

## Lecture : Clustering

Utrecht Summer school on Machine Learning with Python July 2024

## Types of learning



## Why unsupervised learning?

- Clustering: expect groups in our data, but were not able to measure them
  - potential new subtypes of cancer tissue
  - groups of shoppers based on their purchase histories
- We want to summarize features into one feature to use in further decisions/analysis
  - subgrouping customers by their spending types
- Informative way to visualize high-dimensional data
  - Reduce dimensionaltiy of the data to visualise them in a 2-D plot

#### Supervised learning VS Unsupervised





Training set: (x1, y1), (x2,y2)..., (xm,ym)

Training set: (x1, x2..., xm)

#### Supervised learning VS Clustering





Training set: (x1, y1), (x2,y2)..., (xm,ym)

Training set: (x1, x2..., xm) No target/output

#### Clustering

Find subgroups (clusters) if there are similar examples in the dataset

## Applications of clustering

- Recommender systems
  - Cluster users with similar viewing habits
- Medical imaging
  - Cluster images to find patients with similar medical imaging (MRI, X-ray etc)
- Market segmentation
  - Divide consumers based on their spending hapits

## k-Means algorithm

#### K-means

- One of the most popular clustering algorithms
- K-means is an iterative algorithm
- Each data point belongs to the cluster with the nearest mean



#### K-means

• Step 1: Initialization: Choose the number of clusters K and initialize the centroids randomly

• K = 4



#### K-means - Step 2

• Step 2: Assign each data point to the nearest centroid based on a distance metric (e.g., Euclidean distance).



#### K-means - Step 3

• Step 3: Calculate the new centroids as the mean of all points in the cluster.



K-means - Step 4



• Step 4: Repeat the assignment and update clusters until the centroids do not change significantly or a maximum number of iterations is reached

#### K - means

- Input
- K (number of clusters)
- Observations (x1, x2, x3, ..., xm)
- Labels are not required

### K-means algorithm

• Randomly initialize K cluster centroids  $\mu_1, \mu_2, \cdots, \mu_K \in \mathbb{R}^n$ 

Repeat{

for i = 1 to m  $c^{(i)} \coloneqq index$  (from 1 to K) of cluster centroid closest to  $x^{(i)}$ **Cluster assignment step** 

for k = 1 to K  $\mu_k :=$  average (mean) of points assigned to cluster k**Centroid update step** 

#### Random initialization

The initialization of the centroids is random, and this can give different results on the same dataset.



#### Random initialization

- The random initialization can lead us to local optimal solution
- To avoid this inconsistency, we can run it 100 times and select the best performance among these 100
- Performance?

#### Random initialization

- Which of those is the best > measure inertia (or also called distortion cost function)
- The mean square distance between each point and the centroid of its cluster
- The lower the value of distortion the better the solution

#### Choose the number of clusters

- Choose them manually with visualization
- But it is not always easy
- No clear answer



#### Choose k with elbow method

- Run K-means for different values of k
- For each k, calculate the distortion cost function
- Create a plot with the number of clusters k on the x-axis and the distortion on the y-axis.
- Distortion typically decreases as k increases
  - Adding more clusters reduces the distance from each point to its assigned centroid

## Elbow method

- Distortion drops quickly when we increase the number of clusters until a certain point, and then it slows down and decreases more and more slowly.
- Identify the "Elbow": where the rate of decrease sharply slows down.
- The optimal k is typically at the elbow; adding more clusters beyond this point provides does not reduce distortion significantly



#### Silhouette score

- Graphs are not always clear. Sometimes the curve decreases very smoothly.
- In this case, you should try the **Silhouette score**.
- The silhouette value is a measure of how similar an object is to its own cluster (cohesion) compared to other clusters (separation).
- The silhouette ranges from -1 to +1, where a high value indicates that the object is well matched to its own cluster and poorly matched to neighboring clusters.

#### Silhouette score



- To calculate the silhouette score for the whole dataset, you take the mean of silhouette scores over all the instances.
- The higher the score, the better, and it does not constantly decrease as inertia.
- Number of clusters-> the one with the highest silhouette score.

#### K-means in Python

kmeans = KMeans(n\_clusters=2, n\_init= 10)
kmeans.fit(df\_scaled)
print("Centroids of the clusters:", kmeans.cluster\_centers\_)





#### Exercise

- Number of observations = 6
- Number of clusters = 2
- Distance metric: Euclidean distance

Step-1: Choose random K points and set as cluster centers C1 = (2,2) C2 = (3, 3)

Question: Which are going to be the centroids after one iteration?

C1 = (?, ?) C2 = (?, ?)

No	X	Y
1	1	1
2	2	3
3	1	2
4	3	3
5	2	2
6	3	1

#### Exercise

- Number of observations = 6
- Number of clusters = 2
- Distance metric: Euclidean distance

Step-1: Choose random K points and set as cluster centers C1 = (2,2) C2 = (3, 3)

Question: Which are going to be the centroids after one iteration? C1 = (1.75, 1.5) C2 = (2.5, 3)

No	X	Y
1	1	1
2	2	3
3	1	2
4	3	3
5	2	2
6	3	1

### Evaluation of clusters

How to evaluate clustering results

- 1. Use of external information
- 2. Visual exploration
- 3. Stability assessment / sensitivity analysis

#### External validation

- 1. External validation
- Are the clusters associated with external feature *Y*?
- "Making unsupervised supervised"
- Examples:
- Are my customer segments based on spending associated with the demographics of the customers?

#### Visual exploration

- Problem: Kind of hard to see already...
- Wait till you get 1000 variables!
- New idea: Reduce variables into 2D "manifold" for visualization
- Popular techniques: PCA, t-SNE, Discriminant Coordinates

## Cluster stability

- Three "stabilities"
- How much does clustering change when:
- 1. Changing some hyperparameters (distance metric, K, ...)
- 2. Changing some observations (bootstrapping, Hennig, 2007)
- 3. Changing some features

Check if observations are classified into same cluster across choices



#### 10 minutes break

## Hierarchical Clustering

## Hierarchical clustering



- Builds a hierarchy of clusters
- A hierarchy might be more natural to the type of data
- Different users might care about different levels of granularity

## Hierarchical clustering

#### • Top-down (divisive)

- Partition data into 2-groups (e.g., 2-means)
- Recursively cluster each group

- Start with every point in its own cluster.
- Repeatedly merge the "closest" two clusters
- Different definitions of "closest" give different algorithms.

#### Hierarchical clustering



- Compute the distance matrix between the input data points
- Let each cluster be a point
- Repeat:
  - Merge the two closest clusters
  - Update the distance matrix
- Until only one cluster remains

- Compute the distance matrix between the input data points
- Let each cluster be a point
- Repeat:
  - Merge the two closest clusters
  - Update the distance matrix
- Until only one cluster remains
- Distance matrix: Manhattan, Euclidean
- But we need a distance measure for the clusters (for the merging step).

## Single linkage

- Have a distance measure on pairs of objects.
- d(x, y): Distance between x and y
- Single linkage:  $dist(A, B) = \min_{x \in A, x' \in B} d(x, x')$



- Have a distance measure on pairs of objects.
- d(x, y): Distance between x and y
- Complete linkage:  $dist(A, B) = \max_{x \in A, x' \in B} d(x, x')$



- Have a distance measure on pairs of objects.
- d(x, y): Distance between x and y
- Average linkage: dist(A, B) = average d(x, x')



- Have a distance measure on pairs of objects.
- d(x, y): Distance between x and y
- Centroid linkage: dist(A, B) = d(mean(A), mean(B))











#### k-means VS Hierarchical clustering

- k-means clustering is faster and simpler, but requires choosing the number of clusters beforehand and may not capture complex structures
- Hierarchical clustering is more flexible and intuitive, but can be computationally expensive and sensitive to outliers

## Dimensionality Reduction

#### Plot high dimensional data

- Suppose you have data with many dimensions
- How are you going to plot those data?

0r	iginal Data:					
	mean radius	mean texture	mean perimeter	mean area me	an smoothness \	
0	17.99	10.38	122.8	1001.0	0.11840	
1	20.57	17.77	132.9	1326.0	0.08474	
2	19.69	21.25	130.0	1203.0	0.10960	
	mean compact	ness mean con	cavity mean co	ncave points m	ean symmetry \	
0	0.27	7760	0.3001	0.14710	0.2419	
1	0.07	7864	0.0869	0.07017	0.1812	
2	0.15	5990	0.1974	0.12790	0.2069	
	mean fractal	dimension	<ul> <li>worst radius</li> </ul>	worst texture	worst perimeter	\
0		0.07871	. 25.38	17.33	184.6	
1		0.05667	. 24.99	23.41	158.8	
2		0.05999	. 23.57	25.53	152.5	
	worst area w	vorst smoothne	ss worst compa	ctness worst c	oncavity \	
0	2019.0	0.16	22	0.6656	0.7119	
1	1956.0	0.12	38	0.1866	0.2416	

## Dimensionality reduction

- Simplifies models
- Reduces computational cost
- Helps in visualizing high-dimensional data

#### Principal Component Analysis

- Identifies directions (principal components) that maximize variance
- Projects data onto these new directions (PC components) to reduce dimensions while retaining most information of the data

0r	iginal Data:							
	mean radius	mean texture	mean pe	rimeter	mean area	mean smoot	hness	\
0	17.99	10.38		122.8	1001.0	0.	11840	
1	20.57	17.77		132.9	1326.0	0.	08474	
2	19.69	21.25		130.0	1203.0	0.	10960	
	mean compact	ness mean cor	ncavity i	mean con	cave points	s mean symm	netry	\
0	0.2	7760	0.3001		0.14710	) 0.	2419	
1	0.0	7864	0.0869		0.07017	<i>.</i> 0.	1812	
2	0.1	5990	0.1974		0.12790	) 0.	2069	
	mean fractal	dimension	. worst	radius	worst text	ure worst	perime	ter
0		0.07871		25.38	17	.33	18	4.6
1		0.05667		24.99	23	8.41	15	8.8
2		0.05999		23.57	25	<b>.</b> 53	15	2.5
	worst area	worst smoothne	ess wors	t compac	tness wors	st concavity	/ \	
0	2019.0	0.16	522	0	.6656	0.7119	)	
1	1956.0	0.12	238	0	.1866	0.2416	<b>;</b>	



#### Standarize each feature



Find the line that goes through the origins and:

- 1. minimizes the sum of the distances from the points to the line
- 2. Or maximises the distances of the projected points on the line to the origin





#### Not linear regression

PCA







The line is called PC1

- PC1 has a slope of 0.25
- For every 1 unit of x1 axis, we go up 0,25 unit of x2 axis (green arrows)
- Data are mostly spread out on x1 axis
- PC1 is a linear combination of x1 and x2
- -> to make PC1, mix 1 unit1 of x1 and 0,25 units of x2



PCA

- $-a^2 = b^2 + c^2 -> a = 1,03$
- PC1 is scaled so that its length is 1, so we can divide by 1,03
  - a = 1,
  - x1 = 1/ 1,03 = 0,97,
  - x2 = 0,25/ 1,03 = 0,3

The one unit vector that consists of 0,97 of x1 and 0,3 of x2 is called **singular vector** or **eigenvector** for PC1

The proportions are called **loading scores** 



- PC2 is the line that is perpendicular to PC1 and goes through the origin
- Loading scores for PC2:
  - 0.97 for x2
  - - 0.3 for x1
- In terms of how the values are projected in PC2, x2 is 4 times as important as x1



#### Final PCA plot

• Project the points on PC1 and PC2 and rotatate to PC1 is horizontal (x-axis) and PC2 is the y-axis



#### Final PCA plot

• Use the projected points to find the new coordinates



## Loadings

- How much each variable contributes to each principal component (correlation between the original variables and the principal components)
- High absolute values: the original variable strongly influences the principal component.



## Loadings

- Positive values: the variable and the PC are positively correlated
- Negative values: the variable and the PC are negatively correlated.
- A negative loading indicates that its absence contributes to some degree to the principal component



## Summary

#### Unsupervised learning

- Advantages:
  - Requires less manual data preparation (i.e., no hand labeling) than supervised machine learning.
  - Capable of finding previously unknown patterns in data, which is impossible with supervised machine learning models.
- Disadvantages:
  - Results may be unpredictable or difficult to understand.
  - Difficult to measure accuracy or effectiveness due to lack of predefined answers during training.

#### Conclusion: clustering

- Clustering looks for "similar" groups of observations
- k-means is simple but needs to predefine k
- Hierarchical Clustering: no need to predefine the number of clusters. Doesn't work well on vast amounts of data or huge datasets.
- PCA will reduce the dimensionality of your data into principal components. Very popular way to visualise high-dimensional data into a 2-d plot

# Practical 7