



Machine Learning #18

▼ 1. What is the difference between supervised and unsupervised learning? Give some examples to illustrate your point.

Supervised learning and unsupervised learning are two major categories of machine learning techniques.

Supervised learning involves training a model on a labeled dataset, where the model learns to predict an output variable based on one or more input variables. The labeled data contains both the input and the desired output, and the model is trained to make accurate predictions based on this data. Some examples of supervised learning include:

- Image classification: Given a labeled dataset of images, the model learns to classify new images into one of several pre-defined categories.
- Regression analysis: Given a labeled dataset of numerical inputs and outputs, the model learns to predict new output values based on new input values.
- Natural language processing: Given a labeled dataset of text and their corresponding categories, the model learns to classify new text into one of the categories.

Unsupervised learning, on the other hand, **involves training a model on an unlabeled dataset, where the model must identify patterns or relationships in the data without any pre-existing labels.** The model must identify and group similar data points together based on their characteristics. Some examples of unsupervised learning include:

- Clustering: Given an unlabeled dataset, the model must group similar data points together into clusters based on their similarities.
- Dimensionality reduction: Given a high-dimensional dataset, the model must identify the most important features or variables and reduce the dimensionality of the dataset.
- Anomaly detection: Given a dataset of mostly normal data, the model must identify rare or anomalous data points.

Overall, the key difference between supervised and unsupervised learning is that supervised learning involves learning from labeled data with a known output, while unsupervised learning involves learning from unlabeled data without a known output.

▼ 2. Mention a few unsupervised learning applications.

Unsupervised learning has various applications in the field of machine learning, some of which include:

1. **Clustering:** Grouping similar objects or data points together based on some common characteristics. Example: Customer segmentation in marketing, image segmentation in computer vision.
2. **Anomaly detection:** Identifying rare or unusual observations in data that do not conform to the expected pattern. Example: Fraud detection in finance, network intrusion detection in cybersecurity.

3. **Dimensionality reduction:** Reducing the number of features or variables in the dataset by extracting the most important information. Example: Image compression, data visualization.
4. **Association rule learning:** Identifying interesting relationships or patterns in data. Example: Market basket analysis to identify product associations in retail.
5. **Generative models:** Creating new data that resembles the original data based on the patterns in the dataset. Example: Image generation, music generation.

▼ 3. What are the three main types of clustering methods? Briefly describe the characteristics of each.

The three main types of clustering methods are hierarchical clustering, partitioning clustering, and density-based clustering.

1. **Hierarchical clustering:** *This method involves creating a tree-like structure of clusters by iteratively merging or dividing clusters based on their similarity or distance.* It can be either agglomerative or divisive. In agglomerative clustering, each data point initially represents a cluster, which are then recursively merged into larger clusters until all points belong to one cluster. In divisive clustering, all data points initially belong to a single cluster, which is then recursively divided into smaller clusters until each point represents its cluster. The main advantage of hierarchical clustering is that it provides a visual representation of the clustering process.
2. **Partitioning clustering:** *This method partitions the data set into k clusters, where k is a predefined number. It aims to minimize the sum of the squared distances between the data points and their assigned cluster centers.* K-means is a popular partitioning clustering algorithm that initializes k cluster centers randomly, assigns each data point to the nearest center, recalculates the centers' positions based on the points assigned to them, and repeats the process until convergence.
3. **Density-based clustering:** *This method identifies clusters as dense regions of points separated by low-density regions.* The most popular density-based clustering algorithm is DBSCAN (Density-Based Spatial Clustering of Applications with Noise), which groups points that are close to each other in terms of distance and density, while considering points that do not belong to any cluster as noise. It can detect clusters of arbitrary shapes and is less sensitive to the initial conditions than other clustering algorithms.

▼ 4. Explain how the k-means algorithm determines the consistency of clustering.

The k-means algorithm is a popular clustering technique used in unsupervised machine learning. It aims to partition a given dataset into k clusters, where k is a predetermined number. The algorithm works by iterating through two steps: assigning data points to the closest cluster center and then updating the cluster centers based on the new assignments. The algorithm repeats these steps until convergence.

The consistency of clustering in the k-means algorithm is measured by the sum of squared errors (SSE). The SSE represents the sum of the squared distances between each data point and its assigned cluster center. The goal of the k-means algorithm is to minimize the SSE value, which indicates a more consistent and compact clustering.

The algorithm starts by selecting k random initial cluster centers. Each data point is then assigned to the closest cluster center based on the Euclidean distance metric. After all data points have been assigned, the cluster centers are updated by taking the mean of all data points belonging to each

cluster. The algorithm repeats these two steps until convergence or a maximum number of iterations is reached.

The consistency of clustering is determined by calculating the SSE at each iteration. As the k-means algorithm proceeds, the SSE value decreases with each iteration until convergence. The ideal clustering is achieved when the SSE value reaches its minimum and the resulting clusters are consistent and compact.

▼ 5. With a simple illustration, explain the key difference between the k-means and k-medoids algorithms.

Sure, here is an example that illustrates the difference between k-means and k-medoids clustering algorithms using Python.

First, let's generate some random data points:

```
import numpy as np

# Generate random data points
np.random.seed(1)
X = np.random.randn(50, 2)
```

Next, let's apply k-means clustering algorithm to the data points:

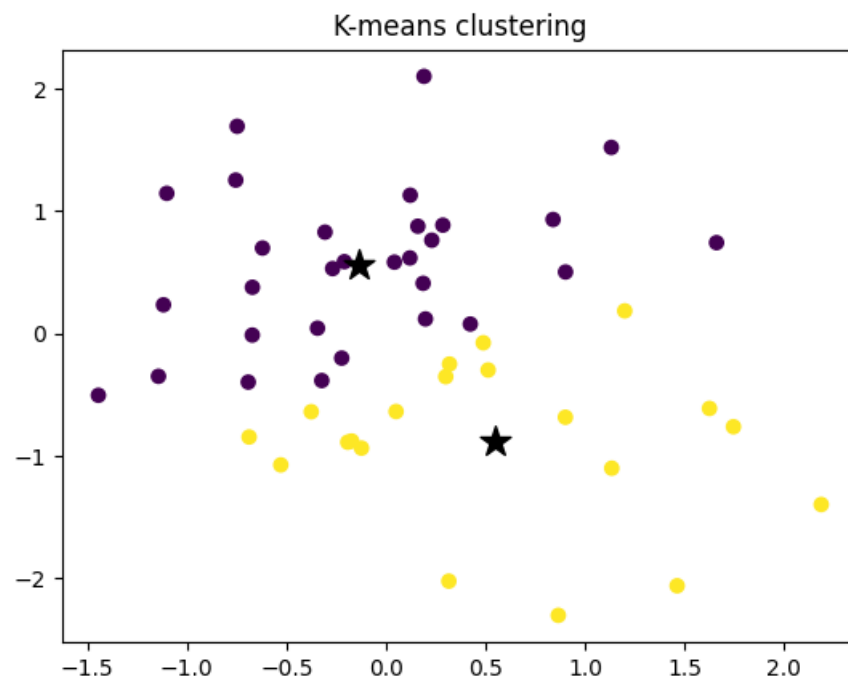
```
from sklearn.cluster import KMeans

# Apply k-means clustering with k=2
kmeans = KMeans(n_clusters=2, random_state=1).fit(X)

# Plot the results
import matplotlib.pyplot as plt

plt.scatter(X[:, 0], X[:, 1], c=kmeans.labels_, cmap='viridis')
plt.scatter(kmeans.cluster_centers_[:, 0], kmeans.cluster_centers_[:, 1], marker='*', s=200, color='black')
plt.title('K-means clustering')
plt.show()
```

This will produce a plot where the data points are colored based on their cluster assignments, and the centroids of the clusters are marked with black stars:



Now, let's apply k-medoids clustering algorithm to the same data points:

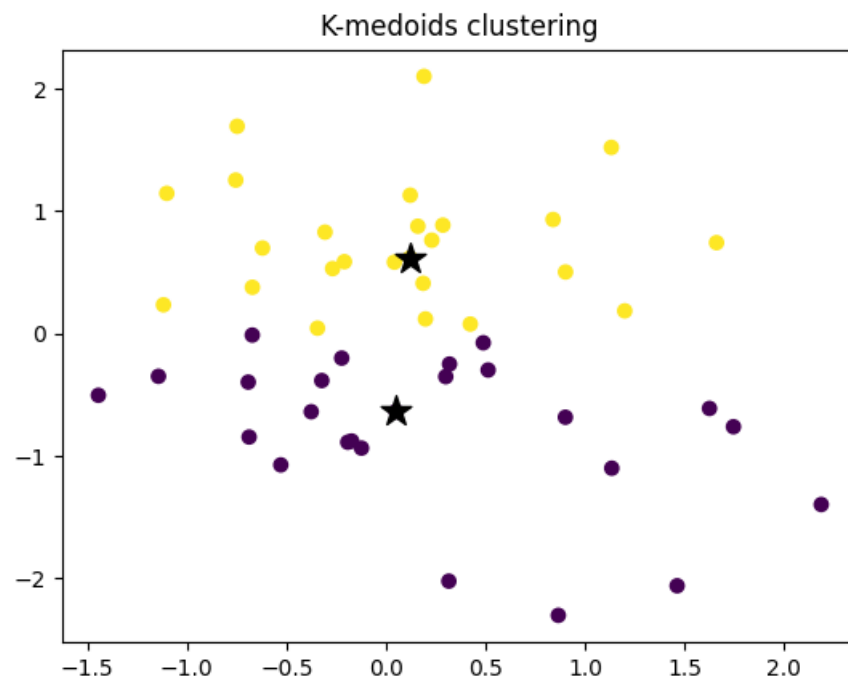
```
!pip install scikit-learn-extra

from sklearn_extra.cluster import KMedoids

# Apply k-medoids clustering with k=2
kmedoids = KMedoids(n_clusters=2, random_state=1).fit(X)

# Plot the results
plt.scatter(X[:, 0], X[:, 1], c=kmedoids.labels_, cmap='viridis')
plt.scatter(kmedoids.cluster_centers_[0], kmedoids.cluster_centers_[1], marker='*', s=200, color='black')
plt.title('K-medoids clustering')
plt.show()
```

This will produce a plot where the data points are colored based on their cluster assignments, and the medoids (data points closest to the centroid) of the clusters are marked with black stars:



The key difference between the two algorithms is that k-means uses the mean (centroid) of the cluster to update the cluster assignment and centroid position, while k-medoids uses the medoid (data point closest to the centroid) of the cluster to update the cluster assignment and centroid position. This means that k-medoids is more robust to outliers and noise in the data, as it selects the medoid based on the distance to other data points in the cluster, while k-means can be influenced by outliers that are far away from the centroid.

▼ 6. What is a dendrogram, and how does it work? Explain how to do it.

A dendrogram is a tree-like diagram that shows the hierarchical relationship between objects or groups of objects. It is often used in hierarchical clustering to represent the order and distances between clusters. The key idea is to represent each observation as a single point, and then merge them iteratively based on their similarity.

To construct a dendrogram, the following steps are typically followed:

1. Start with each observation as a single cluster.
2. Measure the distance between all pairs of clusters.
3. Merge the two closest clusters into a single cluster.
4. Recalculate the distance between the new cluster and all remaining clusters.
5. Repeat steps 3-4 until all observations are merged into a single cluster.

The resulting dendrogram is a hierarchical tree-like diagram that shows the order and distances between the clusters. The height of the vertical lines in the dendrogram represents the distance between the clusters or subclusters being merged, and the horizontal axis represents the observations or groups of observations being merged.

Dendrograms can be visualized using various software tools such as Python's `scipy` library, which provides a `dendrogram()` function that takes in a distance matrix and plots the dendrogram. Here is an

example code snippet to illustrate:

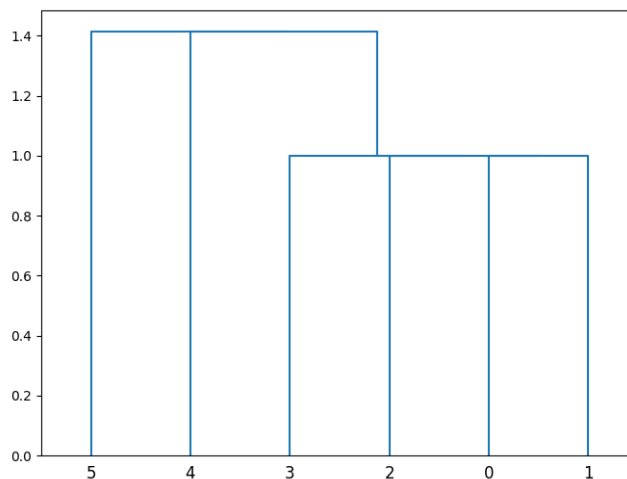
```
import scipy.cluster.hierarchy as sch
import matplotlib.pyplot as plt

# Generate sample data
X = [[0, 0], [0, 1], [1, 0], [1, 1], [2, 2], [3, 3]]

# Compute the distance matrix using Euclidean distance
d = sch.distance.pdist(X)

# Compute the hierarchical clustering
Z = sch.linkage(d, method='single')

# Plot the dendrogram
fig, ax = plt.subplots(figsize=(8, 6))
sch.dendrogram(Z, ax=ax)
plt.show()
```



This code generates a dendrogram for a small sample dataset `x` with six observations. The distance matrix is computed using Euclidean distance, and the hierarchical clustering is performed using the `single` linkage method. Finally, the dendrogram is plotted using the `dendrogram()` function. The resulting plot shows the hierarchical structure of the data, where the first two observations are grouped together at the bottom, and the remaining observations are grouped together at the top.

▼ 7. What exactly is SSE? What role does it play in the k-means algorithm?

SSE stands for Sum of Squared Errors, and *it is a measure of how much variance exists within the clusters formed by the k-means algorithm*. In the k-means algorithm, the objective is to minimize the SSE, which is calculated as the *sum of the squared distance between each data point and its assigned centroid*.

The SSE is used as an optimization metric for the k-means algorithm because the algorithm aims to minimize the distance between data points and their assigned centroids while also minimizing the number of centroids used. As a result, minimizing the SSE leads to the formation of tight, compact clusters with minimal variance.

In each iteration of the k-means algorithm, the SSE is calculated and used as a convergence criterion. The algorithm iteratively adjusts the positions of the centroids to minimize the SSE until convergence

or a maximum number of iterations is reached. By minimizing the SSE, the k-means algorithm can identify the optimal number of clusters and their respective centroids to best represent the underlying data structure.

▼ **8. With a step-by-step algorithm, explain the k-means procedure.**

The k-means algorithm is a popular clustering algorithm used to partition a dataset into k clusters. The steps involved in the k-means algorithm are:

1. Choose the number of clusters k that the data should be divided into.
2. Initialize k cluster centroids randomly from the data points.
3. For each data point, calculate the distance to each of the k centroids, and assign the data point to the cluster with the closest centroid.
4. For each cluster, calculate the mean of the data points assigned to that cluster, and update the centroid of the cluster to be the mean.
5. Repeat steps 3-4 until the cluster assignments no longer change or a maximum number of iterations is reached.

▼ **9. In the sense of hierarchical clustering, define the terms single link and complete link.**

Single link and complete link are two different types of linkage methods used in hierarchical clustering.

In **single link clustering**, also known as the nearest neighbor method, *the distance between two clusters is defined as the distance between the closest pair of objects in the two clusters*. This method tends to produce long, elongated clusters.

In **complete link clustering**, also known as the farthest neighbor method, *the distance between two clusters is defined as the distance between the farthest pair of objects in the two clusters*. This method tends to produce compact, spherical clusters.

▼ **10. How does the apriori concept aid in the reduction of measurement overhead in a business basket analysis? Give an example to demonstrate your point.**

In a business basket analysis, **Apriori is a data mining algorithm used for identifying frequently occurring sets of items (known as itemsets) in a transactional database**. This algorithm helps to identify the relationships between different items and determine which items are frequently purchased together.

The Apriori algorithm reduces measurement overhead by using a minimum support threshold to filter out infrequent itemsets. This means that only the itemsets that meet the minimum support threshold are considered for analysis, reducing the amount of data that needs to be processed. This, in turn, leads to faster processing times and more efficient analysis.

For example, let's say we want to analyze the purchasing patterns of customers in a grocery store. We have a transactional database that records all customer purchases, and we want to identify which items are frequently purchased together. We can use the Apriori algorithm to help us do this.

First, we set a minimum support threshold, say 5%, which means that we only consider itemsets that occur in at least 5% of transactions. Then we apply the Apriori algorithm to the data to identify all itemsets that meet this threshold. These frequent itemsets can then be used to identify which items are commonly purchased together, which can help the grocery store to optimize their product placement and promotions.

By using the Apriori algorithm to filter out infrequent itemsets, we can reduce the amount of data that needs to be processed and make our analysis more efficient.