

Clustering

Agenda

Key Takeaways-

- Clustering and its types
- K-Means
- Evaluation measures for K
- K-Means limitations
- Overcoming K-Means limitations

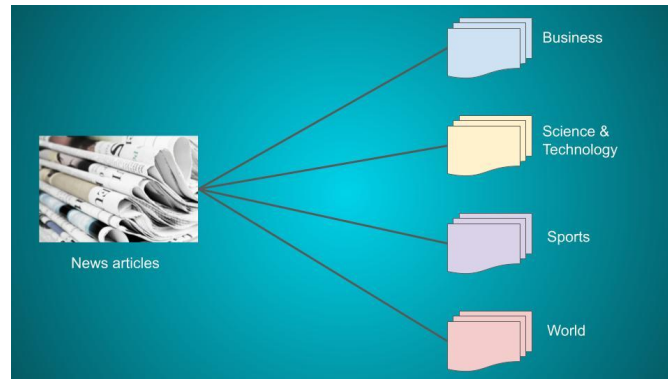
Unsupervised Learning

Unsupervised Learning has no explicit output/target variable, i.e. works without supervision.

So it discovers knowledge, hidden structures or relationship in the unlabelled data. For e.g. it can learn to group or organize data in such a way that similar objects are in the same group.

Similarly, news articles can be put together based on the topics like sports, business, technology, politics etc. This approach is known as Clustering.

Clustering is a technique to group a set of objects in such a way that objects in the same group are much more similar to each other than to those in other groups.



Clustering Use-cases

1. Clustering based on topics

Text documents (such as news articles, white papers, research papers, reports etc.), images are published/available in good quantities so clustering can be used to group or visualize them together based on topics.

2. Text Summarization

Summarizing the text documents to ensure good coverage and avoid redundancy (reducing the size of large documents).

3. Anomaly Detection

Anomaly detection aims to find out the objects that are significantly different from others. E.g. detecting fraudulent transactions in banking and finance.

Types of Clustering Techniques

Most commonly used clustering techniques are -

1. K-Means
2. Hierarchical Clustering
 - Agglomerative
 - Divisive

K-Means

Choose the **k** (number of clusters).

Step 1. **Initialization :**

Randomly select **k data points** as **initial centroids**.

Step 2. **Assignment :**

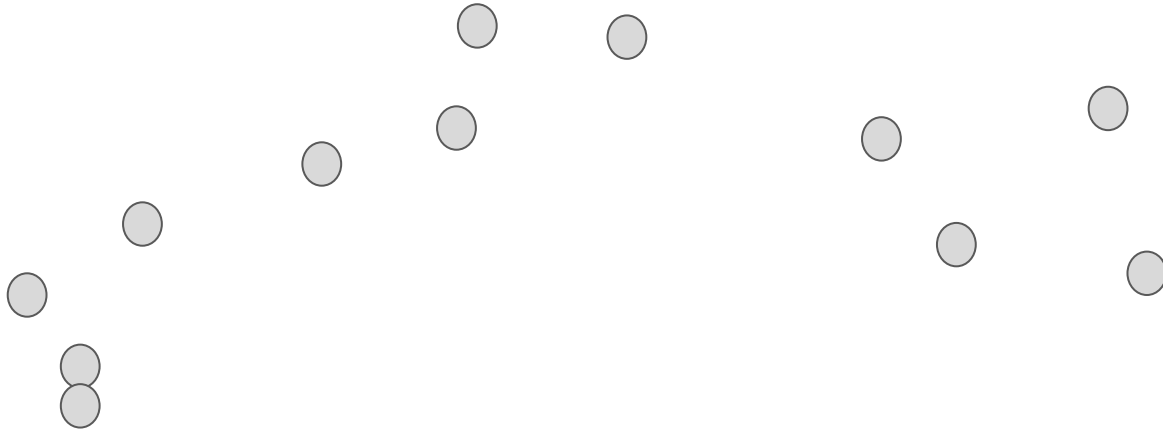
Assign **each data point** to the **closest centroids**, that **forms k clusters**.

Step 3. **Recompute Centroid :**

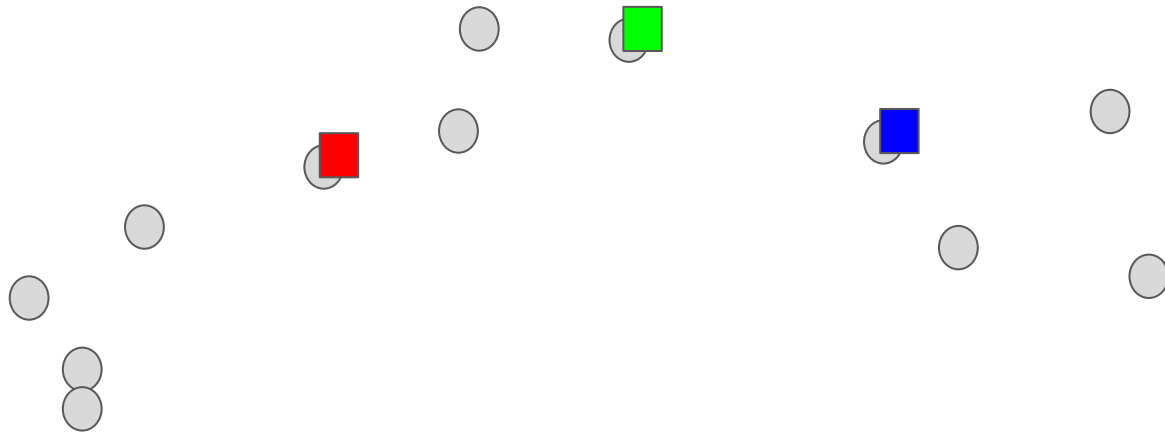
Calculate **new cluster centroid** for **each cluster**.

Step 4. **Repeat Step 2 and 3 until convergence criterion is met.**

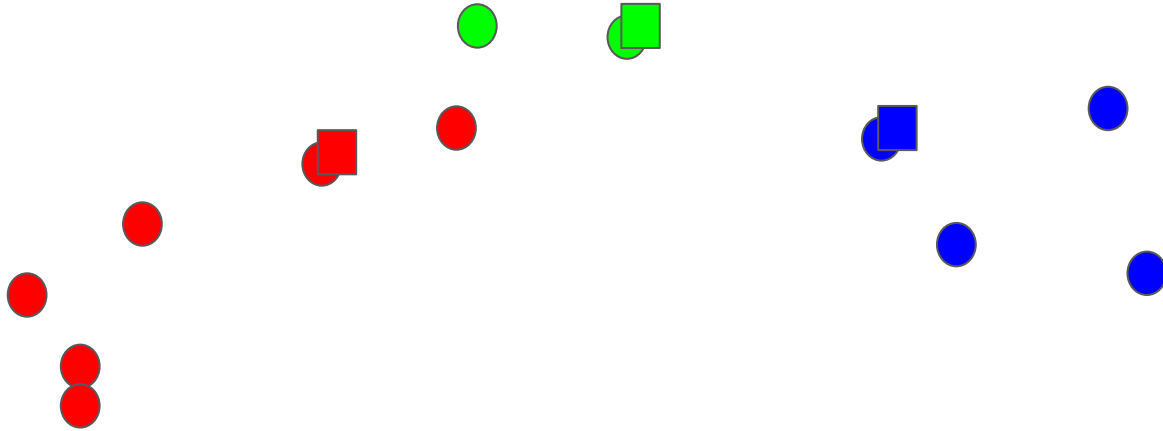
Consider the below data points.



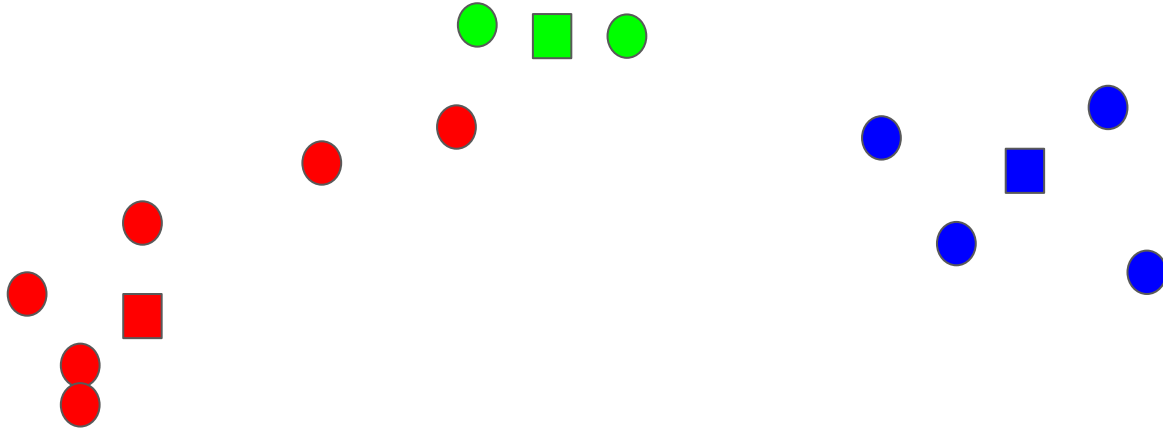
Initialize centroids randomly.



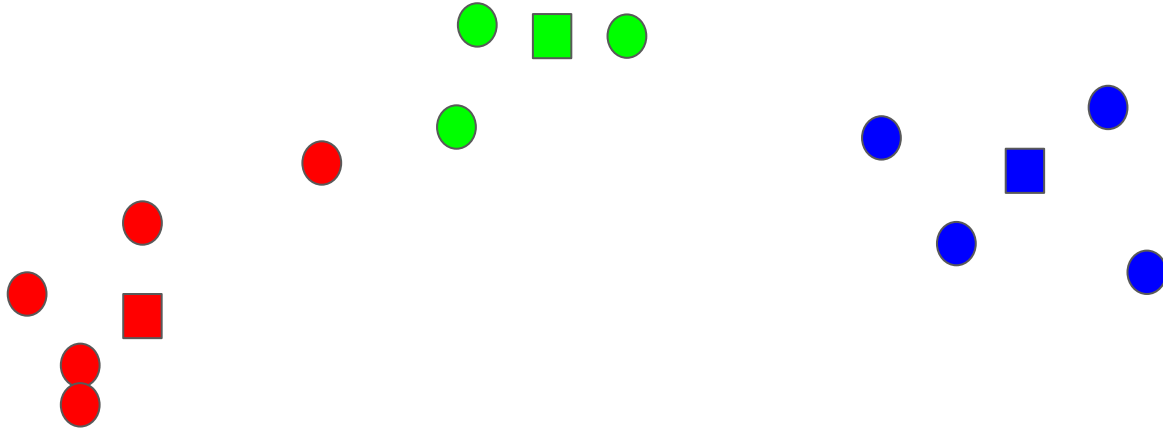
Assign data points to the closest centroid.



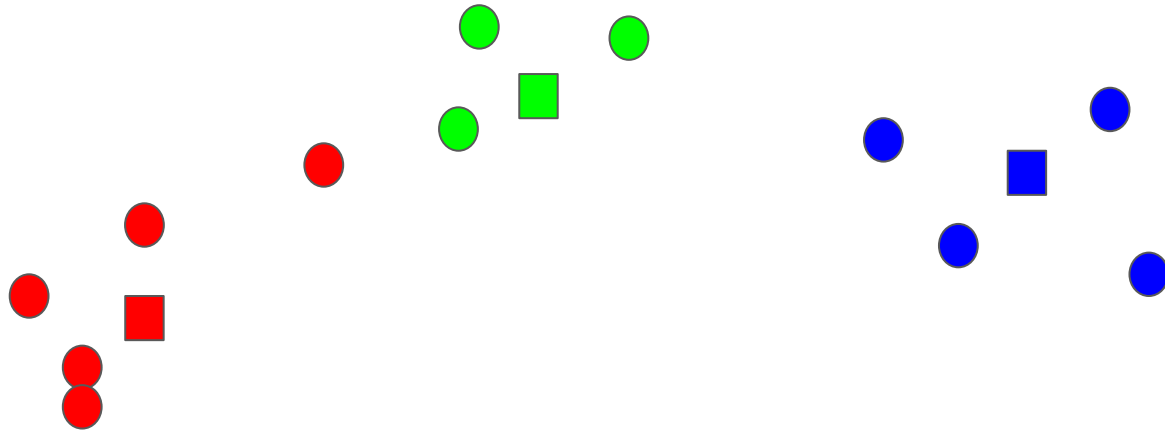
Recompute centroids.



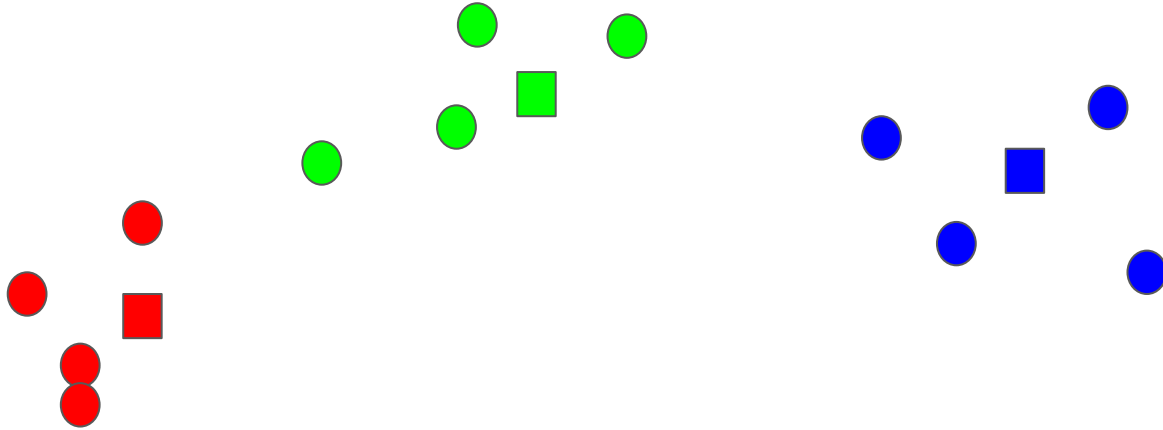
Re-assign data points to the closest centroid.



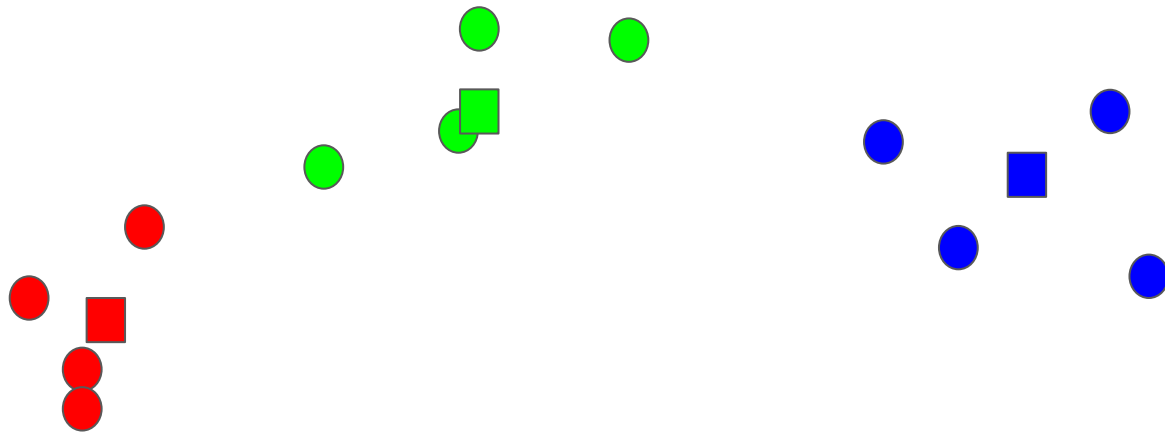
Recompute centroids.



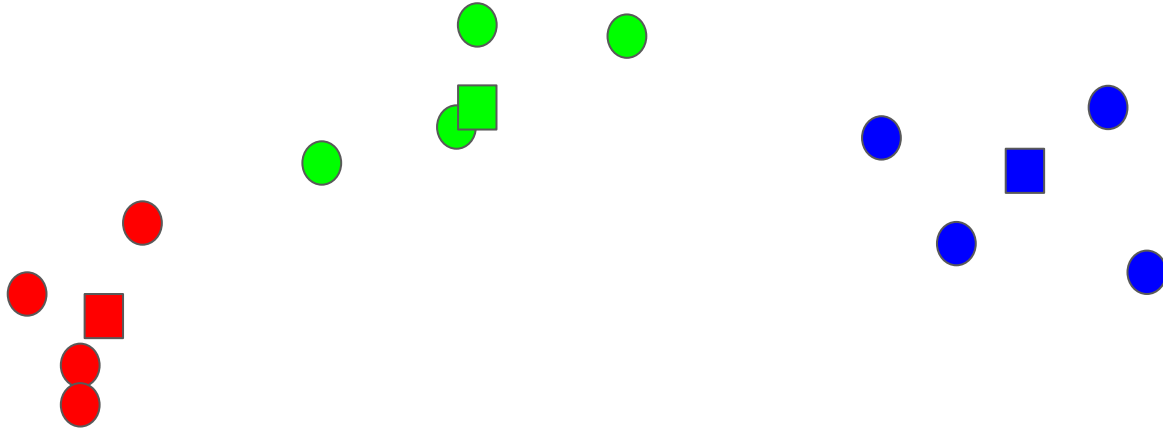
Re-assign data points to the closest centroid.



Recompute centroids.



Re-assign data points to the closest centroid.



Convergence criterion is met.

Quiz 1

Choose the correct statement(s) w.r.t K-Means clustering.

- It is often used for unlabelled data.
- It can be used to segment customers based on their past behaviour/characteristic.
- It puts two dissimilar points in same cluster.
- All of the above

Objective function

The **objective** of **K-Means clustering** is to **minimize** the **total intra-cluster distance** (**squared error**).

The diagram shows the objective function $J = \sum_{j=1}^k \sum_{i=1}^n \|x_i^{(j)} - c_j\|^2$ with several annotations: an arrow points from 'objective function' to J ; an arrow points from 'number of clusters' to k ; an arrow points from 'number of cases' to n ; an arrow points from 'case i ' to $x_i^{(j)}$; an arrow points from 'centroid for cluster j ' to c_j ; and a bracket under the distance term is labeled 'Distance function'.

$$\text{objective function} \leftarrow J = \sum_{j=1}^k \sum_{i=1}^n \underbrace{\|x_i^{(j)} - c_j\|^2}_{\text{Distance function}}$$

Where,

$$c_j = \frac{\sum_{x_i \in S_j} x_i}{|S_j|}$$

$|S_j|$ = **Number** of **instances** in **cluster j** .

Quiz 2

K in K-Means clustering stands for -

- Number of nearest neighbors
- Number of samples in each cluster
- Minimum distance between the clusters
- Number of clusters

Quiz 3

Which of the following can act as a termination criterion in K-Means?

- Fixed number of iterations
- Stationary centroids appear between successive iterations.
- The distance between the clusters is minimum.
- None of the above

Finding optimal value of k

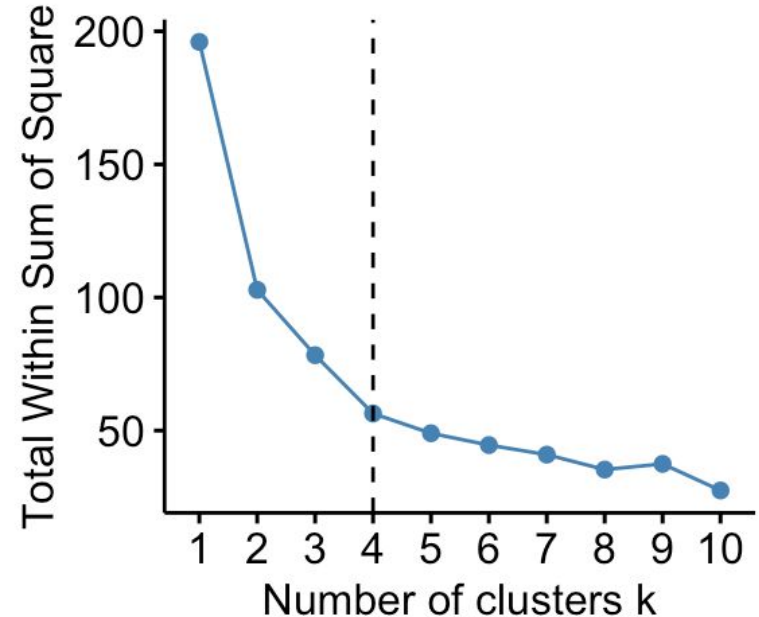
- Data points clustering is a subjective decision as there is no ground truth available. Domain knowledge or better business understanding may help in getting intuition behind right number of clusters.
- Additionally, there are few methods that help in selecting optimal value of k. Most commonly used are -
 1. Elbow method
 2. Average silhouette method

Elbow Method

- The **Elbow Method** is one of the most popular methods to determine this optimal value of k .
- It looks at the inertia for different values of k .

Inertia (or **within cluster sum of squared distance** or **intra cluster distance**) - It is the sum of squared distances of samples to their closest cluster centroid.

- Step 1. Perform k -means clustering for different values of k .
- Step 2. For each k , calculate the inertia.
- Step 3. Plot the curve of inertia according to the number of clusters k .
- Step 4. Choose the k where inertia stops decreasing abruptly.



NOTE - As k increases, the inertia tends towards zero.

Average Silhouette Method

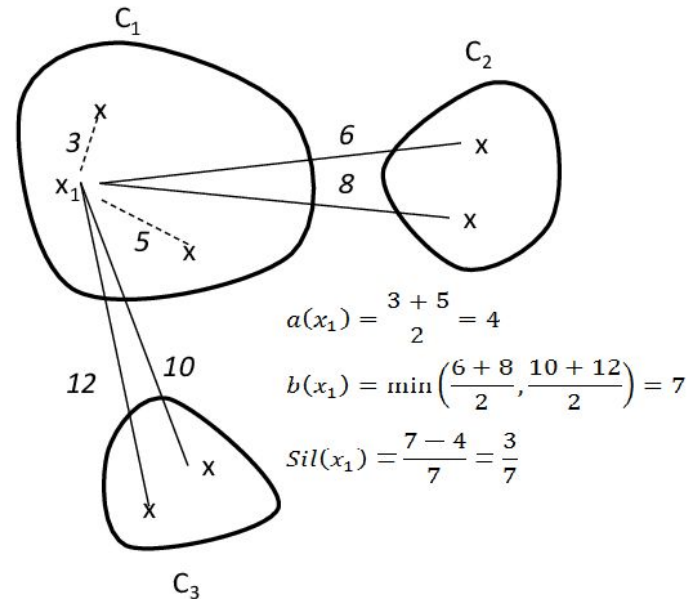
- Silhouette method is used to determine the degree of separation between clusters.

For each sample, it computes -

- mean intra-cluster distance (a) : average distance from all data points in the same cluster.
- mean nearest-cluster distance (b) : average distance from all data points in the closest cluster.
- Compute the Silhouette coefficient :

$$\frac{(b^i - a^i)}{\max(b^i, a^i)}$$

Like this, mean Silhouette Coefficient over all samples is calculated.

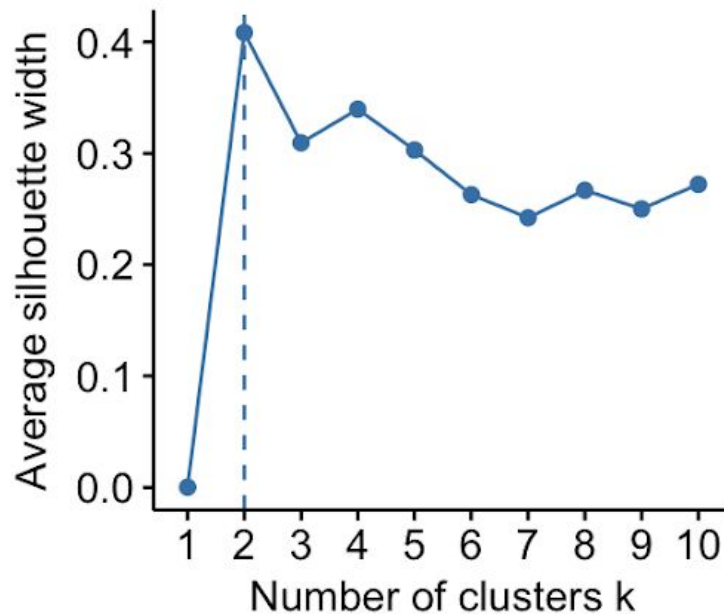


Average Silhouette Method [Contd.]

- The average Silhouette coefficient ranges from -1 to 1.
- If value = 0 -> the sample is very close to the neighboring clusters.
- If value = 1 -> the sample is far away from the neighboring clusters.
- If value = -1 -> the sample is assigned to the wrong clusters.

Finding optimal value of k using average Silhouette method

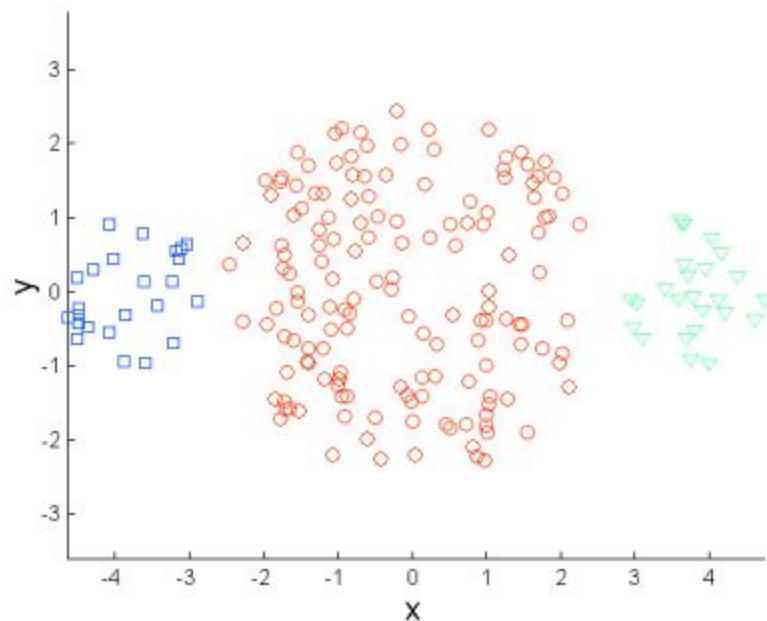
- Step 1. Perform k-means clustering for different values of k.
- Step 2. For each k, calculate the average Silhouette coefficient.
- Step 3. Plot the curve of average Silhouette coefficient according to the number of clusters k.
- Step 4. The location of the maximum is considered as the appropriate number of clusters.



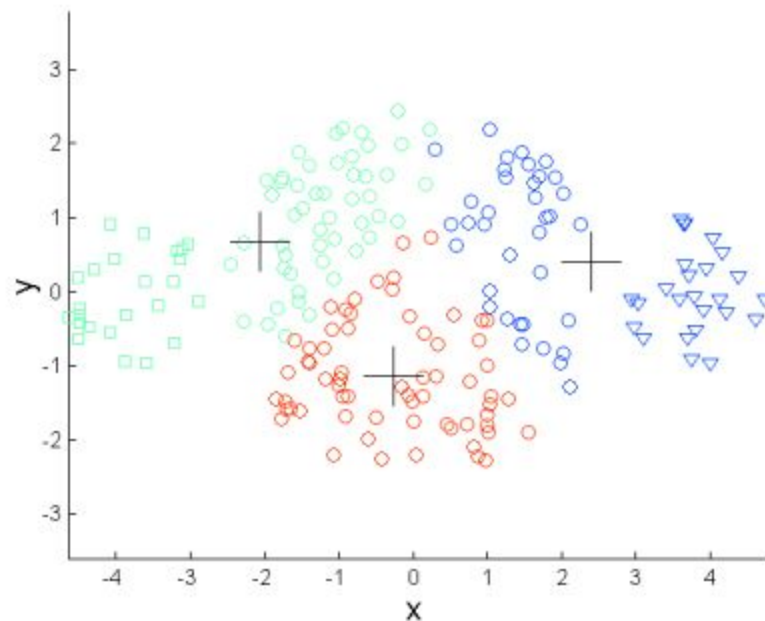
K-Means Limitations

- K-Means has problems when clusters are of different
 - Sizes
 - Densities
 - Non-globular shapes
- K-Means has problems when the data contains outliers.

K-Means Limitations - Differing Sizes

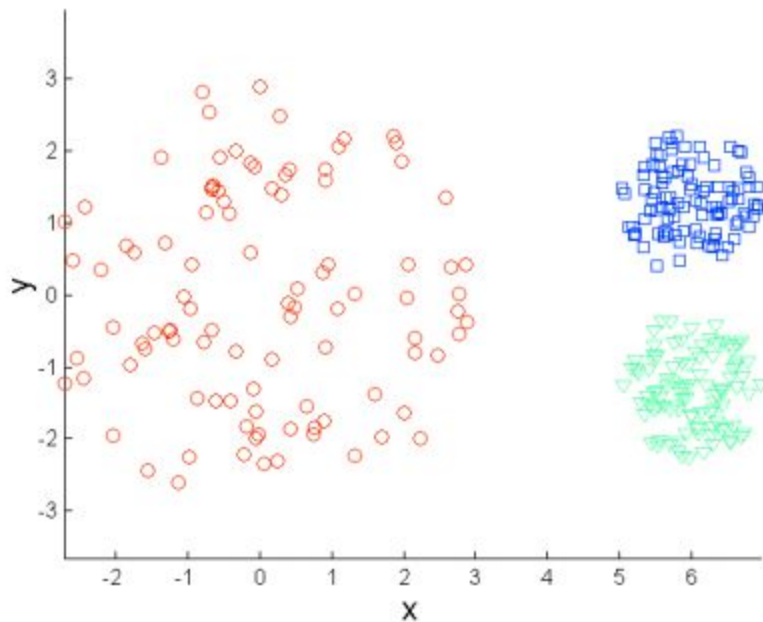


Original Points

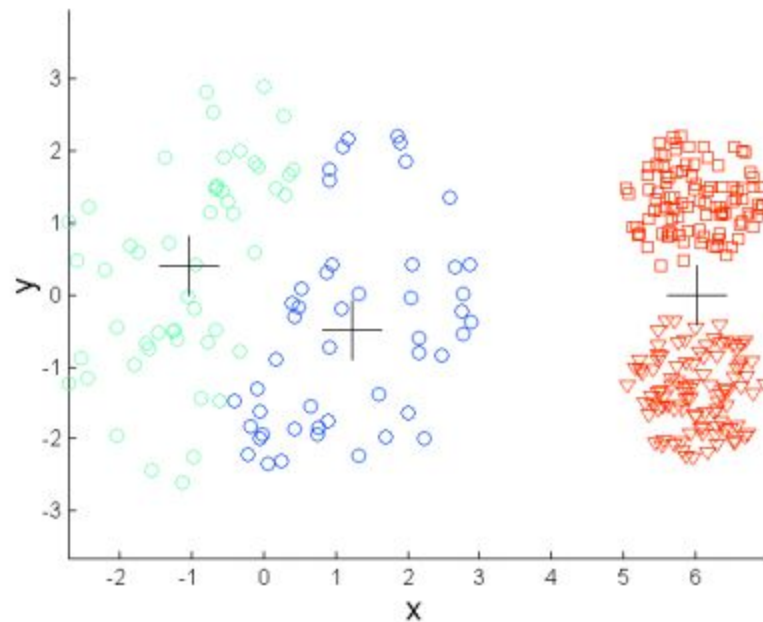


K-Means Clusters

K-Means Limitations - Differing Densities

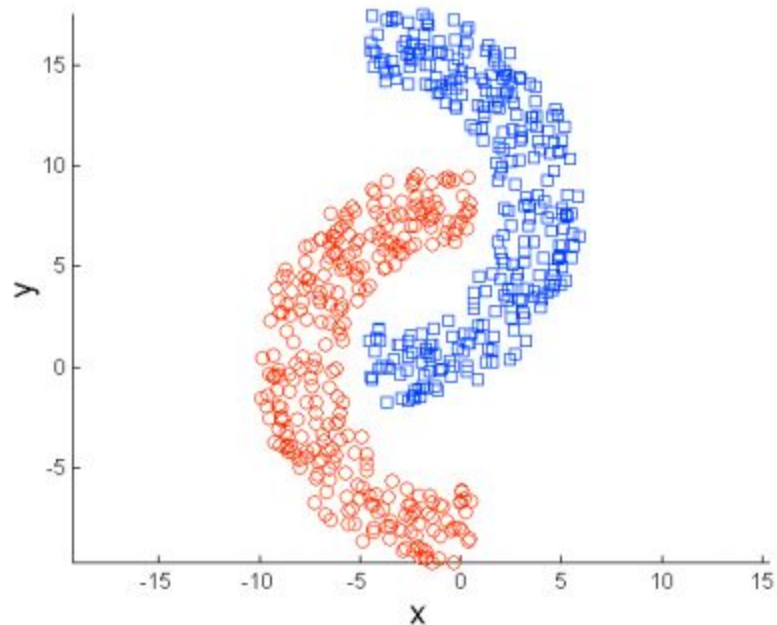


Original Points

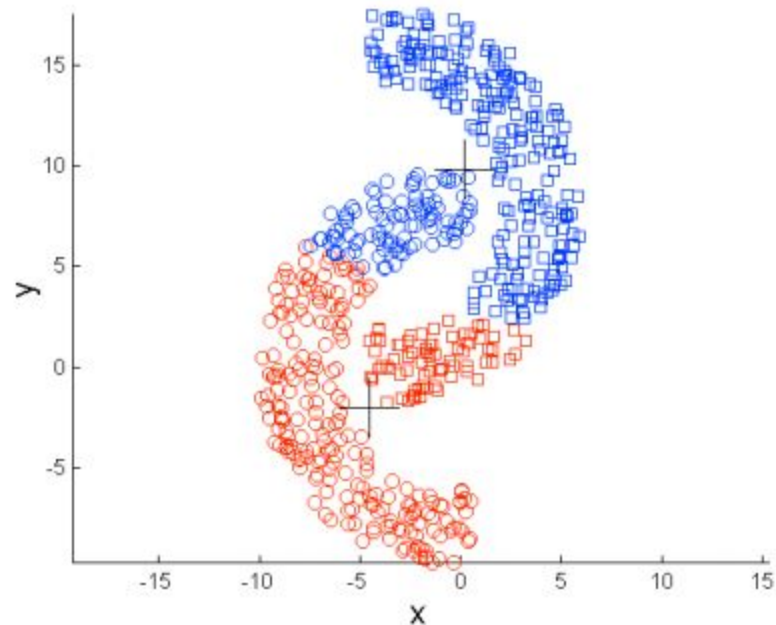


K-Means Clusters

K-Means Limitations - Non-globular shapes

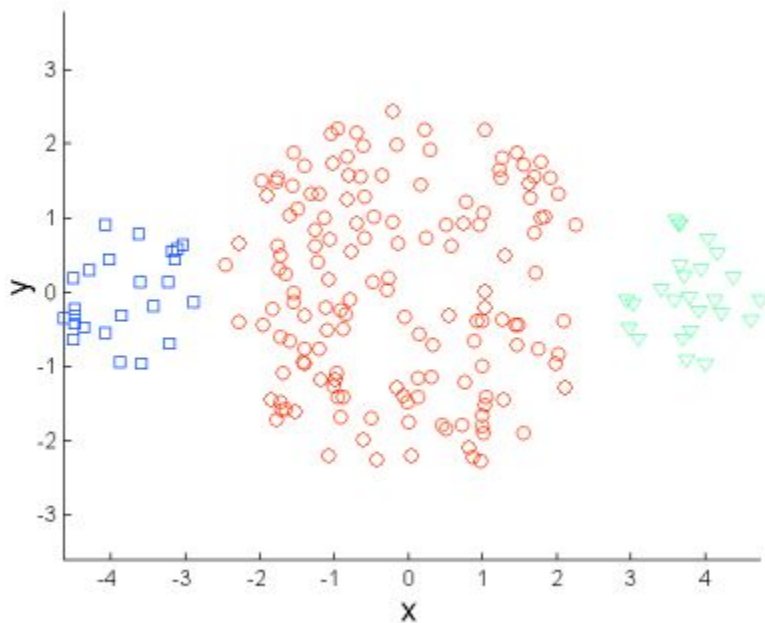


Original Points

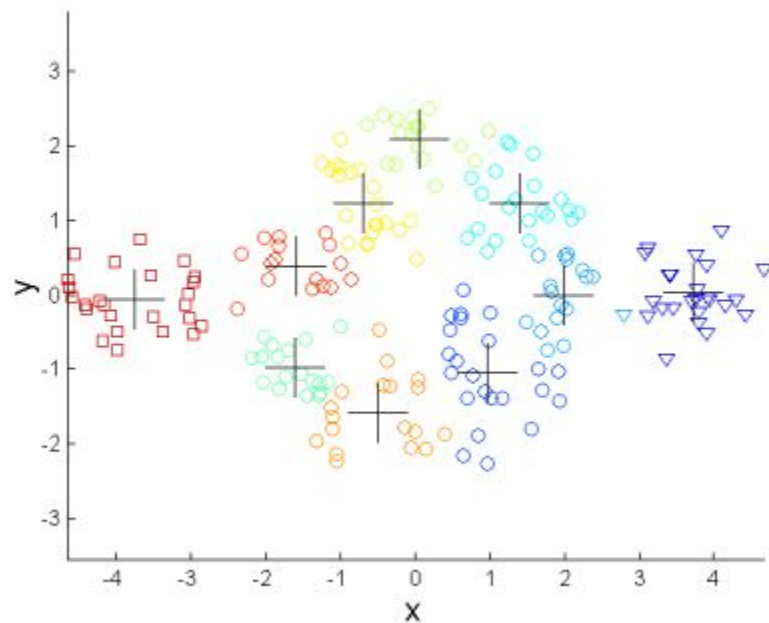


K-Means Clusters

Overcoming K-Means Limitation - Differing Sizes



Original Points

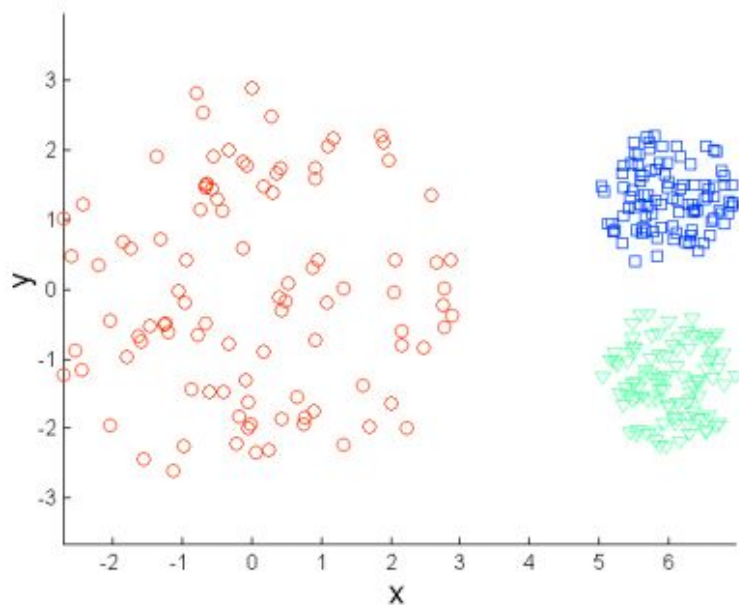


K-Means Clusters

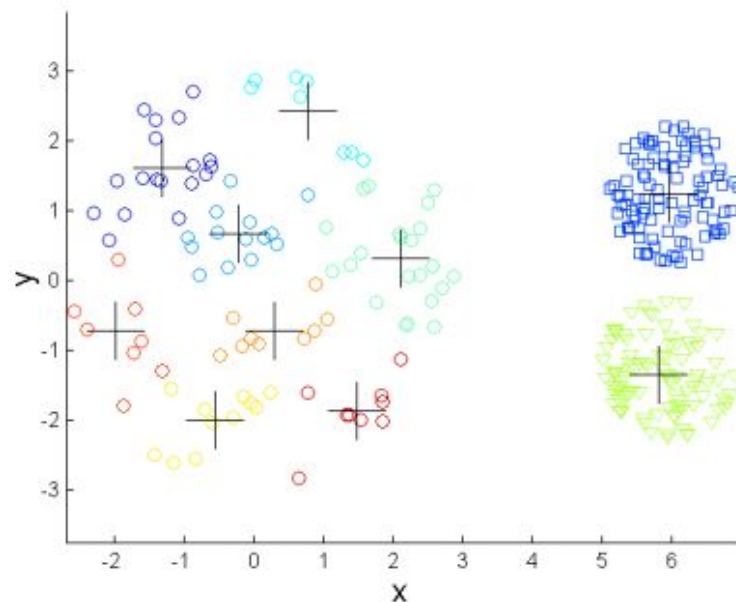
- One solution is trying larger value of k (say 10) and then group smaller clusters into larger.

NOTE - It's not a perfect solution but works in some cases.

Overcoming K-Means Limitation - Differing Densities



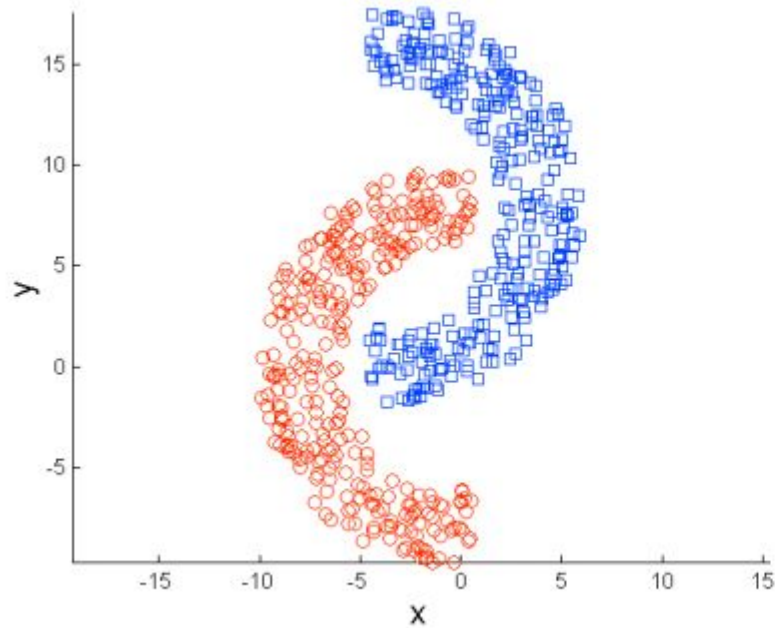
Original Points



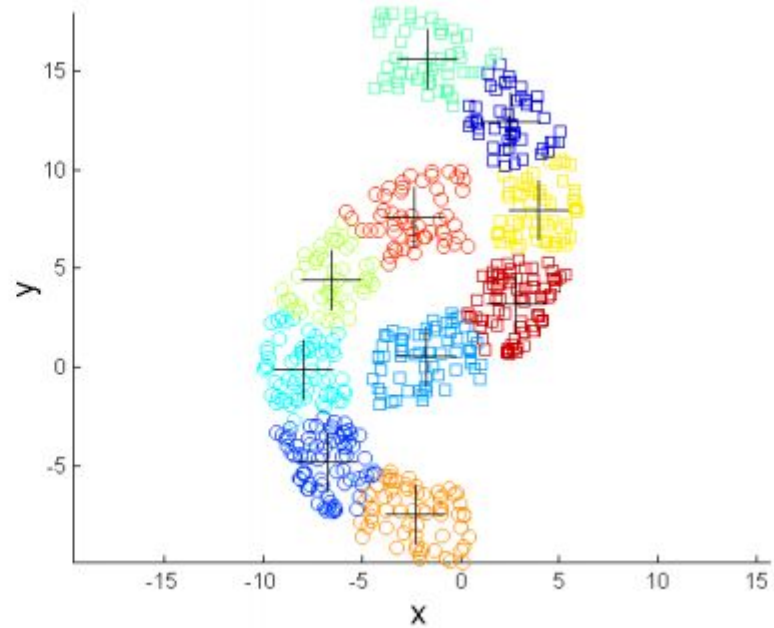
K-Means Clusters

- One solution is trying larger value of k (say 10) and then group smaller clusters into larger.

Overcoming K-Means Limitation - Non-globular shapes



Original Points



K-Means Clusters

- One solution is trying larger value of k (say 10) and then group smaller clusters into larger.

K-Means Variations

K-Medoid

Similar to K-Means but the centroid of the cluster is defined to be one of the actual data points in the cluster (the medoid).