

Click to verify



While trying to analyze data, one approach might be to look for meaningful groups or clusters. Clustering is dividing data into groups based on similarity. And K-means is one of the most commonly used methods in clustering. Why? The main reason is its simplicity. In this tutorial, we'll start with the theoretical foundations of the K-means algorithm, we'll discuss how it works and what pitfalls to avoid. Then, we'll see a practical application of the K-means algorithm with Python using the sklearn library. So, let's get right into it and see what K-means is all about. How Does K-means Clustering Work? Let's say we'd like to divide the following points into clusters. First, we must choose how many clusters we'd like to have. The K in K-means stands for the number of clusters we're trying to identify. In fact, that's where this method gets its name from. We can start by choosing two clusters. The second step is to specify the cluster seeds. A seed is basically a starting cluster centroid. It is chosen at random or is specified by the data scientist based on prior knowledge about the data. One of the clusters will be the green cluster, and the other one - the orange cluster. And these are the seeds. The next step is to assign each point on the graph to a seed. Which is done based on proximity. For instance, this point is closer to the green seed than to the orange one. Therefore, it will belong to the green cluster. The next step is to reassign the seeds. The new seeds will be the centroid of each cluster. The green seed will move closer to the green points, and the orange seed will move closer to the orange points. From here, we can repeat the last step. We can do that 10, 15 or 1000 times, until we've reached a clustering solution where we can no longer adjust any of the clusters. Sounds simple, right? However, due to its simplicity, K-means is prone to some issues. What Are the Disadvantages of K-means? One disadvantage arises from the fact that in K-means we have to specify the number of clusters before starting. In fact, this is an issue that a lot of the clustering algorithms share. In the case of K-means if we choose K too small, the cluster centroid will not lie inside the clusters. As we can see, in this example, this is not representative of the data. In cases where K is too large, some of the clusters may be split into two. Reasonable enough, right? Another important issue is that K-means enforces clusters with a spherical shape or blobs. The reason is that we are trying to minimize the distance from the centroids in a straight line. So, if we have clusters, which are more elongated, K-means will have difficulty separating them. Now that you're familiar with what K-means is and how it works, let's bring theory into practice with an example. What's K-Means Clustering's Application? One of K-means most important applications is dividing a data set into clusters. So, as an example, we'll see how we can implement K-means in Python. To do that, we'll use the sklearn library, which contains a number of clustering modules, including one for K-means. Let's say we have our segmentation data in a csv file. After we've read the file (in our case using the pandas method) we can proceed with implementing K-means. You'll need to import K-means from sklearn.cluster. As mentioned above, K-means doesn't tell us how many clusters there are. Instead, it minimizes the Euclidean norm. However, we usually want to be able to determine the number of clusters, too, right? Well, the way if we assume there are 2 clusters and we calculate the sum of squared distances (inertia) from each point to its assigned centroid against the number of clusters. As the number of clusters increases, the inertia decreases. The point where the decrease sharply slows down (forming an elbow) indicates the optimal number of clusters. This method's popularity stems from its simplicity and visual appeal. It allows data scientists and analysts to quickly identify a suitable number of clusters by visual inspection. Pros: Simple to understand and implement. Provides a visual indication of the optimal number of clusters. Cons: The elbow point can sometimes be ambiguous or not very pronounced, making it hard to decide the optimal number of clusters. Silhouette Method The Silhouette Score or Silhouette Coefficient evaluates the consistency within clusters and the separation between clusters. It ranges from -1 to 1, with higher values indicating better clustering. The method is in detail described on Wikipedia. The score is calculated as follows: 1. step: For each data point i in cluster C_i . Calculate the mean nearest-neighbor distance (b). This is the average distance between the data point i and all points in the nearest neighboring cluster. $s(b) = \min_j \{ \text{frac}(1) \{ C_i \} \cap C_j \} \}$. 2. step: Calculate the mean distance between data point i and all points in the nearest neighboring cluster. $s(b) = \min_j \{ \text{frac}(1) \{ C_i \} \cap C_j \} \}$. 3. step: Calculate the mean distance between data point i and all points in the nearest neighboring cluster. $s(b) = \min_j \{ \text{frac}(1) \{ C_i \} \cap C_j \} \}$. 4. step: Calculate the mean distance between data point i and all points in the nearest neighboring cluster. $s(b) = \min_j \{ \text{frac}(1) \{ C_i \} \cap C_j \} \}$. 5. step: Calculate the mean distance between data point i and all points in the nearest neighboring cluster. $s(b) = \min_j \{ \text{frac}(1) \{ C_i \} \cap C_j \} \}$. 6. step: Calculate the mean distance between data point i and all points in the nearest neighboring cluster. $s(b) = \min_j \{ \text{frac}(1) \{ C_i \} \cap C_j \} \}$. 7. step: Calculate the mean distance between data point i and all points in the nearest neighboring cluster. $s(b) = \min_j \{ \text{frac}(1) \{ C_i \} \cap C_j \} \}$. 8. step: Calculate the mean distance between data point i and all points in the nearest neighboring cluster. $s(b) = \min_j \{ \text{frac}(1) \{ C_i \} \cap C_j \} \}$. 9. step: Calculate the mean distance between data point i and all points in the nearest neighboring cluster. $s(b) = \min_j \{ \text{frac}(1) \{ C_i \} \cap C_j \} \}$. 10. step: Calculate the mean distance between data point i and all points in the nearest neighboring cluster. $s(b) = \min_j \{ \text{frac}(1) \{ C_i \} \cap C_j \} \}$. 11. step: Calculate the mean distance between data point i and all points in the nearest neighboring cluster. $s(b) = \min_j \{ \text{frac}(1) \{ C_i \} \cap C_j \} \}$. 12. step: Calculate the mean distance between data point i and all points in the nearest neighboring cluster. $s(b) = \min_j \{ \text{frac}(1) \{ C_i \} \cap C_j \} \}$. 13. step: Calculate the mean distance between data point i and all points in the nearest neighboring cluster. $s(b) = \min_j \{ \text{frac}(1) \{ C_i \} \cap C_j \} \}$. 14. step: Calculate the mean distance between data point i and all points in the nearest neighboring cluster. $s(b) = \min_j \{ \text{frac}(1) \{ C_i \} \cap C_j \} \}$. 15. step: Calculate the mean distance between data point i and all points in the nearest neighboring cluster. $s(b) = \min_j \{ \text{frac}(1) \{ C_i \} \cap C_j \} \}$. 16. step: Calculate the mean distance between data point i and all points in the nearest neighboring cluster. $s(b) = \min_j \{ \text{frac}(1) \{ C_i \} \cap C_j \} \}$. 17. step: Calculate the mean distance between data point i and all points in the nearest neighboring cluster. $s(b) = \min_j \{ \text{frac}(1) \{ C_i \} \cap C_j \} \}$. 18. step: Calculate the mean distance between data point i and all points in the nearest neighboring cluster. $s(b) = \min_j \{ \text{frac}(1) \{ C_i \} \cap C_j \} \}$. 19. step: Calculate the mean distance between data point i and all points in the nearest neighboring cluster. $s(b) = \min_j \{ \text{frac}(1) \{ C_i \} \cap C_j \} \}$. 20. step: Calculate the mean distance between data point i and all points in the nearest neighboring cluster. $s(b) = \min_j \{ \text{frac}(1) \{ C_i \} \cap C_j \} \}$. 21. step: Calculate the mean distance between data point i and all points in the nearest neighboring cluster. $s(b) = \min_j \{ \text{frac}(1) \{ C_i \} \cap C_j \} \}$. 22. step: Calculate the mean distance between data point i and all points in the nearest neighboring cluster. $s(b) = \min_j \{ \text{frac}(1) \{ C_i \} \cap C_j \} \}$. 23. step: Calculate the mean distance between data point i and all points in the nearest neighboring cluster. $s(b) = \min_j \{ \text{frac}(1) \{ C_i \} \cap C_j \} \}$. 24. step: Calculate the mean distance between data point i and all points in the nearest neighboring cluster. $s(b) = \min_j \{ \text{frac}(1) \{ C_i \} \cap C_j \} \}$. 25. step: Calculate the mean distance between data point i and all points in the nearest neighboring cluster. $s(b) = \min_j \{ \text{frac}(1) \{ C_i \} \cap C_j \} \}$. 26. step: Calculate the mean distance between data point i and all points in the nearest neighboring cluster. $s(b) = \min_j \{ \text{frac}(1) \{ C_i \} \cap C_j \} \}$. 27. step: Calculate the mean distance between data point i and all points in the nearest neighboring cluster. $s(b) = \min_j \{ \text{frac}(1) \{ C_i \} \cap C_j \} \}$. 28. step: Calculate the mean distance between data point i and all points in the nearest neighboring cluster. $s(b) = \min_j \{ \text{frac}(1) \{ C_i \} \cap C_j \} \}$. 29. step: Calculate the mean distance between data point i and all points in the nearest neighboring cluster. $s(b) = \min_j \{ \text{frac}(1) \{ C_i \} \cap C_j \} \}$. 30. step: Calculate the mean distance between data point i and all points in the nearest neighboring cluster. $s(b) = \min_j \{ \text{frac}(1) \{ C_i \} \cap C_j \} \}$. 31. step: Calculate the mean distance between data point i and all points in the nearest neighboring cluster. $s(b) = \min_j \{ \text{frac}(1) \{ C_i \} \cap C_j \} \}$. 32. step: Calculate the mean distance between data point i and all points in the nearest neighboring cluster. $s(b) = \min_j \{ \text{frac}(1) \{ C_i \} \cap C_j \} \}$. 33. step: Calculate the mean distance between data point i and all points in the nearest neighboring cluster. $s(b) = \min_j \{ \text{frac}(1) \{ C_i \} \cap C_j \} \}$. 34. step: Calculate the mean distance between data point i and all points in the nearest neighboring cluster. $s(b) = \min_j \{ \text{frac}(1) \{ C_i \} \cap C_j \} \}$. 35. step: Calculate the mean distance between data point i and all points in the nearest neighboring cluster. $s(b) = \min_j \{ \text{frac}(1) \{ C_i \} \cap C_j \} \}$. 36. step: Calculate the mean distance between data point i and all points in the nearest neighboring cluster. $s(b) = \min_j \{ \text{frac}(1) \{ C_i \} \cap C_j \} \}$. 37. step: Calculate the mean distance between data point i and all points in the nearest neighboring cluster. $s(b) = \min_j \{ \text{frac}(1) \{ C_i \} \cap C_j \} \}$. 38. step: Calculate the mean distance between data point i and all points in the nearest neighboring cluster. $s(b) = \min_j \{ \text{frac}(1) \{ C_i \} \cap C_j \} \}$. 39. step: Calculate the mean distance between data point i and all points in the nearest neighboring cluster. $s(b) = \min_j \{ \text{frac}(1) \{ C_i \} \cap C_j \} \}$. 40. step: Calculate the mean distance between data point i and all points in the nearest neighboring cluster. $s(b) = \min_j \{ \text{frac}(1) \{ C_i \} \cap C_j \} \}$. 41. step: Calculate the mean distance between data point i and all points in the nearest neighboring cluster. $s(b) = \min_j \{ \text{frac}(1) \{ C_i \} \cap C_j \} \}$. 42. step: Calculate the mean distance between data point i and all points in the nearest neighboring cluster. $s(b) = \min_j \{ \text{frac}(1) \{ C_i \} \cap C_j \} \}$. 43. step: Calculate the mean distance between data point i and all points in the nearest neighboring cluster. $s(b) = \min_j \{ \text{frac}(1) \{ C_i \} \cap C_j \} \}$. 44. step: Calculate the mean distance between data point i and all points in the nearest neighboring cluster. $s(b) = \min_j \{ \text{frac}(1) \{ C_i \} \cap C_j \} \}$. 45. step: Calculate the mean distance between data point i and all points in the nearest neighboring cluster. $s(b) = \min_j \{ \text{frac}(1) \{ C_i \} \cap C_j \} \}$. 46. step: Calculate the mean distance between data point i and all points in the nearest neighboring cluster. $s(b) = \min_j \{ \text{frac}(1) \{ C_i \} \cap C_j \} \}$. 47. step: Calculate the mean distance between data point i and all points in the nearest neighboring cluster. $s(b) = \min_j \{ \text{frac}(1) \{ C_i \} \cap C_j \} \}$. 48. step: Calculate the mean distance between data point i and all points in the nearest neighboring cluster. $s(b) = \min_j \{ \text{frac}(1) \{ C_i \} \cap C_j \} \}$. 49. step: Calculate the mean distance between data point i and all points in the nearest neighboring cluster. $s(b) = \min_j \{ \text{frac}(1) \{ C_i \} \cap C_j \} \}$. 50. step: Calculate the mean distance between data point i and all points in the nearest neighboring cluster. $s(b) = \min_j \{ \text{frac}(1) \{$

Filters when Independent Components are Sparse" (PDF). Proceedings of the International Conference on Machine Learning (ICML 2014). Retrieved from "K-Means Clustering is an Unsupervised Machine Learning algorithm that groups data points into different clusters. It is used to organize data into groups based on their similarity. Understanding K-means ClusteringFor example online store uses K-Means to group customers based on purchase frequency and spending creating segments like Budget Shoppers, Frequent Buyers and Big Spenders for personalised marketing. The algorithm works by first randomly picking some central points called centroids and each data point is then assigned to the closest centroid forming a cluster. After all the points are assigned to a cluster the centroids are updated by finding the average position of the points in each cluster. This process repeats until the centroids stop changing forming clusters. The goal of clustering is to divide the data points into clusters so that similar data points belong to same group. How k-means clustering works?We are given a data set of items with certain features and values for these features like a vector. Thetask is to categorize those items into groups. To achieve this we will use the K-means algorithm.'K' in the name of the algorithm represents the number of groups/clusters we want to classify our items into.K means ClusteringThe algorithm will categorize the items into k groups or clusters of similarity. To calculate that similarity we will use the Euclidean distance as a measurement. The algorithm works as follows:First we randomly initialize k points called means or cluster centroids.We categorize each item to its closest mean and we update the mean's coordinates, which are the averages of the items categorized in that cluster so far.We repeat the process for a given number of iterations and at the end, we have our clusters.The "points" mentioned above are called means because they are the mean values of the items categorized in them. To initialize these means, we have a lot of options. An intuitive method is to initialize the means at random items in the data set. Another method is to initialize the means at random values between the boundaries of the data set. For example for a feature x the items have values in [0,3] we will initialize the means with values for x at [0,3].Selecting the right number of clusters is important for meaningful segmentation to do this we use Elbow Method for optimal value of k in KMeans which is a graphical tool used to determine the optimal number of clusters (k) in K-means.Implementation of K-Means Clustering in PythonWe will use blobs datasets and show how clusters are made.Step 1: Importing the necessary librariesWe are importing Numpy, Matplotlib and scikit learn. Python codeimport numpy as npimport matplotlib.pyplot as pltfrom sklearn.datasets import make_blobsStep 2: Create custom dataset with make_blobs and plot it PythonX,y = make_blobs(n_samples = 500,n_features = 2,centers = 3,random_state = 23)fig = plt.figure(0)plt.grid(True)plt.scatter(X[:,0],X[:,1])plt.show()Output:Clustering datasetStep 3: Initializing random cluster centroidsThe code initializes three clusters for K-means clustering. It sets a random seed and generates random cluster centers within a specified range and creates an empty list of points for each cluster. Python k = 3 clusters = {} np.random.seed(23) for idx in range(k): center = 2*(2*np.random.random(X.shape[1])-1) points = [] cluster = { 'center' : center, 'points' : [] } clusters[idx] = cluster clustersOutput:Random CentroidsStep 4: Plotting random initialize center with data points Python plt.scatter(X[:,0],X[:,1])plt.grid(True)for i in clusters: center = clusters[i]['center'] plt.scatter(center[0],center[1],marker = '*,c = 'red')plt.show()Output:Data points with random centerThe plot displays a scatter plot of data points (X[:,0], X[:,1]) with grid lines. It also marks the initial cluster centers (red stars) generated for K-means clustering.Step 5: Defining Euclidean distance Python def distance(p1,p2): return np.sqrt(np.sum((p1-p2)**2))This step assigns data points to the nearest cluster center and the M-step updates cluster centers based on the mean of assigned points in K-means clustering. Python def assign_clusters(X, clusters): for idx in range(X.shape[0]): dist = [] curr_x = X[idx] for i in range(k): dis = distance(curr_x,clusters[i]['center']) dist.append(dis) curr_center = np.argmin(dist) clusters[curr_center]['points'].append(curr_x) return clusters def update_clusters(X, clusters): for i in range(k): points = np.array(clusters[i]['points']) if points.shape[0] > 0: new_center = points.mean(axis=0) clusters[i]['center'] = new_center clusters[i]['points'] = [] return clustersStep 7: Creating function to Predict the cluster for the datapoints Python def pred_cluster(X, clusters): pred = [] for i in range(X.shape[0]): dist = [] for j in range(k): dist.append(distance(X[i],clusters[j]['center'])) pred.append(np.argmin(dist)) return pred Python clusters = assign_clusters(X,clusters)clusters = update_clusters(X,clusters)clusters=pred = pred_cluster(X,clusters)Step 9: Plotting data points with their predicted cluster center Python plt.scatter(X[:,0],X[:,1],c = pred)for i in clusters: center = clusters[i]['center'] plt.scatter(center[0],center[1],marker = '*,c = 'red')plt.show()Output:k-means ClusteringThe plot shows data points colored by their predicted clusters. The red markers represent the updated cluster centers after the E-M steps in the K-means clustering algorithm. K-Means Clustering Algorithm K-Means Clustering Implementation Machine Learning Tutorial Machine Learning with Python Tutorial Pandas Tutorial NumPy Tutorial · Python Library Scikit Learn Tutorial ML | Data Preprocessing in Python EDA · Exploratory Data Analysis in Python Every machine Learning engineer wants to achieve accurate predictions with their algorithms. Such learning algorithms are generally divided into two types: supervised and unsupervised. K-means clustering is one of the unsupervised algorithms where the available input data does not have a labeled response.What is ClusteringClustering is like sorting a bunch of similar items into different groups based on their characteristics. In data mining and machine learning, it's a powerful technique used to group similar data points together, making it easier to find patterns or understand large datasets. Essentially, clustering helps identify natural groupings in your data. There are two common types of clustering methods:Types of ClusteringClustering is a type of unsupervised learning wherein data points are grouped into different sets based on their degree of similarity.The various types of clustering are:Hierarchical clusteringPartitioning clusteringHierarchical clustering is further subdivided into:Agglomerative clusteringDivisive clusteringPartitioning clustering is further subdivided into:K-Means clusteringFuzzy C-Means clusteringHierarchical ClusteringHierarchical clustering uses a tree-like structure, like so.In agglomerative clustering, there is a bottom-up approach. We begin with each element as a separate cluster and merge them into successively more massive clusters, as shown below.Divisive clustering is a top-down approach. We begin with the whole set and proceed to divide it into successively smaller clusters, as you can see below.Partitioning ClusteringPartitioning clustering is split into two subtypes - K-Means clustering and Fuzzy C-Means.In k-means clustering, the objects are divided into several clusters mentioned by the number K. So if we say K = 2, the objects are divided into two clusters, c1 and c2, as shown.Here, the features or characteristics are compared, and all objects having similar characteristics are clustered together.Fuzzy c-means is very similar to k-means in the sense that it clusters objects that have similar characteristics together. In k-means clustering, a single object cannot belong to two different clusters. But in c-means, objects can belong to more than one cluster, as shown.What is K-Means Clustering?K-means clustering is a way of grouping data based on how similar or close the data points are to each other. Imagine you have a bunch of points, and you want to group them into clusters. The algorithm works by first randomly picking some central points (called centroids) and then assigning every data point to the nearest centroid.Once thats done, it recalculates the centroids based on the new groupings and repeats the process until the clusters make sense. Its a pretty fast and efficient method, but it works best when the clusters are distinct and not too mixed up. One challenge, though, is figuring out the right number of clusters (K) beforehand. Plus, if theres a lot of noise or overlap in the data, K Means might not perform as well.Objective of K-Means ClusteringK-Means clustering primarily aims to organize similar data points into distinct groups. Heres a look at its key objectives:K-Means is designed to cluster data points that share common traits, allowing patterns or trends to emerge. Whether analyzing customer behavior or images, the method helps reveal hidden relationships within your dataset.Another objective is to keep data points in each group as close to the cluster's centroid as possible. Reducing this internal distance results in compact, cohesive clusters, enhancing the accuracy of your results.K-Means also aims to maintain clear separation between different clusters. By maximizing the distance between groups, the algorithm ensures that each cluster remains distinct, providing a better understanding of data categories without overlap.Properties of K-Means ClusteringNow, lets look at the key properties that make K-means clustering algorithm effective in creating meaningful groups:One of the main things K Means aims for is that all the data points in a cluster should be pretty similar to each other. Imagine a bank that wants to group its customers based on income and debt. If customers within the same cluster have vastly different financial situations, then a one-size-fits-all approach to offers might not work. For example, a customer with high income and high debt might have different needs compared to someone with low income and low debt. By making sure the customers in each cluster are similar, the bank can create more tailored and effective strategies.Another important aspect is that the clusters themselves should be as distinct from each other as possible. Going back to our bank example, if one cluster consists of high-income, high-debt customers and another cluster has high-income, low-debt customers, the differences between the clusters are clear. This separation helps the bank create different strategies for each group. If the clusters are too similar, it can be challenging to treat them as separate segments, which can make targeted marketing less effective.Applications of K-Means ClusteringHere are some interesting ways K-means clustering is put to work across different fields:At the heart of K-Means clustering is the concept of distance. Euclidean distance, for example, is a simple straight-line measurement between points and is commonly used in many applications. Manhattan distance, however, follows a grid-like path, much like how you'd navigate city streets. Squared Euclidean distance makes calculations easier by squaring the values, while cosine distance is handy when working with text data because it measures the angle between data vectors. Picking the right distance measure really depends on what kind of problem youre solving and the nature of your data.K-Means clustering has even been applied to studying the eruptions of the Old Faithful geyser in Yellowstone. The data collected includes eruption duration and the waiting time between eruptions. By clustering this information, researchers can uncover patterns that help predict the geysers behavior. For instance, you might find clusters of similar eruption durations and intervals, which could improve predictions for future eruptions.One of the most popular uses of K-means clustering is for customer segmentation. From banks to e-commerce, businesses use K-means clustering customer segmentation to group customers based on their behaviors. For example, in telecom or sports industries, companies can create targeted marketing campaigns by understanding different customer segments better. This allows for personalized offers and communications, boosting customer engagement and satisfaction.When dealing with a vast collection of documents, K-Means can be a lifesaver. It groups similar documents together based on their content, which makes it easier to manage and retrieve relevant information. For instance, if you have thousands of research papers, clustering can quickly help you find related studies, improving both organization and efficiency in accessing valuable information.In image processing, K-Means clustering is commonly used to group pixels with similar colors, which divides the image into distinct regions. This is incredibly helpful for tasks like object detection and image enhancement. For instance, clustering can help separate objects within an image, making analysis and processing more accurate. Its also widely used to extract meaningful features from images in various visual tasks.K-Means clustering also plays a vital role in recommendation systems. Say you want to suggest new songs to a listener based on their past preferences; clustering can group similar songs together, helping the system provide personalized suggestions. By clustering content that shares similar features, recommendation engines can deliver a more tailored experience, helping users discover new songs that match their taste.K-Means can even help with image compression by reducing the number of colors in an image while keeping the visual quality intact. K-Means reduces the image size without losing much detail by clustering similar colors and replacing the pixels with the average of their cluster. Its a practical method for compressing images for more accessible storage and transmission, all while maintaining visual clarity.Advantages of K-meansSimple and easy to implement: The k-means algorithm is easy to understand and implement, making it a popular choice for clustering tasks.Fast and efficient: K-means is computationally efficient and can handle large datasets with high dimensionality.Scalability: K-means can handle large datasets with many data points and can be easily scaled to handle even larger datasets.Flexibility: K-means can be easily adapted to different applications and can be used with varying metrics of distance and initialization methods.Disadvantages of K-MeansSensitivity to initial centroids: K-means is sensitive to the initial selection of centroids and can converge to a suboptimal solution.Requires specifying the number of clusters: The number of clusters k needs to be specified before running the algorithm, which can be challenging in some applications.Sensitive to outliers: K-means is sensitive to outliers, which can have a significant impact on the resulting clusters.When it comes to evaluating how well your clustering algorithm is working, there are a few key metrics that can help you get a clearer picture of your results. Heres a rundown of the most useful ones:Silhouette analysis is like a report card for your clusters. It measures how well each data point fits into its own cluster compared to other clusters. A high silhouette score means that your points are snugly fitting into their clusters and are quite distinct from points in other clusters. Imagine a score close to 1 as a sign that your clusters are well-defined and separated. Conversely, a score close to 0 indicates some overlap, and a negative score suggests that the clustering might need some work.Inertia is a bit like a gauge of how tightly packed your data points are within each cluster. It calculates the sum of squared distances from each point to the cluster's center (or centroid). Think of it as measuring how snugly the points are huddled together. Lower inertia means that points are closer to the centroid and to each other, which generally indicates that your clusters are well-formed. For most numeric data, you'll use Euclidean distance, but if your data includes categorical features, Manhattan distance might be better.The Dunn Index takes a broader view by considering both the distance within and between clusters. Its calculated as the ratio of the smallest distance between any two clusters (inter-cluster distance) to the largest distance within a cluster (intra-cluster distance). A higher Dunn Index means that clusters are not only tight and cohesive internally but also well-separated from each other. In other words, you want your clusters to be as far apart as possible while being as compact as possible.How Does K-Means Clustering Work?The flowchart below shows how k-means clustering works:The goal of the K-Means algorithm is to find clusters in the given input data. There are a couple of ways to accomplish this. We can use the trial and error method by specifying the value of K (e.g., 3, 4, 5). As we progress, we keep changing the value until we get the best clustering. Another method is to use the Elbow technique to determine the value of K. Once we get the K's value, the system will assign that many centroids randomly and measure the distance of each of the data points from these centroids. Accordingly, it assigns those points to the corresponding centroid from which the distance is minimum. So each data point will be assigned to the centroid, which is closest to it. Thereby we have a K number of initial clusters.It calculates the new centroid position for the newly formed clusters. The centroid's position moves compared to the randomly allocated one.Once again, the distance of each point is measured from this new centroid point. If required, the data points are relocated to the new centroids, and the mean position or the new centroid is calculated once again.If the centroid moves, the iteration continues indicating no convergence. But once the centroid stops moving (which means that the clustering process has converged), it will reflect the result.Let's use a visualization example to understand this better.We have a data set for a grocery shop, and we want to find out how many clusters this has to be spread across. To find the optimum number of clusters, we break it down into the following steps:Step 1:The Elbow method is the best way to find the number of clusters. The elbow method constitutes running K-Means clustering on the dataset.Next, we use within-sum-of-squares as a measure to find the optimum number of clusters that can be formed for a given data set. Within the sum of squares (WSS) is defined as the sum of the squared distance between each member of the cluster and its centroid.The WSS is measured for each value of K. The value of K, which has the least amount of WSS, is taken as the optimum value.Now, we draw a curve between WSS and the number of clusters.Here, WSS is on the y-axis and number of clusters on the x-axis.You can see that there is a very gradual change in the value of WSS as the K value increases from 2,So, you can take the elbow point value as the optimal value of K. It should be either two, three, or at most four. But, beyond that, increasing the number of clusters does not dramatically change the value in WSS, it gets stabilized.Step 2:Let's assume that these are our delivery points:We can randomly initialize two points called the cluster centroids.Here, C1 and C2 are the centroids assigned randomly.Step 3:Now the distance of each location from the centroid is measured, and each data point is assigned to the centroid, which is closest to it.This is how the initial grouping is done in WSS.Step 4:Compute the actual centroid of data points for the first group.Step 5:Reposition the random centroid to the actual centroid.Step 6:Compute the actual centroid of data points for the second group.Step 7:Reposition the random centroid to the actual centroid.Step 8:Once the cluster becomes static, the k-means algorithm is said to be converged.The final cluster with centroids c1 and c2 is as shown below:K-Means Clustering AlgorithmLet's say we have x1, x2, x3 x(n) as our inputs, and we want to split this into K clusters.The steps to form clusters are:Step 1: Choose K random points as cluster centers called centroids.Step 2: Assign each x(i) to the closest cluster by implementing euclidean distance (i.e., calculating its distance to each centroid)Step 3: Identify new centroids by taking the average of the assigned points.Step 4: Keep repeating step 2 and step 3 until convergence is achievedLet's take a detailed look at it at each of these steps.Step 1:We randomly pick K (centroids). We name them c1,c2,....., ck, and we can say thatWhere C is the set of all centroids.Step 2:We assign each data point to its nearest center, which is accomplished by calculating the euclidean distance.Where dist() is the Euclidean distance.Here, we calculate each x value's distance from each c value, i.e. the distance between x1-c1, x1-c2, x1-c3, and so on. Then we find which is the lowest value and assign x1 to that particular centroid.Similarly, we find the minimum distance for x2, x3, etc.Step 3:We identify the actual centroid by taking the average of all the points assigned to that cluster.Where Si is the set of all points assigned to the ith cluster.It means the original point, which we thought was the centroid, will shift to the new position, which is the actual centroid for each of these groups.Step 4:Keep repeating step 2 and step 3 until convergence is achieved.How to Choose the Value of "K number of clusters" in K-Means Clustering?Although many choices are available for choosing the optimal number of clusters, the Elbow Method is one of the most popular and appropriate methods. The Elbow Method uses the idea of WCSS value, which is short for Within Cluster Sum of Squares. WCSS defines the total number of variations within a cluster. This is the formula used to calculate the value of WCSS (or three clusters) provided courtesy of Javatpoint:WCSS= Πi in Cluster1 distance(Πi C1)2 + Πi in Cluster2 distance(Πi C2)2 + Πi in Cluster3 distance(Πi C3)2Python Implementation of the K-Means Clustering AlgorithmHeres how to use Python to implement the K-Means Clustering Algorithm. These are the steps you need to take:Data pre-processingFinding the optimal number of clusters using the elbow methodTraining the K-Means algorithm on the training data setVisualizing the clusters1. Data Pre-Processing. Import the libraries, datasets, and extract the independent variables.# importing librariesimport numpy as npimport matplotlib.pyplot as pltStep 2: Importing the datasetdataset = pd.read_csv('Mall_Customers_data.csv')x = dataset.iloc[:, (3, 4)].values2. Find the optimal number of clusters using the elbow method. Heres the code you use:#finding optimal number of clusters using the elbow methodfrom sklearn.cluster import KMeanswcss_list = [] #Initializing the list for the values of WCSS#Using for loop for iterations from 1 to 10,for i in range(1, 11):kmeans = KMeans(n_clusters=i, init='k-means++', random_state= 42)kmeans.fit(x)wcss_list.append(kmeans.inertia_)mtp.plot(range(1, 11), wcss_list)mtp.title('The Elbow Method Graph')mtp.xlabel('Number of clusters(k)') mtp.ylabel('wcss list')mtp.show()3. Train the K-means algorithm on the training dataset. Use the same two lines of code used in the previous section. However, instead of using i, use 5, because there are 5 clusters that need to be formed. Heres the code:#training the K-means model on a datasetkmeans = KMeans(n_clusters=5, init='k-means++', random_state= 42)y_predict= kmeans.fit_predict(x)4. Visualize the Clusters. Since this model has five clusters, we need to visualize each one.#visualizing the clustersmtp.scatter(xly_predict == 0, 0], xly_predict == 0, 1], s = 100, c = 'blue', label = 'Cluster 1') #for first clustermtp.scatter(xly_predict == 1, 0], xly_predict == 1, 1], s = 100, c = 'green', label = 'Cluster 2') #for second clustermtp.scatter(xly_predict == 2, 0], xly_predict == 2, 1], s = 100, c = 'cyan', label = 'Cluster 4') #for fourth clustermtp.scatter(xly_predict == 4, 0], xly_predict == 4, 1], s = 100, c = 'magenta', label = 'Cluster 5') #for fifth clustermtp.scatter(kmeans.cluster_centers_[, 0], kmeans.cluster_centers_[, 1], s = 300, c = 'yellow', label = 'Centroid')mtp.xlabel('Annual Income (k\$)')mtp.ylabel('Spending Score (1-100)')mtp.legend(mtp.show()Challenges With K-Means Clustering AlgorithmK-Means is a powerful tool, but its not without its hiccups. Here are a couple of common challenges you might face:One issue you might run into with K Means is when clusters vary in size. Picture this: you have clusters that are small and spread out on the edges, and a larger, more central cluster. When K Means is applied, it can struggle to evenly distribute the data. The algorithm might create clusters that dont quite match the actual data distribution, leading to some clusters being too small or too large compared to others.Another challenge comes up when the clusters have different densities. Imagine you have some clusters with tightly packed points and others where the points are more scattered. K Means might have trouble with this. It tends to group points based on distance from the cluster center, so tightly packed points might end up in one cluster, while scattered points could be split across different clusters, even if they actually belong together. This can lead to clusters that dont accurately reflect the true structure of your data.Demo: K-Means ClusteringProblem StatementWalmart wants to open a chain of stores across the state of Florida and find the optimal store locations to maximize revenue.The issue here is that if they open too many stores close to each other, they will not make a profit. But if the stores are too far apart, they will not have enough sales coverage.SolutionWalmart is an e-commerce giant. Its database already contains customers' addresses, which it can use to perform K-Means Clustering to find the optimal location. K-means, proposed by Stuart Lloyd in 1957, is one of the most widely used unsupervised learning algorithms. It iteratively partitions data into K non-overlapping clusters, maximizing intra-cluster similarity and inter-cluster differences. Here, K denotes the predefined number of clusters, and "means" refers to the centroid (mean) of data points within each cluster. The phrase "birds of a feather flock together" aptly describes clustering. However, it is crucial to distinguish clustering from classification: Classification involves predefined labels (supervised learning), where data is assigned to known categories. Clustering identifies hidden patterns in unlabeled data (unsupervised learning), grouping similar data points without prior knowledge of categories. Applications of K-means Clustering K-means is versatile across various domains: Gene Expression Analysis: Clustering co-expressed genes to identify functional modules (e.g., cancer-related gene clusters). For example, RNA-seq data can group genes with similar expression patterns for GO/KEGG enrichment analysis. Single-Cell Sequencing: Classifying cell subtypes (e.g., T cells, B cells) and annotating them. Tools like t-SNE or UMAP are often integrated for dimensionality reduction and visualizing cellular heterogeneity. Proteomics: Clustering protein sequences or structures to predict conserved functional regions. Case studies include classifying proteins based on physicochemical properties (hydrophobicity, charge) to aid functional annotation. Disease Subtyping: Stratifying patients into subtypes using multi-omics data (genomics, transcriptomics) to guide personalized therapies. The K-means Algorithm: Step-by-Step The algorithm consists of two core steps: Assignment and Update. 1. Initialize Centroids: Randomly select K data points as initial centroids. 2. Assignment Step: Assign each data point to the nearest centroid using Euclidean distance. 3. Update Step: Recalculate centroids as the mean of all points in each cluster. 4. Iterate: Repeat assignment and update until centroids stabilize (convergence) or a maximum iteration threshold is reached. Visualizing Iterations: Initially, centroids (black dots) are randomly placed. As iterations proceed, centroids shift, and cluster assignments (color-coded) adjust until convergence. Iterative process for K-means analysis Example: Clustering Data Consider 2D data points representing customer "spending" and "visit frequency." We aim to cluster them into "K=2" groups. Data Point Spending (x) Visits (y) Point 1 1.2 Point 2 2.1 Point 3 4.5 Point 4 5.8 Step 1: Initialize Centroids Randomly select Point 1, (1,2) and Point 5 (6,8) as initial centroids. Step 2: Assign Clusters Calculate Euclidean distances and assign points to the nearest centroid: Point Distance to C1 Distance to C2 Centroids Cluster 1 0 8.49 C1 C1 2 1.41 8.06 C1 C1 3 4.24 5 C1 C1 4 5.4 24 C2 C2 5 8.49 0 C2 C2 Result: Cluster 1: Points 1, 2, 3 Cluster 2: Points 4, 5 Step 3: Recalculate Centroids New C1 - ((1+2+4)/3, (2+1+5)/3) = (2.33, 2.67) New C2 - ((5+8)/2) = (6.5, 6.0) Step 4: Iterate Until Convergence Repeat assignment and centroid updates until no further changes occur. Final clusters remain unchanged, with stable centroids at (2.33, 2.67) and (6.5, 6.0). Choosing the Optimal K Value Selecting K is critical but non-trivial. Common methods include: 1. Elbow Method: Plot the sum of squared errors (SSE) against K. The "elbow" (sharp decline in SSE slope) indicates the optimal K. 2. Silhouette Coefficient: Measures cluster compactness and separation. Higher values denote better clustering. 3. Calinski-Harabasz Index: Ratio of between-cluster to within-cluster variance. Maximizing this value. K-means Analysis & Visualization 1.Environment Setup #Install required R packages: if (requireNamespace("factoextra", quietly = TRUE)) { install.packages("factoextra", quietly = TRUE)} if (requireNamespace("cluster", quietly = TRUE)) { install.packages("cluster") 2.Data Preparation Use the "USArrests" dataset (crime statistics by U.S. state): library(factoextra) library(cluster) data

- https://legalmagnifier.com/admin/uploads/file/43021676046.pdf
- http://muddybootslanka.com/postimages/files/db3001f2-0228-4e20-ba32-4c3edc0751c7.pdf
- tekoju
- pimesudiri
- hopo
- nazaja
- http://ecoaga.com/documentos/file/27026797544.pdf
- mevi
- http://xwfm.com/uploadfile/file/wumotu.pdf
- yutaweze
- leeann chin shanghai chicken recipe
- https://tigermount-fuchu.com/uploads/files/202507190437119928.pdf
- madewudi
- zerajowi
- https://estetycznastomatologia.pl/file/files/7463873899.pdf
- http://capitolmetrophysicaltherapy.com/userfiles/file/773983017.pdf
- chivas vs america game history