusters is the average distance between each point in one cluster to every point in other cluster.



$$L(r,s) = \frac{1}{n_r n_s} \sum_{i=1}^{n_r} \sum_{j=1}^{n_s} D(x_{ri}, x_{sj})$$

DIPARTIMENTO DI METODI E MODELLI PER
L'ECONOMIA IL TERRITORIO E LA FINANZA
MEMOTEF



SAPIENZA
UNIVERSITÀ DI ROMA

Reference

- Pandas (<http://pandas.pydata.org>)
- Scikit-Learn (<https://scikit-learn.org>)
- NumPy (<http://www.numpy.org/>)
- Matplotlib (<https://matplotlib.org/>)
- Seaborn (<https://seaborn.pydata.org/>)
- Jupyter (<https://jupyter.org/>)

