

# 18. Clustering

Corso di Python per il Calcolo Scientifico

# Outline

- Il clustering
- Tipi di clustering
- Workflow del clustering
- L'algoritmo k-means
- Valutazione della bontà del clustering

# Il clustering

- Prevede la suddivisione dei campioni nei dataset **senza che questi abbiano un'etichetta a priori**.
- Ha varie applicazioni, come ad esempio:
  - segmentazione del mercato;
  - individuazione di aree coerenti all'interno di un'immagine;
  - suddivisione delle stelle sulla base delle caratteristiche di magnitudine;
  - definizione delle feature mancanti in un dataset (anche supervisionato);
  - raggruppamento dei film proposti da Netflix.
- Ogni cluster è contraddistinto da un **identificativo**.
  - Può essere usato come ingresso ad un altro algoritmo di machine learning (magari supervisionato).

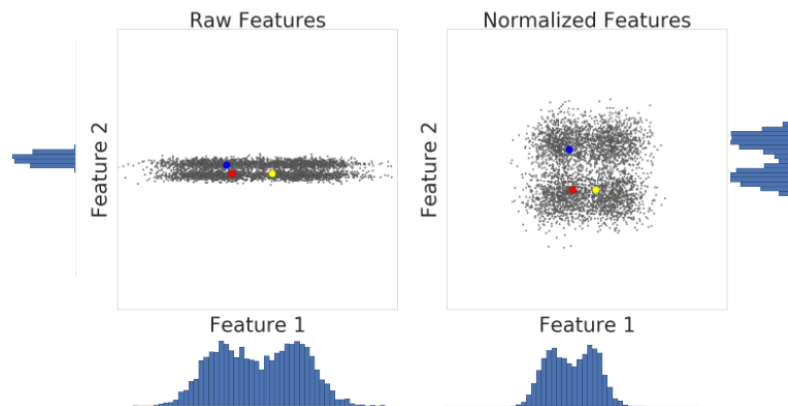
# Tipi di clustering

- Ogni tipo di algoritmo di clustering ha una diversa applicazione e complessità computazionale.

Tipo di clustering	Descrizione
<b>Centroid – based</b>	Dati organizzati in base alla distanza da un centroide. Efficienti, ma sensibili a condizioni iniziali e presenza di outliers.
<b>Density – based</b>	Dati organizzati in base alla densità. Efficaci nel caso di cluster ad alta densità e per l'outlier detection.
<b>Distribution – based</b>	Dati organizzati secondo la distribuzione, supposta gaussiana. Efficaci soltanto se la tesi di distribuzione gaussiana risulta essere corretta.
<b>Hierarchical</b>	Dati organizzati secondo un albero gerarchico, che può essere tagliato per ridurre il numero complessivo di cluster. Efficaci nel caso di dati di un certo tipo, come le tassonomie.

# Workflow del clustering (1)

- Gli algoritmi di clustering prevedono un workflow, esattamente come quelli di machine learning visti finora.

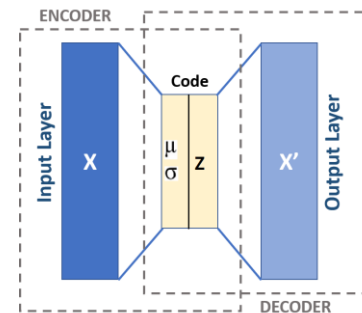


# Workflow del clustering (2)

- Per definire una metrica ci sono due possibilità
- Nel **primo caso**, ci possiamo affidare ad una semplice combinazione di due/tre feature del nostro dato
- Nel secondo caso, dobbiamo usare un **embedding**, ovvero una rappresentazione ridotta di un dato ad alta dimensionalità

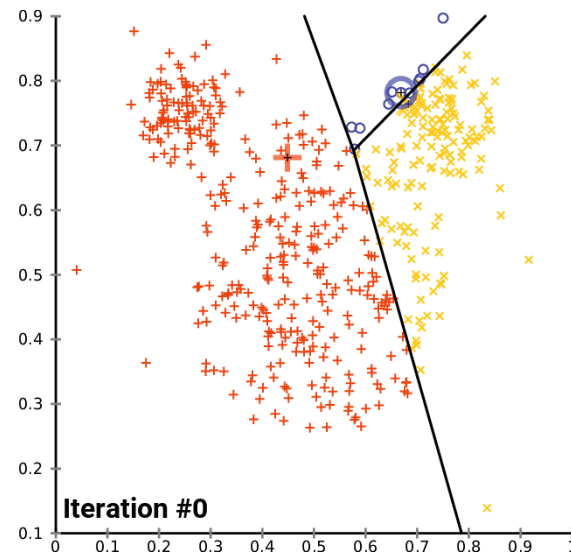


$$d = \sqrt{(x_1 - x_2)^2}$$



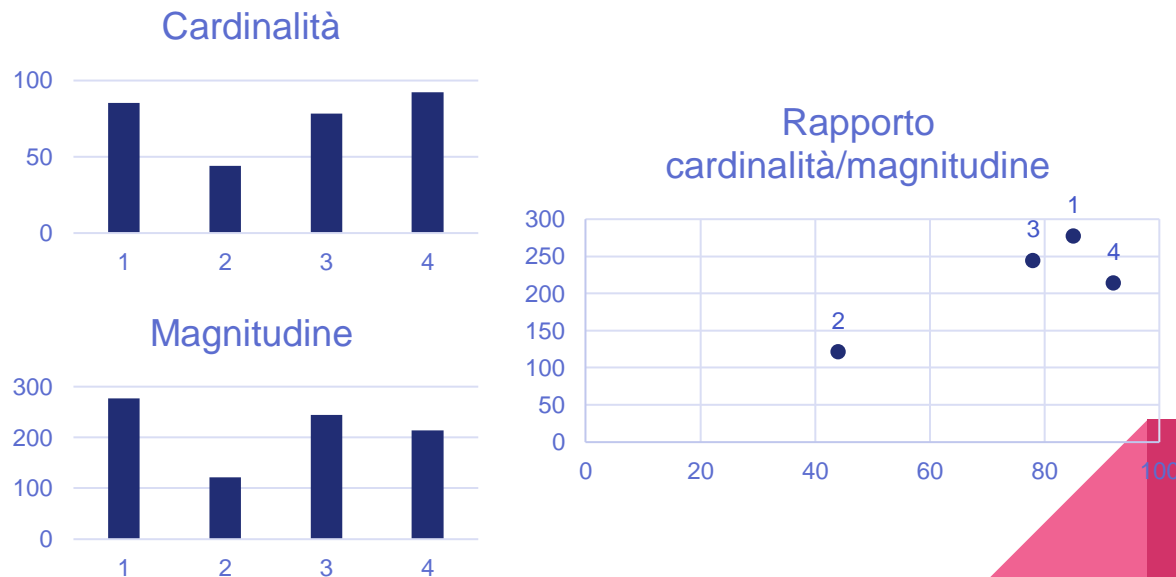
# L'algoritmo K-means

- Il **k-means** è un algoritmo iterativo
- Prevede l'assegnazione a priori del numero di cluster (il valore **k**)
- **Primo step**: determinare i centroidi
- **Secondo step**: calcolare la distanza dai centroidi
- **Terzo step**: aggiornare i centroidi, e ripetere dal secondo step fino a che non si arriva a convergenza
- Implementato in Scikit Learn grazie alla classe `KMeans()`



# Valutazione della bontà del clustering (1)

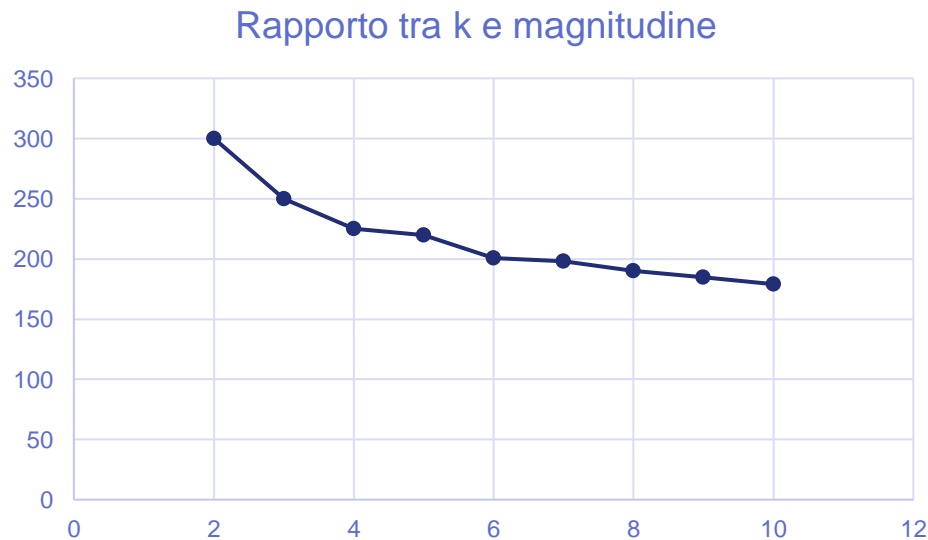
- Valutiamo il rapporto tra **cardinalità** e **magnitudine** per valutare empiricamente la qualità del clustering





# Valutazione della bontà del clustering (2)

- Valutiamo il rapporto tra **k** e **magnitudine** per stabilire un numero *ottimo* di cluster



# Domande?

42