

LECTURER: Nghia Duong-Trung

MACHINE LEARNING

Introduction to Machine Learning

Clustering

Regression

Support Vector Machines

Decision Trees

Genetic Algorithm

1

2

3

4

5

6

UNIT 2

CLUSTERING



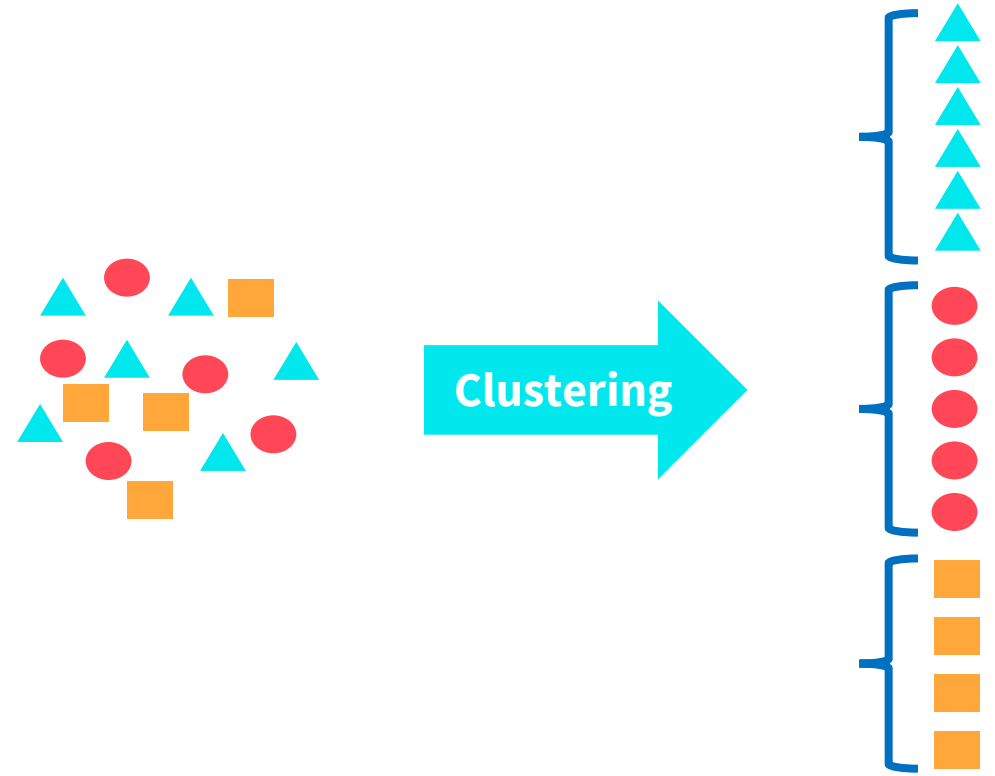
- Know the definitions and terms used for clustering
- Comprehend common applications of clustering analysis
- Understand different methods for clustering analysis
- Analyse the advantages and limitations of the clustering methods
- Implement clustering methods in Python



1. What is clustering and its application?
2. What are the popular methods of clustering, their advantages and limitations?
3. How to implement the clustering methods in python?

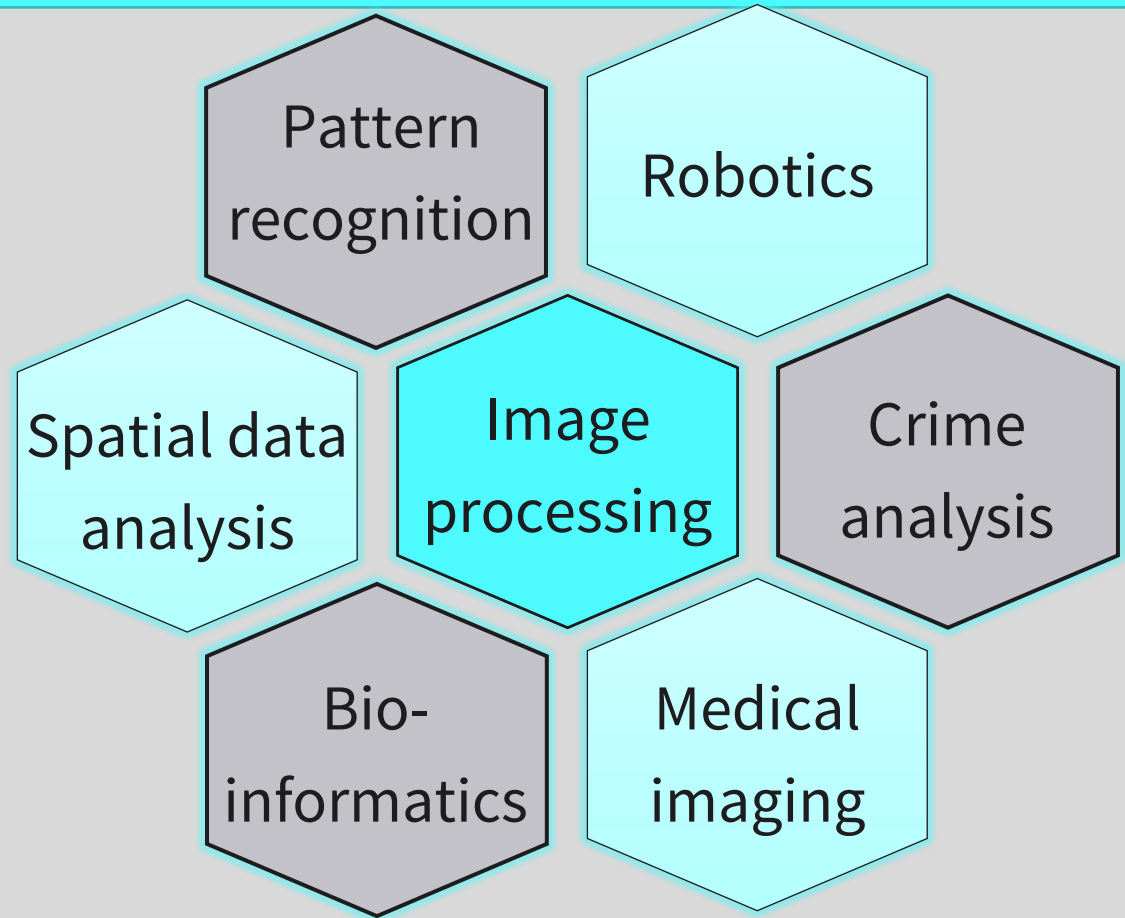
INTRODUCTION

- Clustering is an **unsupervised** learning technique to:
 - reveal meaningful **partitions, hierarchies**
 - find association **rules** of data
 - gather data points into **groups** and extract **useful information**

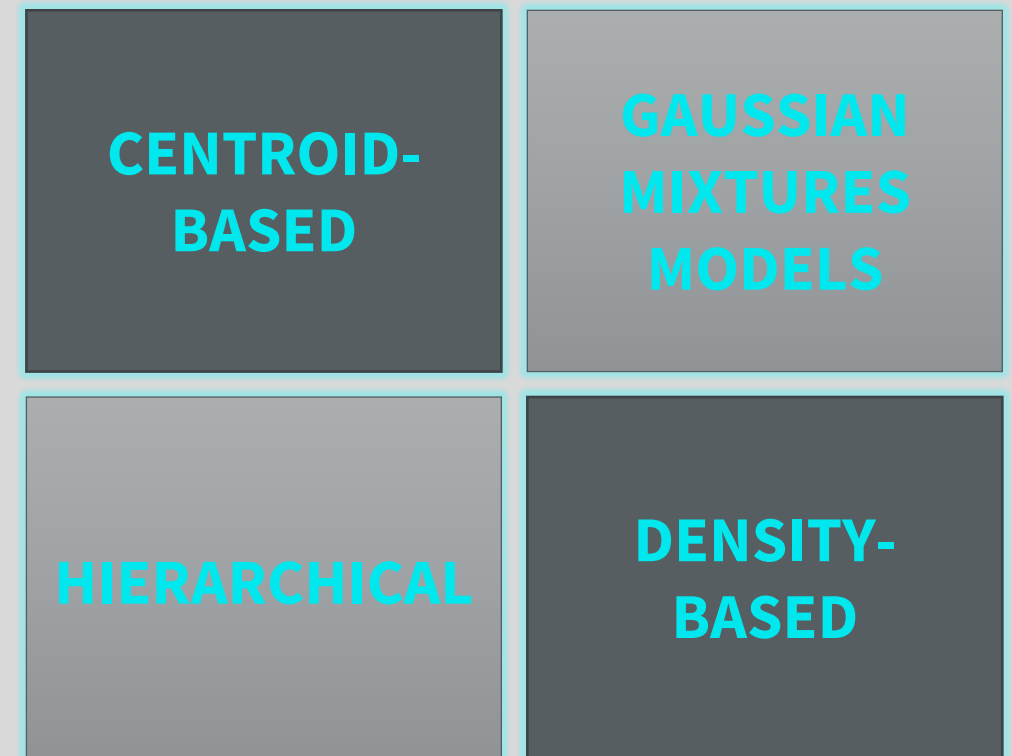


An example of clustering

Clustering applications

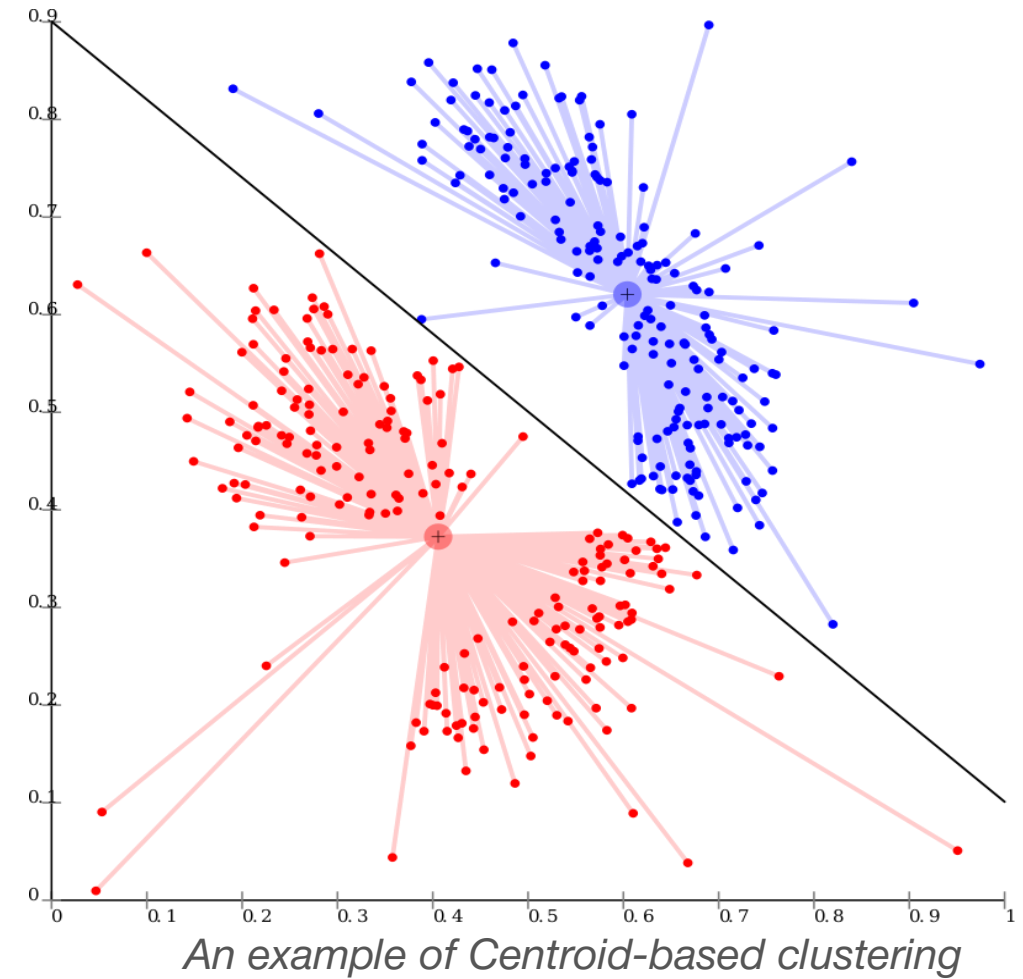


Clustering methods



2.1 CENTROID-BASED CLUSTERING

- **Centroid**: the arithmetic average position of all the points
- Centroid-based clustering searches for a **pre-determined** number of clusters
- Each data point is assigned to the cluster achieving the **minimum distance** from the centroid



DISTANCE: WHO IS CLOSE TO WHOM

- Distance between two numbers (1-D): For real numbers a and b , the distance is the absolute difference: $d(a, b) = |a - b|$
- Example: $d(3, 10) = |3 - 10| = 7$.

- Euclidean distance (2-D, 3-D, n-D)

2-D (points (x_1, y_1) , (x_2, y_2))

$$d = \sqrt{(x_2 - x_1)^2 + (y_2 - y_1)^2}$$

Example: $(1, 2)$ to $(4, 6)$: $d = \sqrt{3^2 + 4^2} = 5$

DISTANCE: WHO IS CLOSE TO WHOM

- 3-D (points (x_1, y_1, z_1) , (x_2, y_2, z_2))

$$d = \sqrt{(x_2 - x_1)^2 + (y_2 - y_1)^2 + (z_2 - z_1)^2}$$

Example: $(1, 2, 3)$ to $(4, 0, 6)$:squares $9 + 4 + 9 = 22 \rightarrow d = \sqrt{22} \approx 4.690$

n-D vectors $x, y \in \mathbb{R}^n$

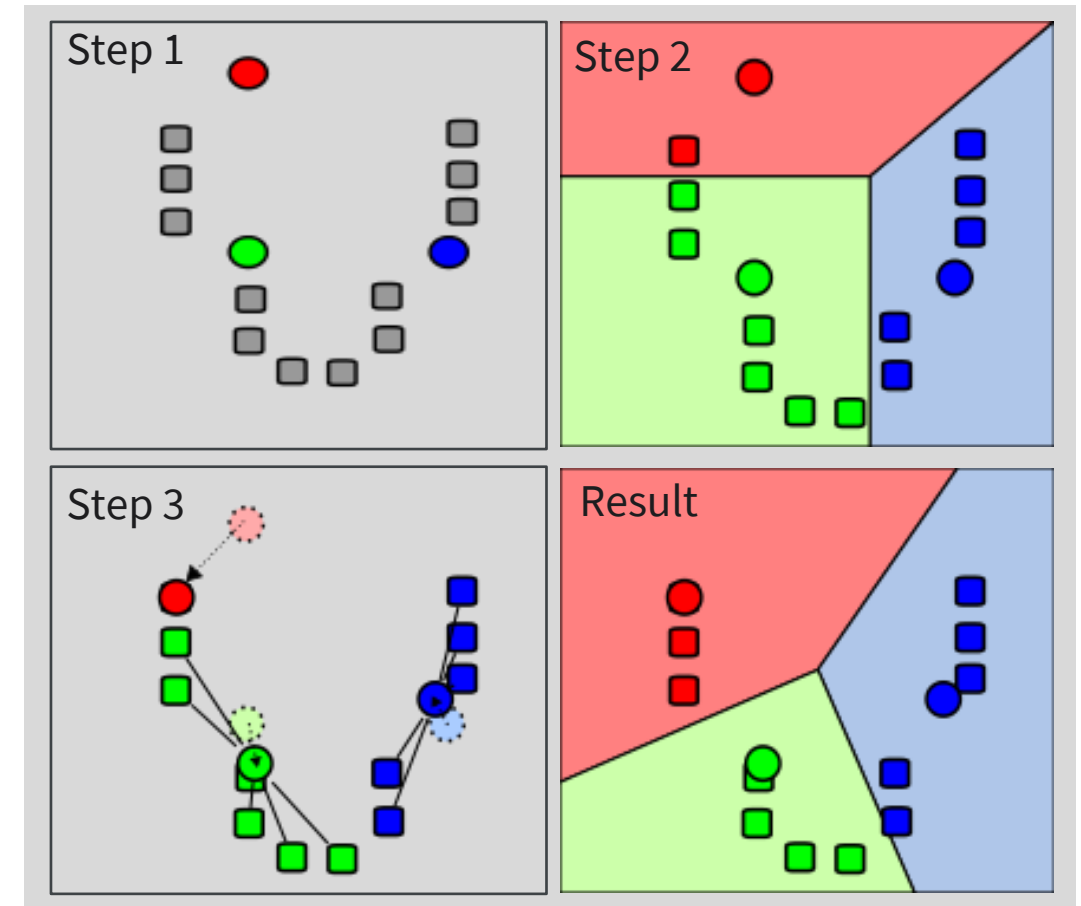
$$d(\mathbf{x}, \mathbf{y}) = \sqrt{\sum_{i=1}^n (x_i - y_i)^2}$$

DISTANCE: WHO IS CLOSE TO WHOM

- Squared Euclidean $d^2(\mathbf{x}, \mathbf{y}) = \sum_{i=1}^n (x_i - y_i)^2$
- Manhattan (L1) $d_1(\mathbf{x}, \mathbf{y}) = \sum_{i=1}^n |x_i - y_i|$
- Cosine similarity $\text{cos_sim}(\mathbf{x}, \mathbf{y}) = \frac{\mathbf{x} \cdot \mathbf{y}}{|\mathbf{x}| |\mathbf{y}|}$
- Cosine distance $\text{cos_dist} = 1 - \text{cos_sim}$
- Example $\text{cos_dist}((1, 2), (4, 6)) \approx 1 - \frac{16}{\sqrt{5} \sqrt{52}} \approx 0.0077$

2.1 CENTROID-BASED CLUSTERING

- **K-Means Clustering** partitions data points into K clusters
- Three main steps: initialization, assignment, update
- Advantages: **simple** to employ, efficient for large datasets
- Limitation: K must be determined, **initialization-dependent**, problems with high-dimensional data and outliers



Source of the text: Data Science Team, 2020; Pragathi, Jayanthi, & Malathi, 2018; Zhu, 2022.

Source of the image: File:K Means Example Step 1.svg (2020).

File:K Means Example Step 2.svg (2020).

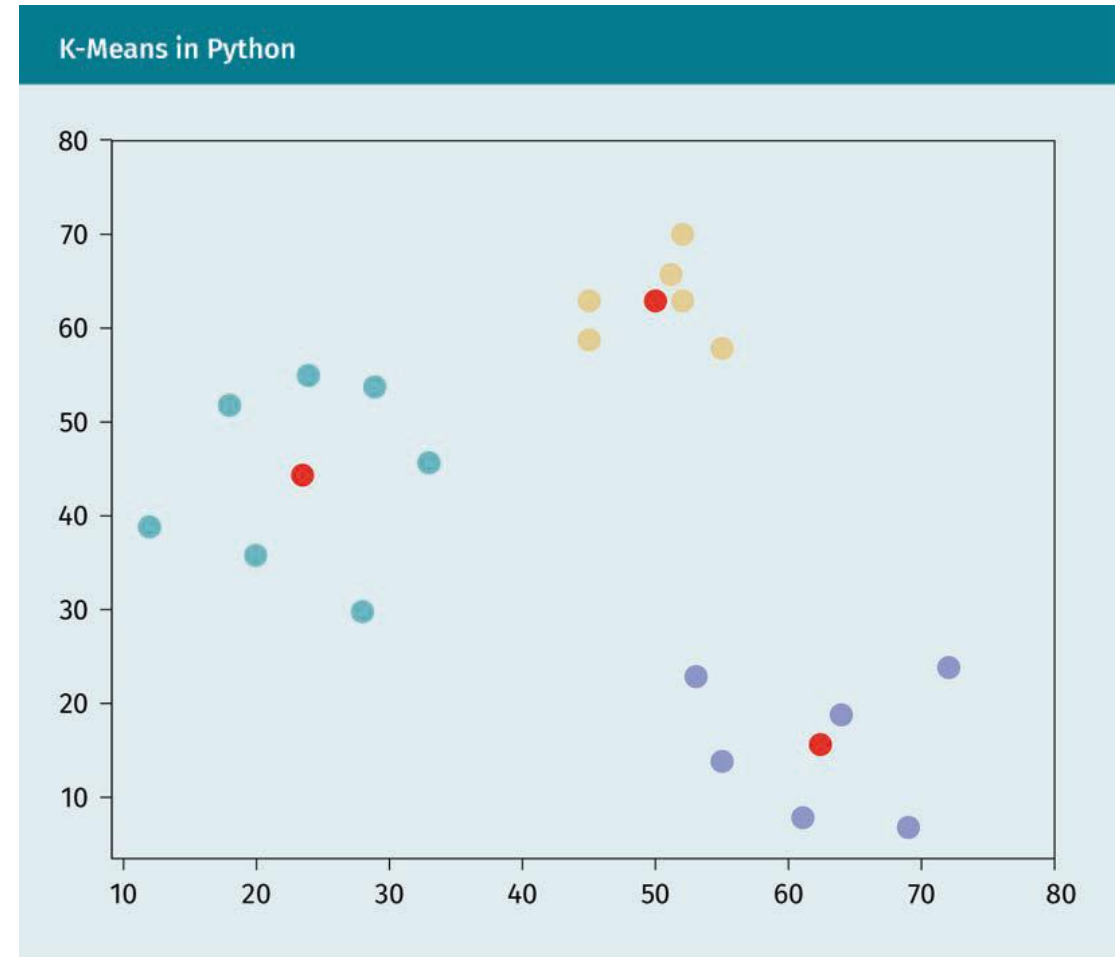
File:K Means Example Step 3.svg (2020).

File:K Means Example Step 4.svg (2020).

Algorithm of K-Means clustering

2.1 CENTROID-BASED CLUSTERING

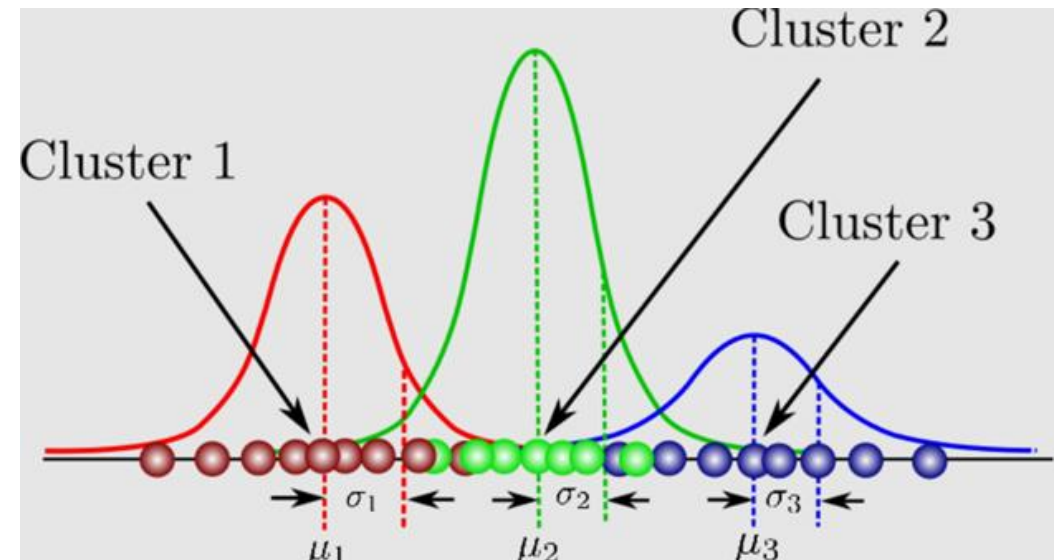
```
1 #K-Means clustering with Python
2 >>> import pandas as pd # For reading datasets
3 >>> import numpy as np # For computations
4 >>> import matplotlib.pyplot as plt # For visualization
5 >>> from pandas import DataFrame # For creating data frame
6 >>> from sklearn.cluster import KMeans
7 >>> Data={
8 'x': [12, 20, 28, 18, 29, 33, 24, 45, 45, 52, 51, 52, 55, 53, 55, 61, 64, 69, 72],
9 'y': [39, 36, 30, 52, 54, 46, 55, 59, 63, 70, 66, 63, 58, 23, 14, 8, 19, 7, 24]
10 }
11 >>> df = DataFrame(Data,columns=['x','y'])
12 # Create and fit the KMeans model
13 >>> kmeans = KMeans(n_clusters=3).fit(df)
14 # Find the centroids of the clusters
15 >>> centroids = kmeans.cluster_centers_
16 # Get the associated cluster for each data record
17 >>> kmeans.labels_
18 # Display the clusters contents and their centroids
19 >>> plt.scatter(df['x'], df['y'], c= kmeans.labels_.astype(float), s=50,
20 alpha=0.5)
21 >>> plt.scatter(centroids[:, 0], centroids[:, 1], c='red', s=50)
22 >>> plt.show()
```



Scatter Plot of Dataset With Clusters Identified Using K-means Clustering

2.2 GAUSSIAN MIXTURE MODELS CLUSTERING

- **GMM**: probabilistic generative models
- Each data point is assigned to the cluster returning the **highest probability**
- **Soft vs. hard clustering**
 - **K-means** does hard assignment: each point belongs to exactly one cluster.
 - **GMM** does **soft assignment**: each point has a *responsibility* (probability) for each cluster.



An example of GMM clustering

2.2 GAUSSIAN MIXTURE MODELS CLUSTERING

A GMM example with **one data point** getting **softly assigned to two clusters**.

Setup (1D, two Gaussians):

Cluster 1: $\mathcal{N}(\mu_1 = 0, \sigma_1^2 = 1)$, Cluster 2: $\mathcal{N}(\mu_2 = 4, \sigma_2^2 = 1)$, Mixing weights: $\pi_1 = \pi_2 = \frac{1}{2}$

Data point: $x = 1.5$

$$\text{Densities at } x: \mathcal{N}(x | \mu \sigma^2) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right)$$

$$N(1.5|0,1)=0.1295, N(1.5|4,1)=0.0175$$

$$\text{Responsibilities (soft assignments): } r_1 = p(z = 1 | x) = \frac{\pi_1 \mathcal{N}(x | \mu_1 \sigma_1^2)}{\pi_1 \mathcal{N}(x | \mu_1 \sigma_1^2) + \pi_2 \mathcal{N}(x | \mu_2 \sigma_2^2)} \approx 0.8808$$

$$r_2 = p(z = 2 | x) = 1 - r_1 \approx 0.1192$$

Interpretation: The single point $x = 1.5$ belongs to **both** clusters with probabilities about **88.1%** (cluster 1) and **11.9%** (cluster 2). That's GMM's “soft” assignment in action.

2.2 GAUSSIAN MIXTURE MODELS CLUSTERING

- The optimum means (μ) and standard deviations (σ) are found by applying the **Expectation maximization** algorithm
- Advantages: **efficient** for datasets with complex data distribution
- Limitation: **slow**, heavy computation, and can get stuck in local-maximum

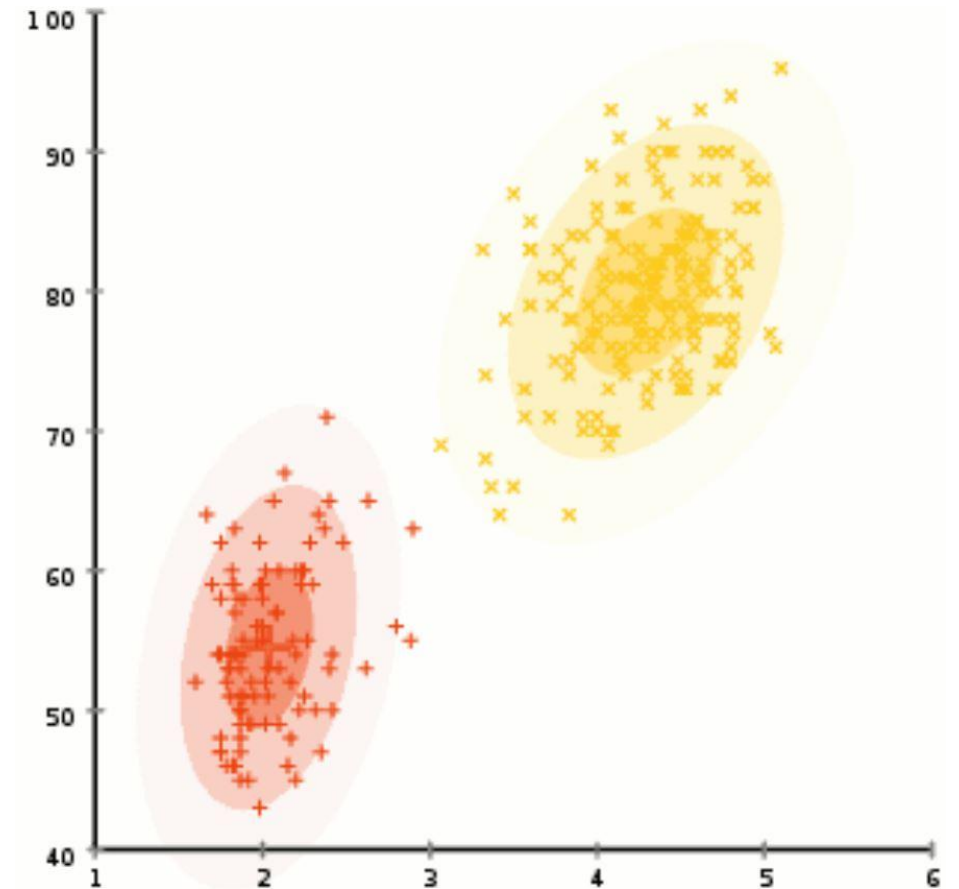
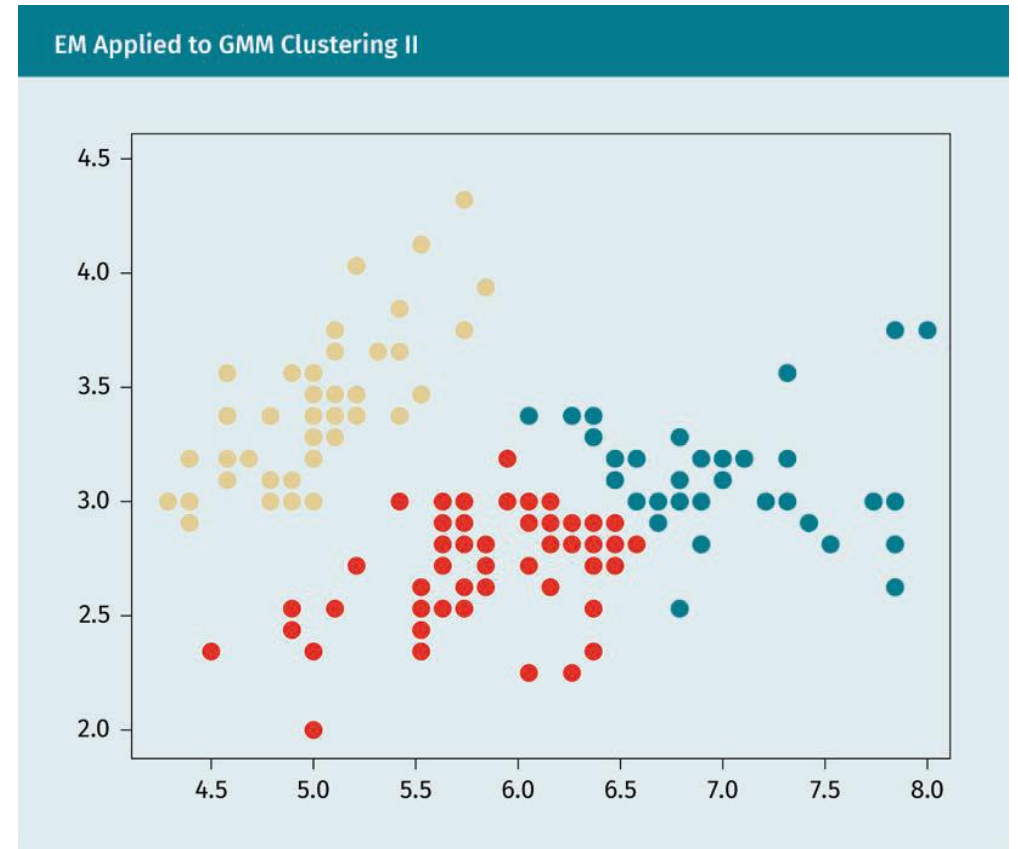


Illustration of expectation-maximization algorithm

2.2 GAUSSIAN MIXTURE MODELS CLUSTERING

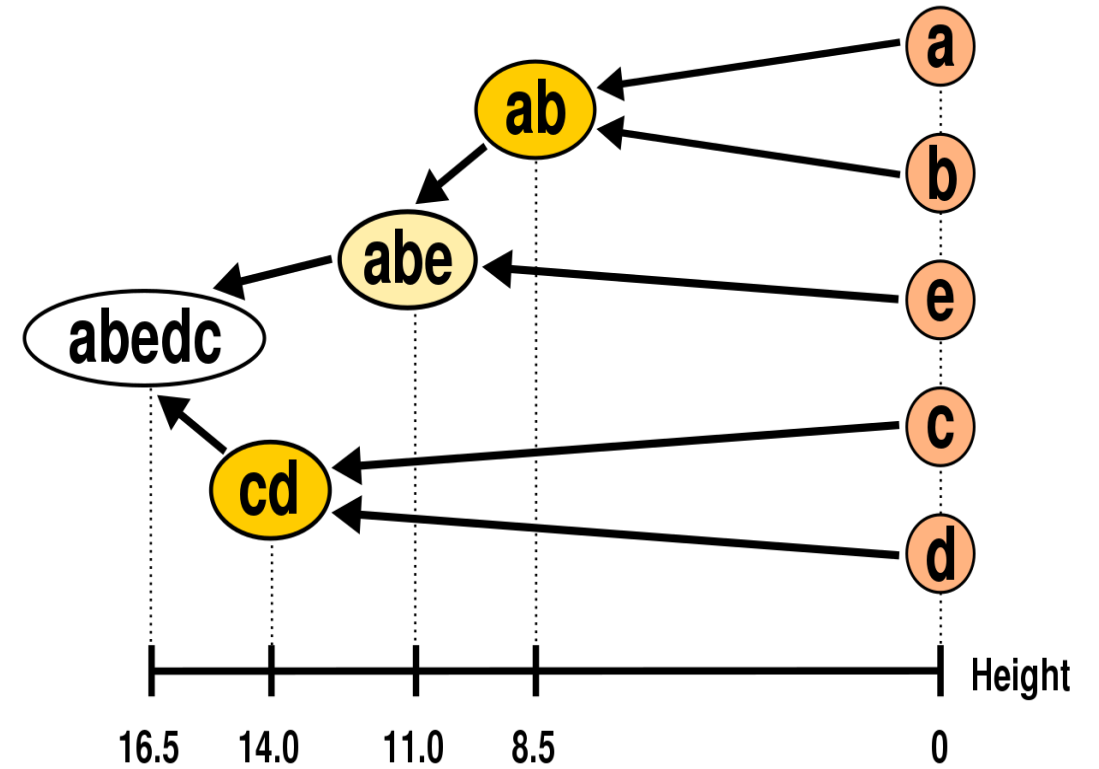
```
1 # GMM clustering with Python
2 >>> import numpy as np
3 >>> import pandas as pd
4 >>> import matplotlib.pyplot as plt
5 >>> from pandas import DataFrame
6 >>> from sklearn import datasets
7 >>> from sklearn.mixture import GaussianMixture
8 >>> iris = datasets.load_iris() # load the iris dataset
9 >>> X = iris.data[:, :2] # select first two columns
10 >>> d = pd.DataFrame(X) # turn it into a dataframe
11 >>> plt.scatter(d[0], d[1])
12 >>> plt.show() # plot the data
13 >>> gmm = GaussianMixture(n_components=3)
14 >>> gmm.fit(d) # fit the data as a mixture of 3 Gaussians
15 >>> labels = gmm.predict(d) # predict the cluster of each data record
16 >>> print('Converged:', gmm.converged_) # Check if the model has converged
17 >>> means = gmm.means_ # get the final "means" for each cluster
18 >>> covariances = gmm.covariances_ # get the final standard deviations
19 >>> d['labels'] = labels
20 >>> d0 = d[d['labels'] == 0]
21 >>> d1 = d[d['labels'] == 1]
22 >>> d2 = d[d['labels'] == 2]
23 >>> plt.scatter(d0[0], d0[1], c='r')
24 >>> plt.scatter(d1[0], d1[1], c='yellow')
25 >>> plt.scatter(d2[0], d2[1], c='g')
26 >>> plt.show() # plot the data records in each clusters in different color
```



Scatter Plot of Dataset With Clusters Identified Using GMM Clustering

2.3 HIERARCHICAL CLUSTERING

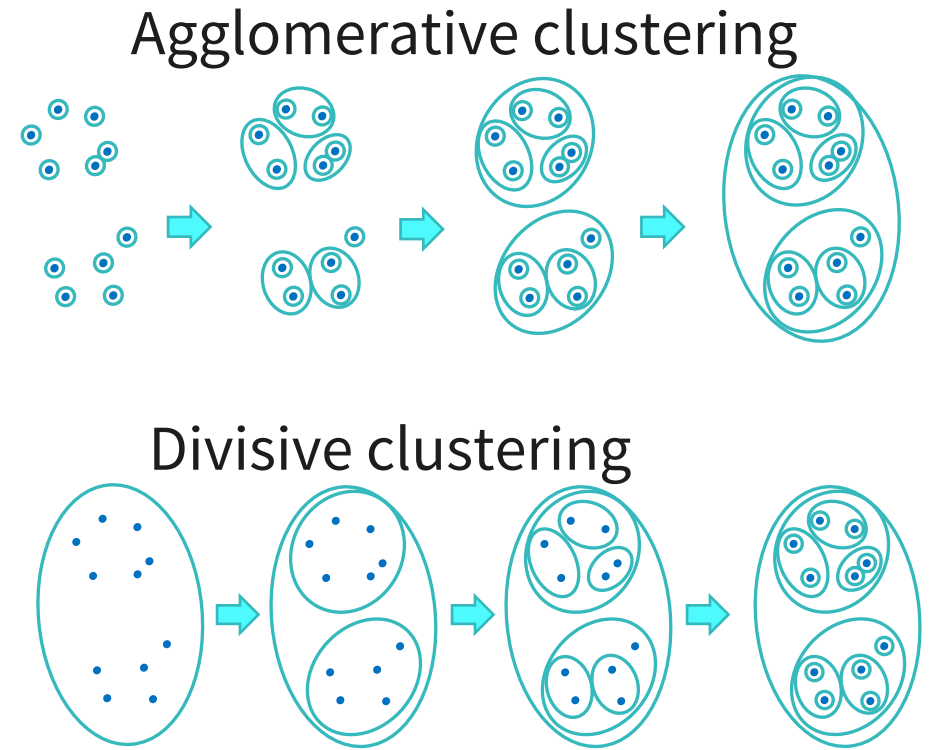
- The clusters are built in a hierarchy as a **tree**
- Three levels of clusters:
 - **Universe:** root
 - **Intermediate**
 - **Single-points.** leaves



Tree-view (dendrogram) of hierarchical clustering

2.3 HIERARCHICAL CLUSTERING

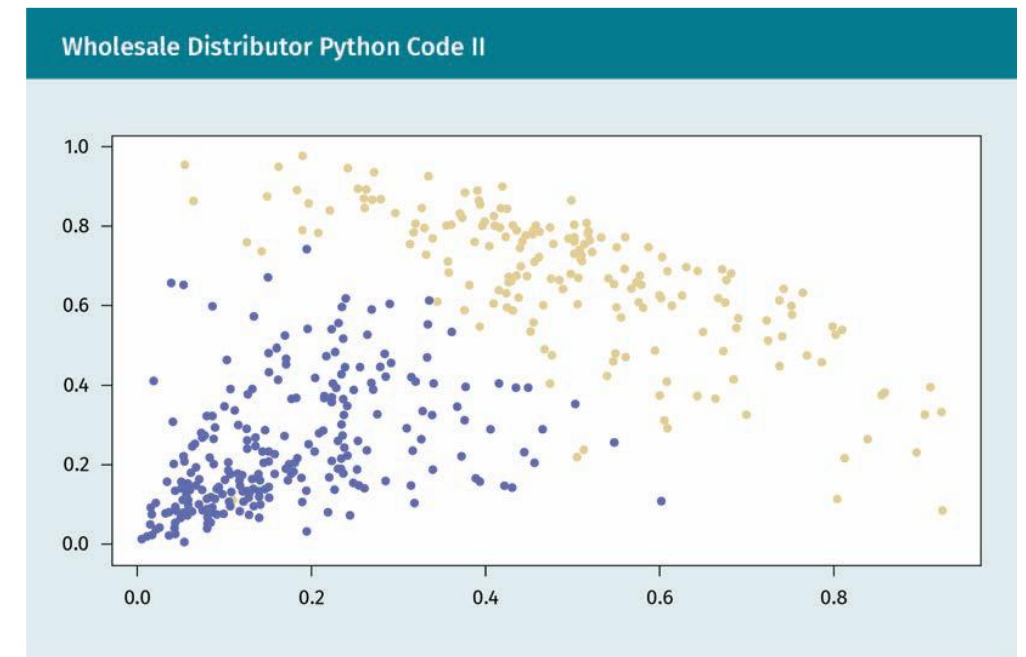
- Two types of clustering:
 - **Agglomerative**: bottom-up
 - **Divisive**: top-down
- Data points are merged/split **based on their distance**
- Advantages: easy to implement, intuitive, no assumption
- Limitation: not suitable for large datasets, sensitive to outliers



An example of hierarchical clustering

2.3 HIERARCHICAL CLUSTERING

```
1 # Agglomerative Clustering with Python
2 >>> import pandas as pd
3 >>> import numpy as np
4 >>> import matplotlib.pyplot as plt
5 >>> data = pd.read_csv('Wholesale customers data.csv')
6 # Normalize the dataset to get all the features at the same scale
7 >>> from sklearn.preprocessing import normalize
8 >>> data_scaled = normalize(data)
9 >>> data_scaled = pd.DataFrame(data_scaled, columns=data.columns)
10 # Draw the dendrogram to find the optimum number of clusters
11 >>> import scipy.cluster.hierarchy as shc
12 >>> plt.figure(figsize=(10, 7))
13 >>> plt.title("Dendrograms")
14 >>> dend = shc.dendrogram(shc.linkage(data_scaled, method='ward'))
15 >>> plt.show()
16 # From the dendrogram, we decide the optimum number of clusters is 2
17 # apply hierarchical clustering for two clusters only
18 >>> from sklearn.cluster import AgglomerativeClustering
19 >>> cluster = AgglomerativeClustering(n_clusters=2, affinity='euclidean', link
20 age='ward')
21 >>> cluster.fit_predict(data_scaled)
22 # visualize the two clusters
23 >>> plt.figure(figsize=(10, 7))
24 >>> plt.scatter(data_scaled['Milk'], data_scaled['Grocery'], c=cluster.labels_)
25 >>> plt.show()
```

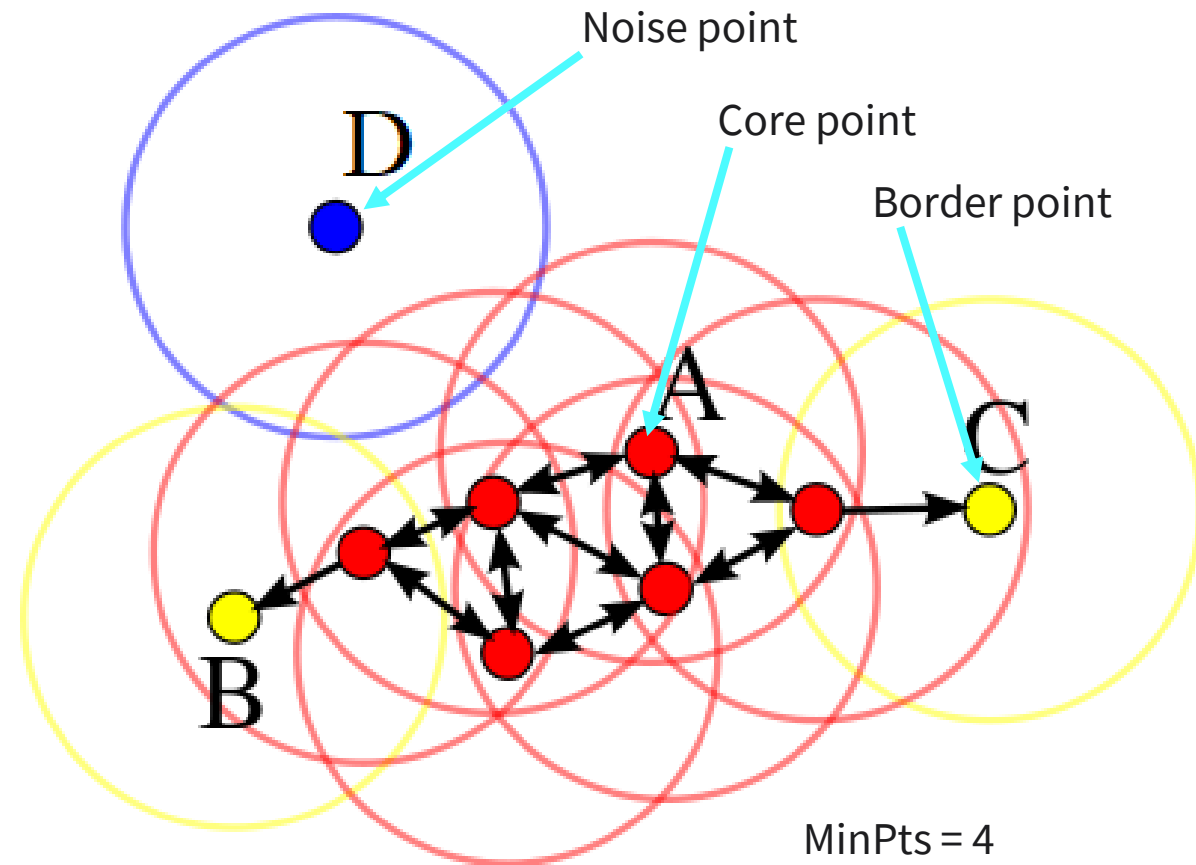


Scatter Plot of Dataset With Clusters Identified Using Agglomerative Clustering

2.4 DENSITY-BASED CLUSTERING

- Clusters are identified by grouping “dense” data points together, which permits:
 - the representation of **arbitrarily shaped clusters**,
 - the learning of **outliers**
- ϵ -Neighborhood of a point \mathbf{p} in the dataset \mathbf{d} :

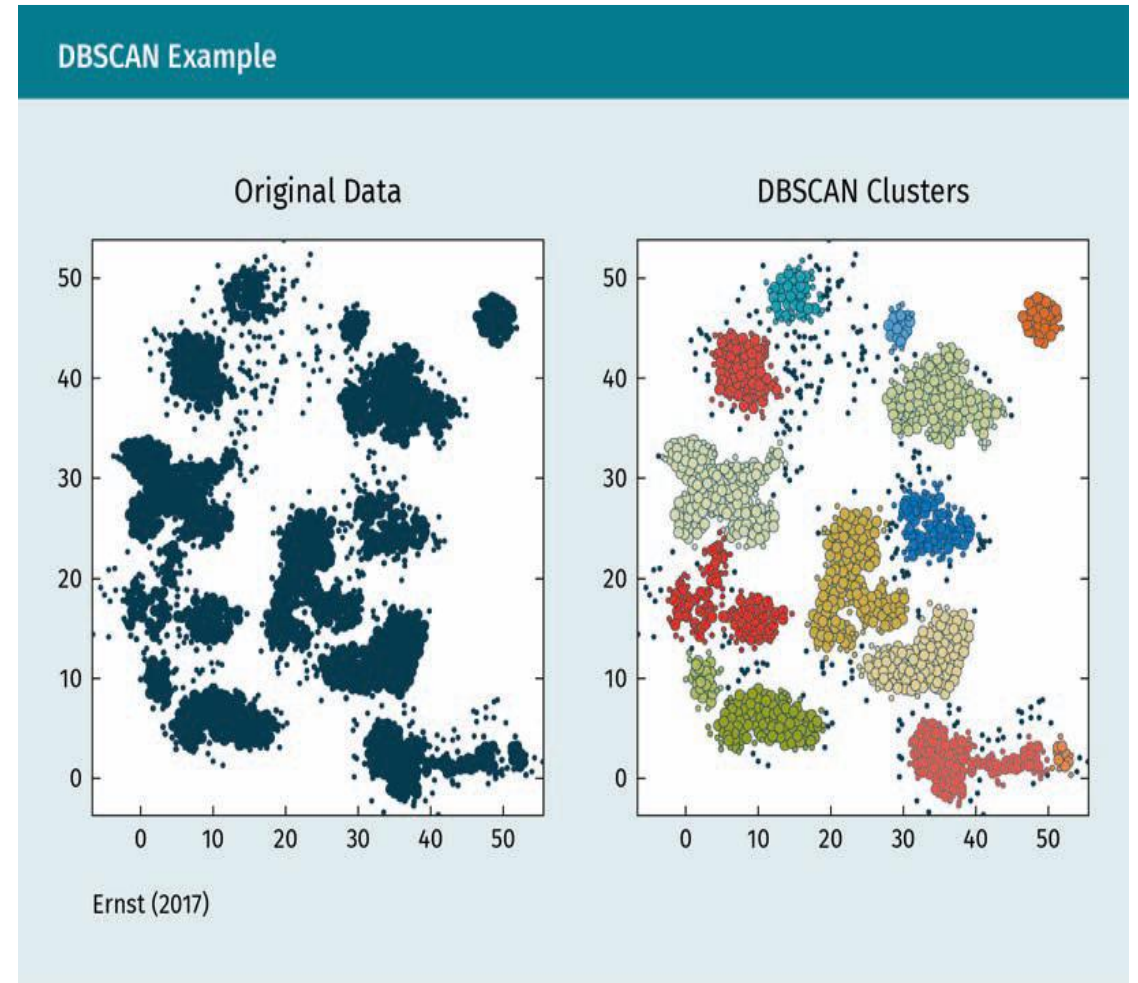
$$N(p) = \{q \in d \mid \text{dist}(p, q) \leq \epsilon\}$$



Algorithm of Density-based Clustering

2.4 DENSITY-BASED CLUSTERING

- Density-based spatial clustering for data with noise:
 1. Initialize MinPts and ϵ
 2. Start at a random data point p
 3. Check if p is a core point
 4. Density-connected clustering
 5. Start at new point, repeat 3 & 4
- Advantages: estimation of cluster number is not required
- Limitation: ϵ and MinPts setting up could be tricky



An example of DBSCAN Clustering

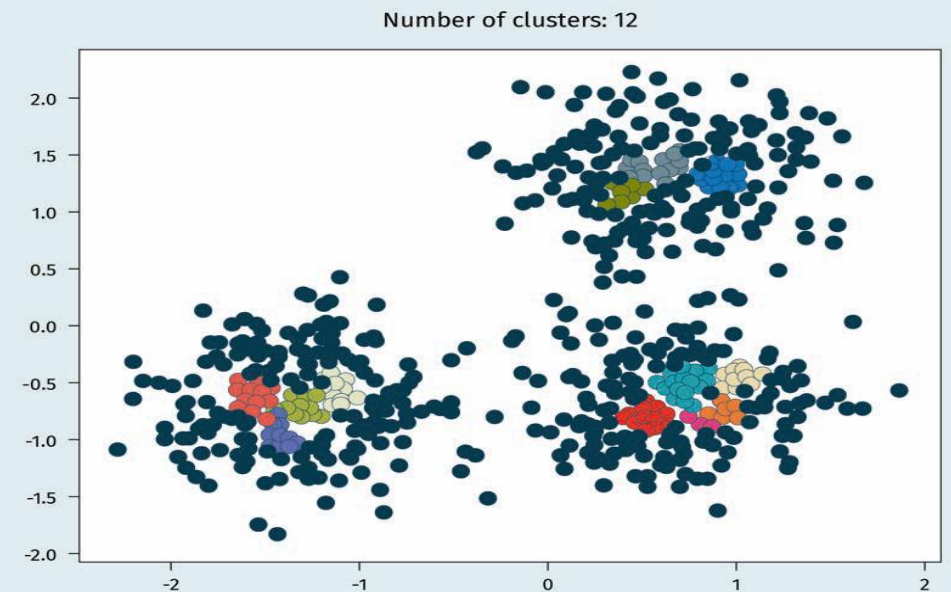
2.4 DENSITY-BASED CLUSTERING

```
1 # DBSCAN clustering with Python
2 >>> import numpy as np
3 >>> from sklearn.cluster import DBSCAN
4 >>> from sklearn import metrics
5 >>> from sklearn.datasets.samples_generator import make_blobs
6 >>> from sklearn.preprocessing import StandardScaler
7 >>> from pylab import *
8 >>> import matplotlib.pyplot as plt
9 >>> # Generate sample data
10 >>> centers = [[1, 1], [-1, -1], [1, -1]]
11 >>> X, labels_true = make_blobs(n_samples=750, centers=centers,
12 cluster_std=0.4, random_state=0)
13 >>> X = StandardScaler().fit_transform(X) # standardize the dataset
14 >>> xx, yy = zip(*X)
15 >>> db = DBSCAN(eps=0.3, min_samples=10).fit(X) # Set up parameters
16 >>> core_samples = db.core_sample_indices_
17 >>> core_samples_mask = np.zeros_like(db.labels_, dtype=bool)
18 >>> core_samples_mask[db.core_sample_indices_] = True
19 >>> n_clusters_ = len(set(labels))- (1 if -1 in labels else 0) # Nr of clusters
20 >>> labels = db.labels_
21 >>> outliers = X[labels == -1] # find the outliers
22 >>> cluster1 = X[labels == 0] # Get the contents of each cluster
23 >>> cluster2 = X[labels == 1]
24 >>> cluster3 = X[labels == 2]
25 >>> unique_labels = set(labels)
26 >>> colors = ['y', 'b', 'g', 'r']
```

Source of the image: Zöller (2022).

```
27 >>> for k, col in zip(unique_labels, colors): if k == -1: col = 'k'
28 class_member_mask = (labels == k)
29 xy = X[class_member_mask & core_samples_mask]
30 plt.plot(xy[:, 0], xy[:, 1], 'o', markerfacecolor=col,
31 markedgecolor='k', markersize=6)
32 xy = X[class_member_mask & ~core_samples_mask]
33 plt.plot(xy[:, 0], xy[:, 1], 'o', markerfacecolor=col,
34 markedgecolor='k', markersize=6)
35 >>> plt.title('number of clusters: %d' % n_clusters_)
36 >>> plt.show()
```

DBSCAN Clustering Example (eps = 0.1, min_samples = 10)



Using DBSCAN Clustering

REVIEW STUDY GOALS



- Know the definitions and terms used for clustering
- Comprehend common applications of clustering analysis
- Understand different methods for clustering analysis
- Analyse the advantages and limitations of the clustering methods
- Implement clustering methods in Python

SESSION 2

TRANSFER TASK

TRANSFER TASKS

Create the dataset using the following commands:

```
>>> from sklearn.datasets import make_blobs
```

```
>>> dataset, _ = make_blobs(n_samples=100, random_state=12)
```

Partition the data points into 3 clusters by using:

1. K-Means clustering
2. GMM clustering
3. Agglomerative clustering
4. DBSCAN with $\epsilon=1$ and MinPts=10

TRANSFER TASKS 2

<https://scikit-learn.org/stable/modules/clustering.html>

Codes\KMeans_on_Digits_StepByStep.ipynb

TRANSFER TASK
PRESENTATION OF THE RESULTS

Please present your
results.

The results will be
discussed in plenary.





1. Which of the following clustering algorithms suffers from the problem of convergence to local optima?

- a) Agglomerative clustering
- b) Only Expectation-Maximization clustering
- c) K- Means clustering and expectation-maximization clustering
- d) DBSCAN



2. For clustering to be applied, what is the minimum number of variables/ features required?

- a) 0
- b) 1
- c) 2
- d) 3



3. Is it possible that assignment of observations to clusters does not change between successive iterations in K-Means?

- a) Yes
- b) No
- c) Cannot say
- d) It depend on the number of data points



4. What could be the possible reason(s) for producing two different dendrograms using the agglomerative clustering algorithm for the same dataset?

- a) Proximity function used
- b) Number of data points used
- c) Number of variables used
- d) All these options are possible



5. K-Means clustering with Euclidean distance measure produces which kind of cluster shapes?

- a) Spherical
- b) Ellipsoid
- c) Arbitrary
- d) Cubic

LIST OF SOURCES

Text:

Carrasco, O.C. (2019). Gaussian Mixture Models Explained. <https://towardsdatascience.com/gaussian-mixture-models-explained-6986aaf5a95>

Das, J., Mukherjee, P., Majumder, S., & Gupta, P. (2014). Clustering-Based recommender system using principles of voting theory. *International Conference on Contemporary Computing and Informatics (IC3I 2014)*.

Data Science Team (2020). A Layman's Guide to K Means Clustering. <https://datascience.eu/machine-learning/k-means-clustering-in-machine-learning/>

Ester, M., Kriegel, H., Sander, J., & Xiaowei, X. (1996). A density-based algorithm for discovering clusters in large spatial databases with noise. *Proceedings of the Second International Conference on Knowledge Discovery and Data Mining*, 226—231. <https://www.aaai.org/Papers/KDD/1996/KDD96-037.pdf>

Pai, P. (2021). Hierarchical clustering explained. <https://towardsdatascience.com/hierarchical-clustering-explained-e59b13846da8>

Pragathi, K., Jayanthi, T., & Malathi V. (2018). Cluster based segmentation using K-Means algorithm. *International Journal of Data Mining Techniques and Applications*, 7(1), 64—67. http://ijdmata.com/abstract_temp.php?id=V7-I1-P13

Zhu, A. (2022). 7 most asked questions on K-Means Clustering. <https://towardsdatascience.com/explain-ml-in-a-simple-way-k-means-clustering-e925d019743b>

Image:

Ernst (2017).

Zöller (2022).

File:KMeans-density-data.svg. (2020, October 10). In *Wikimedia Commons, the free media repository*. Retrieved, January 25, 2023, from <https://commons.wikimedia.org/w/index.php?title=File:KMeans-density-data.svg&oldid=485735167>.

File:EM Clustering of Old Faithful data.gif. (2020, September 23). In *Wikimedia Commons, the free media repository*. Retrieved, January 25, 2023, from https://commons.wikimedia.org/w/index.php?title=File:EM_Clustering_of_Old_Faithful_data.gif&oldid=469170019.

File:UPGMA Dendrogram Hierarchical.svg. (2020, October 31). In *Wikimedia Commons, the free media repository*. Retrieved, January 25, 2023, from https://commons.wikimedia.org/w/index.php?title=File:UPGMA_Dendrogram_Hierarchical.svg&oldid=508481311.

File:DBSCAN-Illustration.svg. (2020, October 7). In *Wikimedia Commons, the free media repository*. Retrieved, January 25, 2023, from <https://commons.wikimedia.org/w/index.php?title=File:DBSCAN-Illustration.svg&oldid=483221815>.

File:K Means Example Step 1.svg. (2020, September 4). In *Wikimedia Commons, the free media repository*. Retrieved, January 26, 2023, from https://commons.wikimedia.org/w/index.php?title=File:K_Means_Example_Step_1.svg&oldid=448124955.

File:K Means Example Step 2.svg. (2020, September 4). In *Wikimedia Commons, the free media repository*. Retrieved, January 26, 2023, from https://commons.wikimedia.org/w/index.php?title=File:K_Means_Example_Step_2.svg&oldid=448124977.

File:K Means Example Step 3.svg. (2020, September 4). In *Wikimedia Commons, the free media repository*. Retrieved, January 26, 2023, from https://commons.wikimedia.org/w/index.php?title=File:K_Means_Example_Step_3.svg&oldid=448124994.

File:K Means Example Step 4.svg. (2020, September 4). In *Wikimedia Commons, the free media repository*. Retrieved, January 26, 2023, from https://commons.wikimedia.org/w/index.php?title=File:K_Means_Example_Step_4.svg&oldid=448125010.

© 2022 IU Internationale Hochschule GmbH

This content is protected by copyright. All rights reserved.

This content may not be reproduced and/or electronically edited, duplicated, or distributed in any kind of form without written permission by the IU Internationale Hochschule GmbH.