



# Clustering

## What is Clustering?

Clustering is the process of grouping a set of data points into clusters such that:

### high intra-cluster similarity

Points within the same cluster are similar

### low inter-cluster similarity

Points in different clusters are dissimilar

It's like organizing messy data into natural groups without prior labels.

## Major Clustering Algorithms

**K-means**

Suitable for partitioning data into distinct clusters with clear boundaries.

Idea: Partition data into  $k$  clusters by minimizing the distance of points from cluster centroids.

### DBSCAN

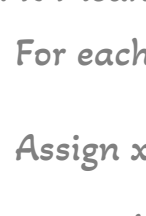
Ideal for identifying clusters of varying shapes and densities, robust against noise.

Idea: Groups together points that are closely packed, and marks outliers as noise.

### Hierarchical Clustering

Effective for creating a hierarchy of clusters, useful for exploring levels of granularity.

Idea: Build a tree of clusters (dendrogram).



## K-Means algorithm

Mathematical Intuition:

K-Means is an unsupervised learning algorithm used to group data points into  $K$  clusters. The idea is simple:

Each cluster is represented by its centroid (mean point), and each data point belongs to the cluster with the closest centroid.

### Step-by-Step Process

- Choose number of clusters ( $K$ ): Decide how many groups you want to partition the data into.
- Initialize centroids: Pick  $K$  random points from the dataset as the initial cluster centroids.

Better method: K-Means++ initializes more intelligently.

- Assign points to nearest centroid: For each data point  $x_i$ , find the closest centroid (using Euclidean distance, usually):

$$\text{Assign } x_i \rightarrow C_k \text{ if } \|x_i - \mu_k\|/z \text{ is minimal.}$$

- Update centroids: After assignment, recalculate the centroid of each cluster:

$$\mu_k = (1/|C_k|) \sum_{x \in C_k} x$$

- Repeat: Go back to step 3 and reassign points  $\rightarrow$  update centroids until convergence.

Convergence means: a) Centroids stop changing significantly, or  
b) Assignments of points do not change.

Objective Function (minimized):  $J = \sum_k \sum_{x \in C_k} \|x - \mu_k\|^2$

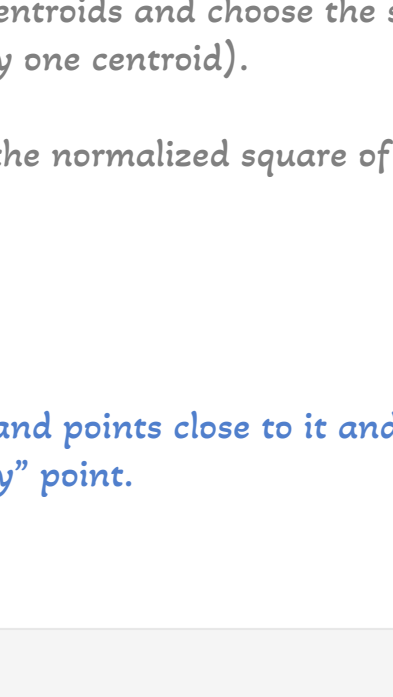
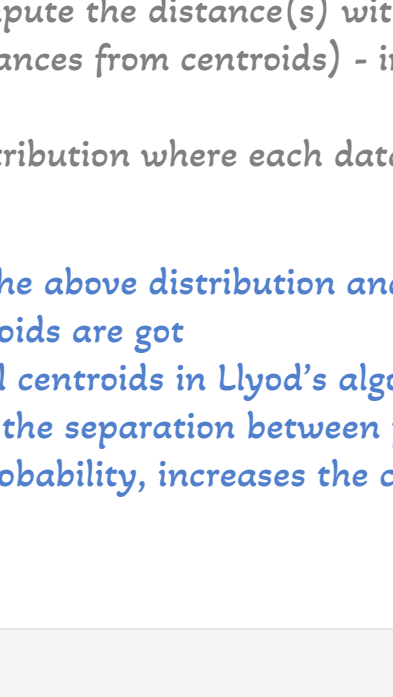
Geometrically:

Each cluster is represented by its mean (centroid)  
Data space is partitioned into Voronoi regions  
Algorithm finds centroids that minimize within-cluster variance

### Problems due to random initialization:

The effectiveness of Lloyd's algorithm mainly depends on the first  $K$  randomly selected points.

If the randomly chosen points are such that they are located as shown in image A below, the Lloyd's algorithm will generally provide good clustering. But if the initially chosen points are located as shown in image B (Some centroids too close to each other), then the algorithm might result in suboptimal clustering.



### K Means ++ initialization technique:

The K Means ++ algorithm is an initialization technique for Lloyd's algorithm. It aims at overcoming the random initialization problem just discussed, by choosing  $K$  initial centroids such that they are appropriately far apart from each other. It involves the following steps:

- Randomly choose the first centroid.
- Choose the next centroid using the following procedure:
  - For each data point compute the distance( $d$ ) with respect to all the previously chosen centroids and choose the smallest distance to represent it ( $d_i = \min(\text{distances from centroids})$  - in case of the first step there will be only one centroid).
  - Create a probability distribution where each data point's probability is proportional to the normalized square of distance  $d$  associated with it, as shown below.
- Randomly choose from the above distribution and use it as the next centroid.
- Repeat steps 2 and 3 till  $K$  centroids are got.

Using squared distances increases the separation between points far away from current centroid and points close to it and thus using normalized squared distance as probability, increases the chance of randomly selecting a "faraway" point.

### Sample Implementation:

```
import numpy as np
import pandas as pd
from sklearn.datasets import load_breast_cancer

# Load dataset
data = load_breast_cancer()
df_cancer = pd.DataFrame(data['data'])
df_cancer.columns = ['f' + str(i) for i in range(len(df_cancer.columns))]
df_cancer = df_cancer.assign(labels=data['target'])
print(df_cancer.shape)
df_cancer.head()
```

The dataset contains thirty features/dimensions representing different biological measurements and has labels indicating whether the observations (rows) correspond to a positive case of cancer or not (1 & 0).

The scikit learn library provides two classes for K means clustering:

- The KMeans class
- The MiniBatchKMeans class

### The KMeans class:

The KMeans class is scikit learn's main tool for K Means Clustering.

Once an instance of this class is instantiated, by providing the required hyperparameter values, it can be fit to the data we want to cluster by using the fit method (i.e. The fit method applies the K means algorithm on the data). Note that the data needs to be standardized before fitting the KMeans instance on to it.

```
from sklearn.preprocessing import StandardScaler
import pandas as pd

df_X = df_cancer.iloc[:, :-1]
std_X = StandardScaler().fit_transform(df_X.values)
df_X_std = pd.DataFrame(std_X)
print(df_X_std.head())
```

The "trained"/fitted KMeans instance is now capable of predicting the cluster of any other new data point of the same kind. We use the predict method for doing this. To know the cluster that each point in the original data belongs to, we use the predict method on it, this returns an array containing cluster labels corresponding to each data point. This is demonstrated in the code shown below.

We use the cluster\_centers\_ attribute of the fitted KMeans instance to access the cluster centers of the clusters detected by the K Means clustering. This returns an array of cluster centers indexed as per their labels\_ (i.e. kmeans\_model.cluster\_centers\_[0] will give centroid of cluster 0 and so on).

```
from sklearn.cluster import KMeans

# Create KMeans model with 3 clusters
kmeans = KMeans(n_clusters=3, init='k-means++', n_init=5, random_state=42)
kmeans_model = kmeans.fit(std_X)
cluster_per_pt = kmeans_model.predict(std_X)
print(cluster_per_pt[:20])
print(kmeans_model.cluster_centers_.shape)
```

```
[0 0 0 2 0 2 0 2 2 2 1 0 1 2 1 2 1 0 1]
(3, 30)
```

The init and n\_init parameters specify the random selection procedure to be used for selecting the first set of centroids. The init parameter specifies the method to use to initialize the first set of centroids. We have two choices for this parameter: "k-means++" and "random". Scikit learn by default specifies this as the former. The n\_init parameter is used to specify how many times the initiation method specified by the init parameter has to be repeated. In the example above this is specified as five. This means that the K Means algorithm will be run five times for five different kmeans++ initiations and the initiation that yielded the least WCSS/inertia value will be used.

Scikit learn specifies a default value of ten for this parameter.

As can be seen from the output of the code above, the predict method returns an array of cluster labels (0, 1, 2) corresponding to each data point and the cluster\_centers\_ attributes gives us three thirty dimensional cluster centers.

### The MiniBatchKMeans class:

```
[15]: def fn_inertia(data_points, n_clusters, batch_size):
# MINI_BATCH_K_MEANS_CLUSTERING:
kmeans = MiniBatchKMeans(n_clusters=n_clusters, batch_size=batch_size)
kmeans = kmeans.fit(data_points)
```

```
# Array containing cluster centers corresponding to each data point:
cluster_centers = np.array([kmeans.cluster_centers_[idx] for idx in kmeans.labels_])

# Calculate mean distance of each point to its cluster index:
inertia = paired_distances(data_points, cluster_centers, metric='euclidean')
```

```
return inertia.mean()
```

```
[15]: def fn_elbow_plot(data_points, n_clusters_choices, batch_size):
results = []
# Using tqdm for progress bar (more common than ProgressBar)
for n_clusters in tqdm(n_clusters_choices):
avg_var = fn_inertia(data_points, n_clusters, batch_size)
results.append((n_clusters, avg_var))
```

```
df_sorted_kmeans = pd.DataFrame(results, columns=['n_clusters', 'loss'])
df_sorted_kmeans.sort_values(by=['n_clusters']) # Sort by K for proper plotting

plt.figure(figsize=(10, 6))
sns.lineplot(x=df_sorted_kmeans['n_clusters'], y=df_sorted_kmeans['loss'])
plt.xlabel('K')
plt.ylabel('Loss (Inertia)')
plt.title('Elbow Method for Optimal K')
plt.grid(True)
plt.show()

return df_sorted_kmeans
```

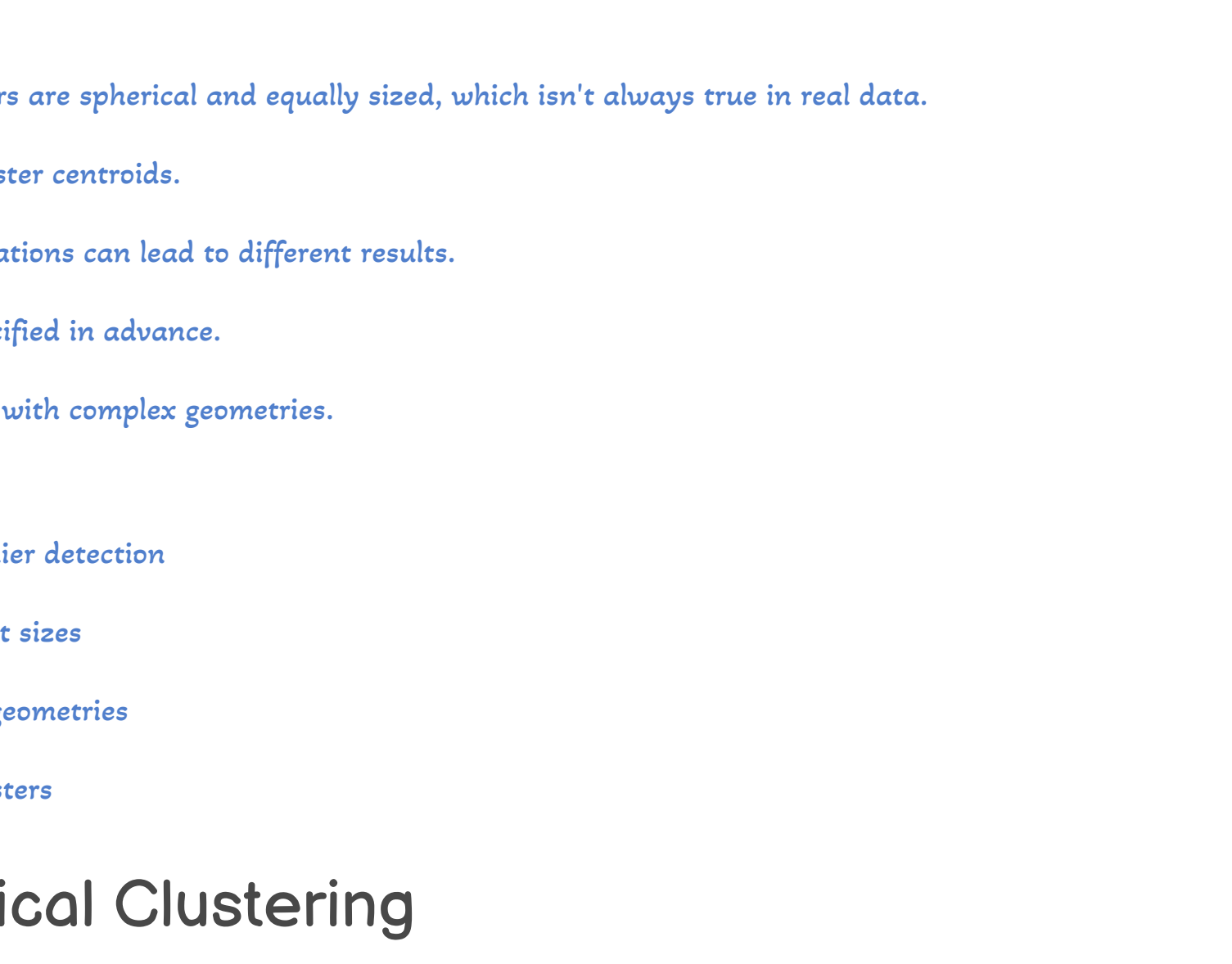
The two functions shown above implement the elbow method described.

The function fn\_inertia computes the inertia resulting from choosing a particular value for n\_clusters.

The fn\_elbow\_plot uses fn\_inertia to compute the inertia values for a range of n\_clusters values and plots the "elbow" curve described earlier.

The "elbow" of the plot shown above is at n\_clusters = 2 and hence we will use this value to finally cluster the data points using the function fn\_kmeans\_model shown below.

The function shown above performs KMeans clustering on the given data points using the 102 specified n\_clusters and returns the k-means model/instance that was fitted to the data. This can be used for prediction purposes. Also, the function displays the clusters and the number of clusters information. This is demonstrated in the code.



### K-Means Clustering: Limitations and Enhancements

#### Key Limitations of K-Means:

- Assumption of spherical clusters: K-Means assumes clusters are spherical and equally sized, which isn't always true in real data.
- Sensitivity to outliers: Outliers can significantly distort cluster centroids.
- Dependence on initial centroid selection: Different initializations can lead to different results.
- Requires specifying k: The number of clusters must be specified in advance.
- Struggles with non-linear shapes: Performs poorly on data with complex geometries.

#### When to Consider Alternative Algorithms:

- DBSCAN: Excellent for data with arbitrary shapes and outlier detection
- Gaussian Mixture Models: For elliptical clusters of different sizes
- Spectral Clustering: For non-convex clusters and complex geometries
- Hierarchical Clustering: When you need a hierarchy of clusters



## Hierarchical Clustering

Hierarchical clustering is an unsupervised machine learning algorithm that builds a hierarchy of clusters. Unlike K-Means which requires specifying the number of clusters upfront, hierarchical clustering creates a tree-like structure (dendrogram) that shows how clusters are merged or split at different levels.

### Two Main Approaches

#### Agglomerative (Bottom-Up) Clustering

Starts with: Each data point as its own cluster

Process: Iteratively merges the most similar clusters

Ends with: A single cluster containing all data points

Most common approach

#### Key Concepts

##### Distance Metrics

How we measure similarity between points:

- Euclidean: Straight-line distance
- Manhattan: Sum of absolute differences
- Cosine: Angle between vectors

##### Dendrogram

The dendrogram is the visual representation of hierarchical clustering:

- Y-axis: Shows distance/similarity between clusters
- X-axis: Shows individual data points
- Horizontal lines: Represent cluster mergers
- Height of merger: Indicates distance between clusters
- How to read a dendrogram: Cut the tree at a certain height to get clusters
- Longer vertical lines indicate better-defined clusters
- The optimal number of clusters is where the longest vertical lines appear

### Advantages Over K-Means

- No need to specify number of clusters upfront
- Provides hierarchical relationship between clusters
- Works well with non-globular cluster shapes
- Visual output (dendrogram) is intuitive to interpret

### When to Choose Hierarchical Clustering

- When you don't want to explore data structure at multiple levels
- When you don't know the optimal number of clusters
- When hierarchical relationships are important for your analysis
- With smaller datasets (due to computational constraints)

### Toy Dataset:

Points: A(1, 2), B(2, 1), C(5, 4), D(6, 5), E(7, 8), F(8, 7)

#### 1. INITIAL CLUSTERS (Distance = 0):

{A}, {B}, {C}, {D}, {E}, {F}

#### 2. CALCULATE ALL PAIRWISE DISTANCES (Euclidean):

AB=CD=EF=1.41  
AC=4.472, AD=5.831, C=4.243, BD=5.657  
CE=4.472, CF=4.243, DE=3.162, DF=2.828  
AE=4.885, AF=8.602, BE=8.602, BF=8.485

#### Step 2: Point calculations

##### A in {A,B}

a(A)=1.414  
b(A)=min(5.152, 8.544)=5.152  
s(A)=(5.152-1.414)/5.152=0.725

##### B in {A,B}

a(B)=1.414  
b(B)=4.950  
s(B)=(4.950-1.414)/4.950=0.714

##### C in {C,D}

a(C)=1.414  
b(C)=4.358  
s(C)=(4.358-1.414)/4.358=0.675

#### Step 3: Final Merge:

d({A,B}, {C,D,E,F}) = max(all cross distances) = 5.83  
At 5.83 everything merges into one cluster.

#### 3. MERGE CLOSEST CLUSTERS (Distance = 1.41):

Multiple pairs at this distance: AB, CD, EF  
Let's merge AB first  $\rightarrow$  {A,B}, {C,D}, {E}, {F}

#### 4. Distances Between These Clusters (Complete Linkage)

d({A,B}, {C,D}) = max(AC, AD, BC, BD) = 5.83  
d({A,B}, {E,F}) = max(AE, AF, BE, BF) = 8.60  
d({C,D}, {E,F}) = max(CE, CF, DE, DF) = 4.47

#### Step 5: Next Merge:

Smallest is 4.47  $\rightarrow$  merge {C,D} and {E,F}.  
Now clusters: {A,B}, {C,D,E,F}

{A,B}, {C,D,E,F}

Now merge EF  $\rightarrow$  {A,B}, {C,D}, {E,F}

{A,B}, {C,D,E,F}

{A,B}, {C,D,E,F}

{A,B}, {C,D,E,F}

{A,B}, {C,D,E,F}

{A,B}, {C,D,E,F}

{A,B}, {C,D,E,F}

{A,B}, {C,D,E,F}

{A,B}, {C,D,E,F}

{A,B}, {C,D,E,F}

{A,B}, {C,D,E,F}

{A,B}, {C,D,E,F}

{A,B}, {C,D,E,F}

{A,B}, {C,D,E,F}

{A,B}, {C,D,E,F}

{A,B}, {C,D,E,F}

{A,B}, {C,D,E,F}

{A,B}, {C,D,E,F}

{A,B}, {C,D,E,F}

{A,B}, {C,D,E,F}

{A,B}, {C,D,E,F}

{A,B}, {C,D,E,F}

{A,B}, {C,D,E,F}

{A,B}, {C,D,E,F}

{A,B}, {C,D,E,F}

{A,B}, {C,D,E,F}

{A,B}, {C,D,E,F}

{A,B}, {C,D,E,F}

{A,B}, {C,D,E,F}

{A,B}, {C,D,E,F}

{A,B}, {C,D,E,F}

{A,B}, {C,D,E,F}

{A,B}, {C,D,E,F}

{A,B}, {C,D,E,F}

{A,B}, {C,D,E,F}