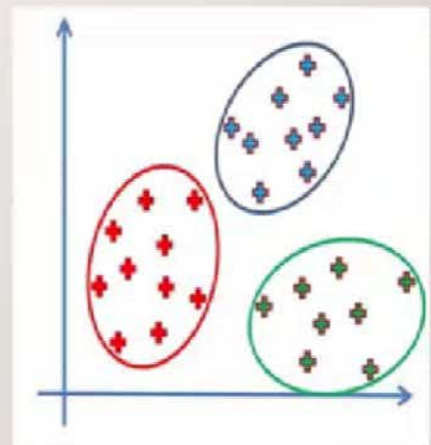


- 
- Clustering is similar to classification, but the basis is different.
  - In Clustering we don't know what we are looking for, and we are trying to identify some segments or clusters in our data.
  - When we use clustering algorithms on your dataset, unexpected things can suddenly pop up like structures, clusters and groupings we would have never thought of otherwise.



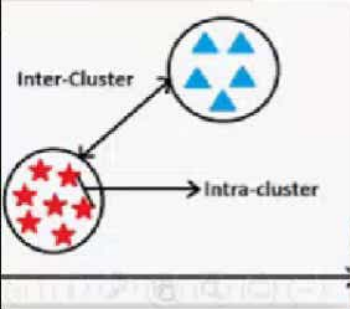
## METRICS FOR CLUSTERING

---

$$\text{Dunn Index} = \frac{\min(\text{Inter cluster distance})}{\max(\text{Intra cluster distance})}$$

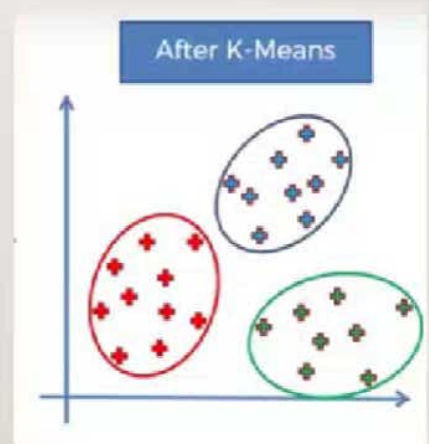
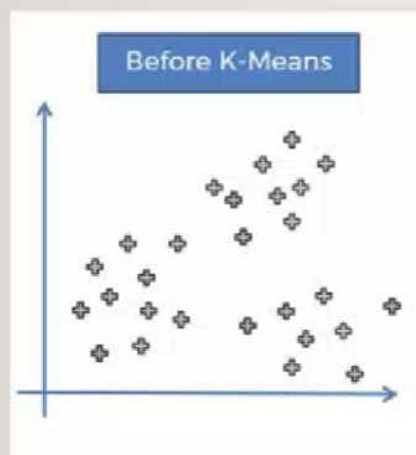
Clusters are far apart

Clusters are compact



CAN WE DIVIDE THE DATA INTO GROUPS??

---



## TYPES OF CLUSTERING

---

1. K –Means clustering
2. Hierarchical clustering
3. DBSCAN (Density based clustering)

# OBJECTIVE OF CLUSTERING

---

- Task ----→ Group similar points in one cluster

1. Points in one cluster are close together
2. Points in different clusters are far away

2.  $\frac{1}{n} \sum_{i=1}^n \|x_i - \mu\|^2$  (within cluster variance)

2.  $\frac{1}{n} \sum_{i=1}^n \|x_i - \mu\|^2$  (within cluster variance)

2.  $\frac{1}{n} \sum_{i=1}^n \|x_i - \mu\|^2$  (within cluster variance)

AutoSave On

Clustering

File Home Insert Draw Design Transitions Animations Slide Show Help

Press Esc to exit full screen

From Beginning From Current Slide Present Online Custom Slide Show Start Slide Show

Rehearse with Coach Rehearse

Set Up Slide Show Hide Slide

Rehearse Record Slide Timings Show Play Narrations Show

Use Timings

Automatic

Always Use Subtitles

Use Presenter View

Subtitle Settings

Monitors

Captions & Subtitles

1 CLUSTERING

2

3

4

5

6

CLUSTERING

DAMANDEEP KAUR

[LINKEDIN.COM/IN/DAMANVIRDI/](https://www.linkedin.com/in/damanvirdi/)

0:07 / 42:09

# K-MEANS IS CENTROID BASED CLUSTERING

---

- K- no of parameters (hyper-parameter)
- HyperParameters are those which can't be deduced from the given data.
- Eg K= 3 ie. 3 Means Clustering
- No of clusters = 3
- For each cluster, No of centroids = 3 ( $C_1$  ,  $C_2$  ,  $C_3$ )
- For each cluster , No of corresponding set of points = 3( $S_1$  ,  $S_2$  ,  $S_3$ )
- $S_1 \cup S_2 \cup S_3 = D$  i.e Each point should belong to any cluster
- $S_1 \cap S_2 = \emptyset$  ,  $S_1 \cap S_3 = \emptyset$  ,  $S_2 \cap S_3 = \emptyset$  i.e No point belongs to more than one cluster

AutoSave On

File Home Insert Draw Design Transitions Animations Slide Show Review View Help

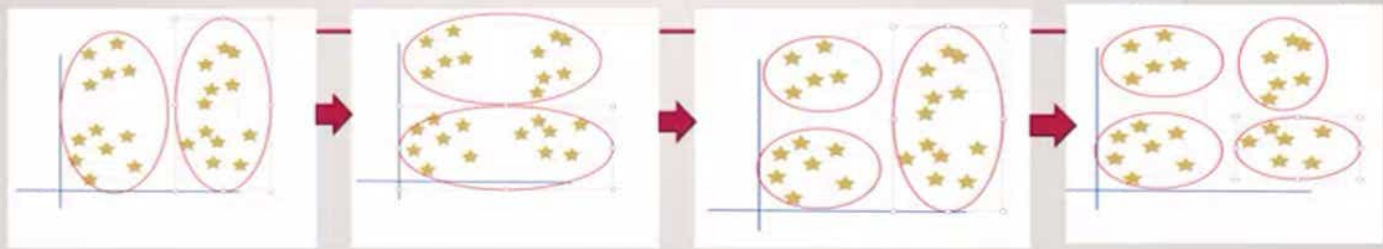
From Beginning From Current Slide Present Online Custom Slide Show Start Slide Show Rehearse with Coach Rehearse Set Up Slide Show Hide Slide Rehearse Record Slide Show Set Up Keep Slides Updated Use Timings Play Narrations Show Media Controls Use Presenter View Always Use Subtitles Subtitle Settings

# APPLICATIONS OF CLUSTERING

- [https://en.wikipedia.org/wiki/Cluster\\_analysis#Applications](https://en.wikipedia.org/wiki/Cluster_analysis#Applications)
- Ecommerce -> To 'gp' similar customers
- Where manual labelling is time consuming, perform clustering as pre processing step and then do manual labelling

Slide 4 of 29 English (India) 17:10 14-05-2021





1. Points in one cluster are close together
2. Points in different clusters are far away

AutoSave Clustering Daman Verdi

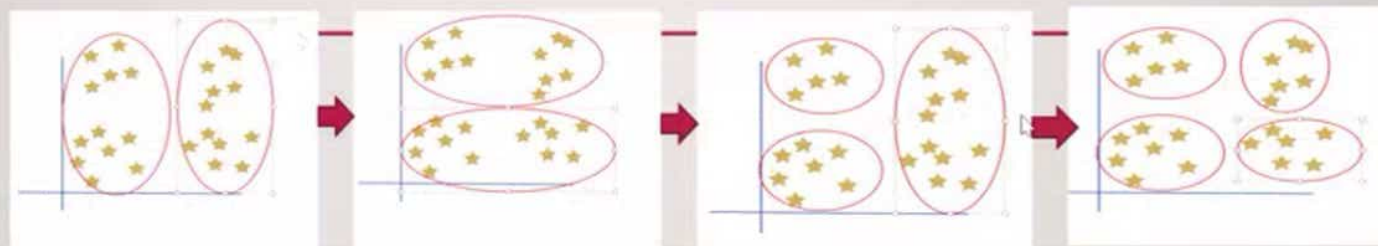
File Home Insert Draw Design Transitions Animations **Slide Show** Review View Help

From Beginning From Current Slide Present Online Custom Slide Show Start Slide Show Rehearse with Coach Rehearse Set Up Slide Show Hide Slide Rehearse Record Slide Show Keep Slides Updated Use Timings Play Narrations Show Media Controls Use Presenter View Always Use Subtitles Subtitle Settings

## K-MEANS ALGORITHM

1. Choose the number  $K$  of clusters.
2. Select random  $K$  points, the **CENTROIDS** (not necessarily from the dataset)
3. Assign each data point to the closest centroid  $\rightarrow$  It form  $K$  clusters
4. Compute and place new centroid of each cluster
5. Reassign each datapoint to the new closest centroid. If any reassignment happens . go to step4 otherwise finish.

Slide 10 of 25 English (India) 17:17 14-05-2021

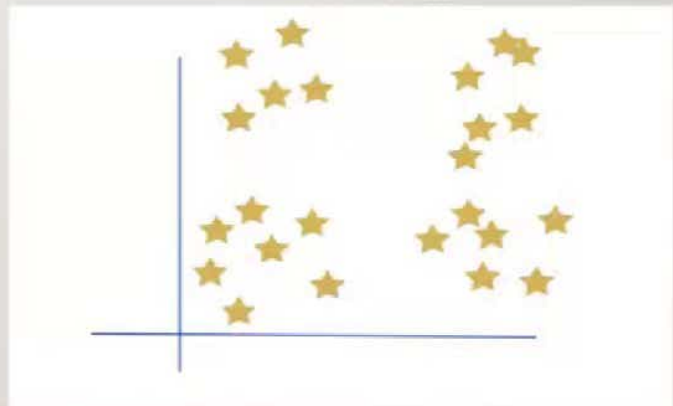


## DETERMINING THE RIGHT K

---

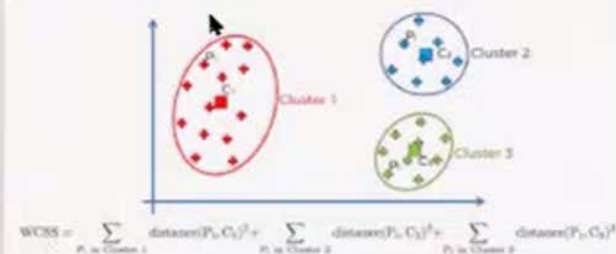
- Once we get the centroids, we can easily describe the datapoint. But first how many clusters should be there?

Q- Find no of clusters?



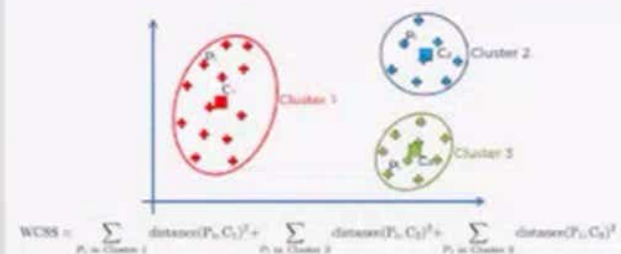
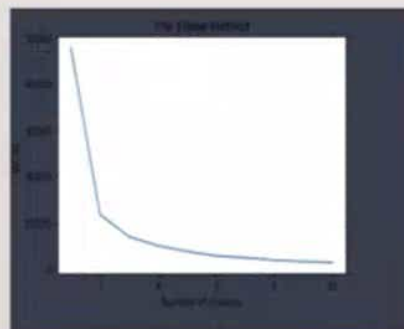
# DETERMINING THE RIGHT K

- K is Hyper Parameter i.e. It's value can't be estimated from the data
- 1. DOMAIN KNOWLEDGE e.g. –Movie Reviews.. Only three categories are possible +ve, neutral, -ve
- 2. Elbow Method :



# DETERMINING THE RIGHT K

- K is Hyper Parameter i.e. It's value can't be estimated from the data
- 1. DOMAIN KNOWLEDGE e.g. –Movie Reviews.. Only three categories are possible +ve, neutral, -ve
- 2. Elbow Method :



# K-MEANS ALGORITHM

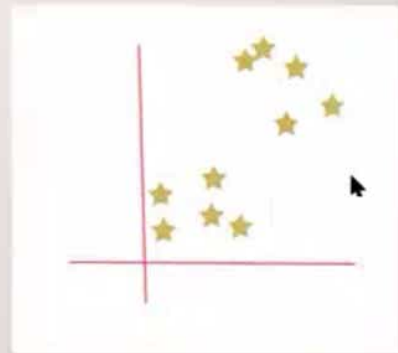
---

1. Choose the number  $K$  of clusters.
2. Select random  $K$  points, the CENTROIDS (not necessarily from the dataset)
3. Assign each data point to the closest centroid  $\rightarrow$  It form  $K$  clusters
4. Compute and place new centroid of each cluster
5. Reassign each datapoint to the new closest centroid. If any reassignment happens , go to step4 otherwise finish.

## STEP 1: CHOOSE THE NUMBER K OF CLUSTERS

---

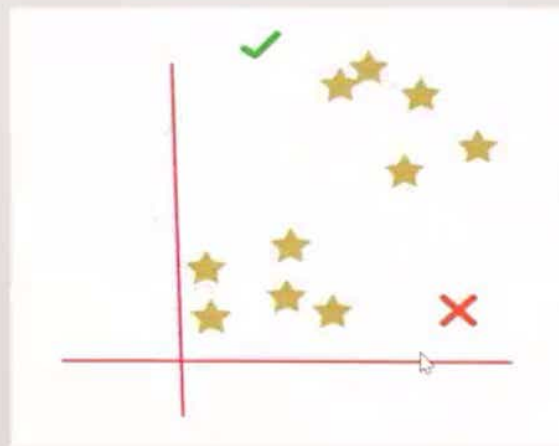
$K = 2$





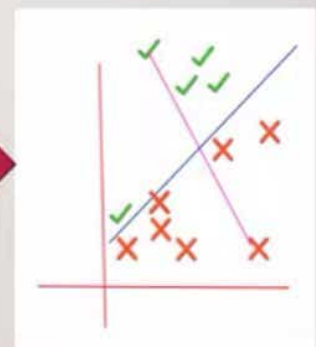
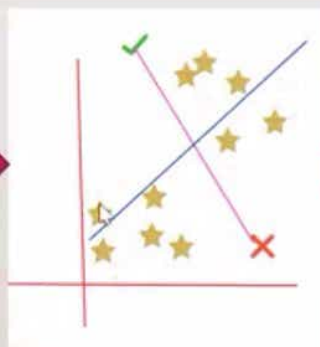
## STEP 2: SELECT RANDOM K POINTS, THE CENTROIDS

---



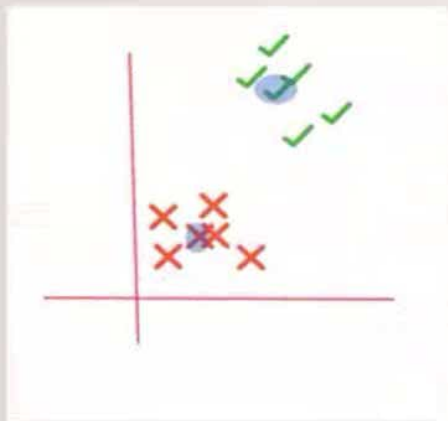
### STEP 3 - ASSIGN EACH DATA POINT TO THE CLOSEST CENTROID

---

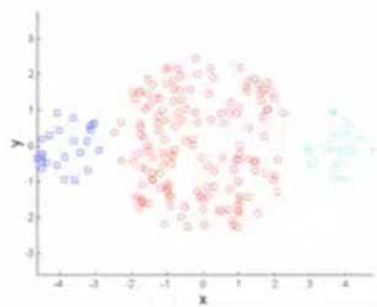


REPEAT STEP 5-REASSIGN EACH DATAPOINT TO THE NEW CLOSEST CENTROID. SINCE NO REASSIGNMENT HAPPENS ,FINISH.

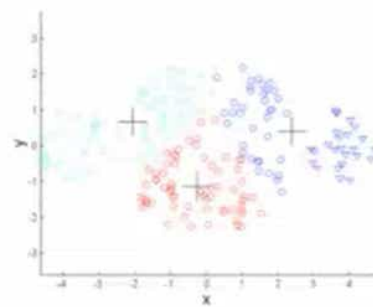
---



## Limitations of K-means: Differing Sizes



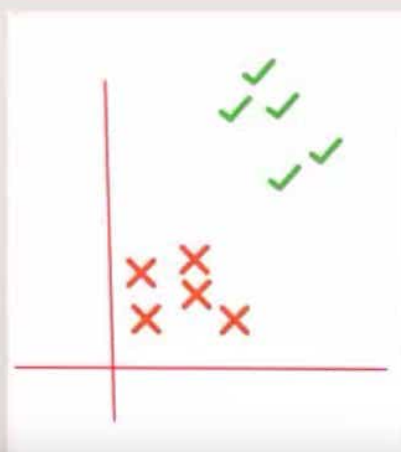
Original Points



K-means (3 Clusters)

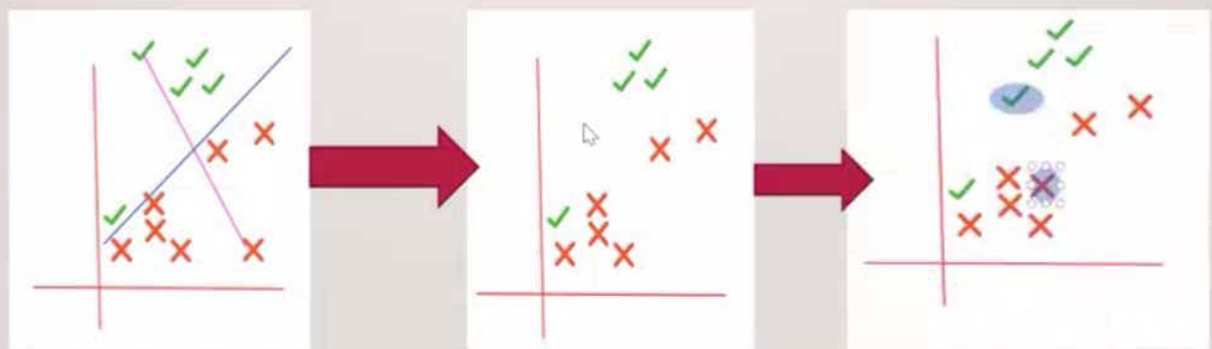
## FINAL CLUSTERS

---



## STEP -4 COMPUTE AND PLACE NEW CENTROID OF EACH CLUSTER

---



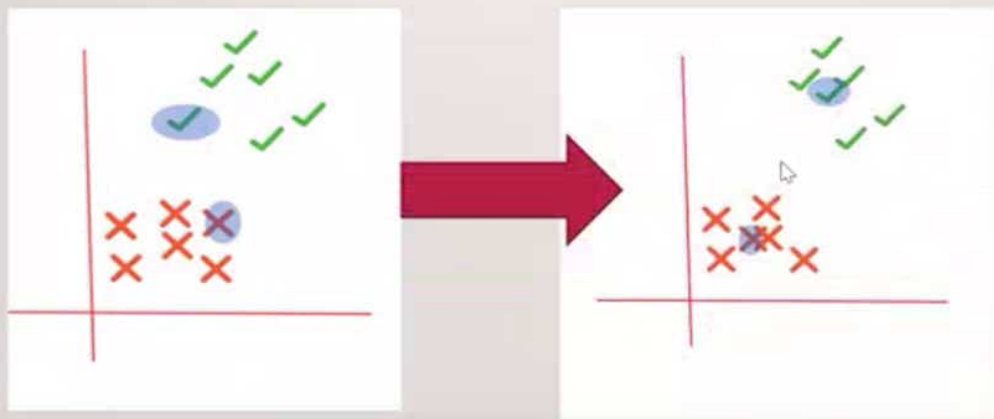
## LIMITATIONS OF K-MEANS

---

- K-Means have problem when clusters are of different
  1. sizes.
  2. Density
  3. Non-globular shapes
- K-Means have problem when data contains outliers

REPEAT STEP -4 COMPUTE AND PLACE NEW  
CENTROID OF EACH CLUSTER

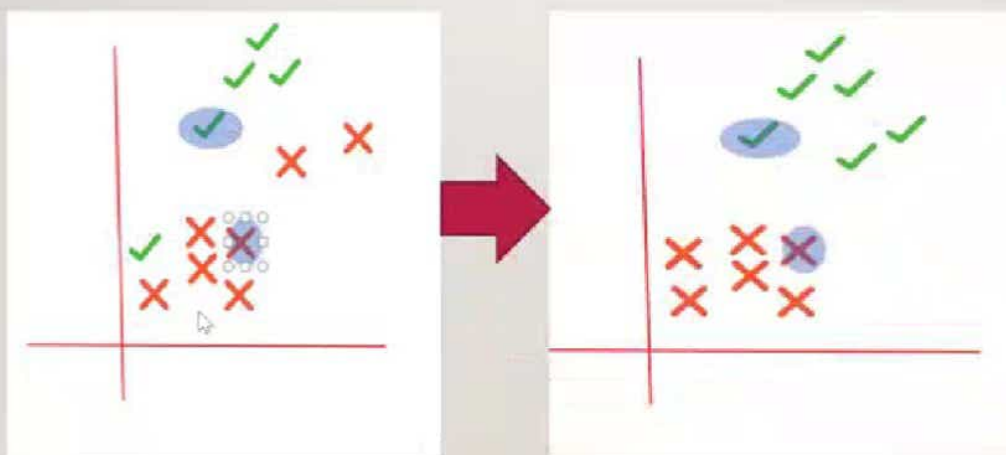
---





## STEP -5 REASSIGN EACH DATAPOINT TO THE NEW CLOSEST CENTROID

---



Home Page - Select or create a... x | W Edge detection - Wikipedia x | Home Page - Select or create a... x | k\_means\_clustering - Jupyter No... x

localhost:8890/notebooks/k\_means\_clustering.ipynb

File Edit View Insert Cell Kernel Widgets Help Trusted Python 3

### Using the elbow method to find the optimal number of clusters

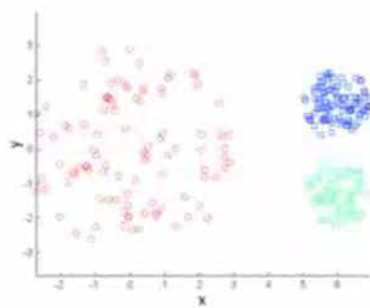
```
1 from sklearn.cluster import KMeans
2 wcss = []
3 for i in range(1, 11):
4     kmeans = KMeans(n_clusters = i, init = 'k-means++', random_state = 42)
5     kmeans.fit(X)
6     wcss.append(kmeans.inertia_)
7
8
9 print(wcss)
10 plt.plot(range(1, 11), wcss)
11 plt.title('The Elbow Method')
12 plt.xlabel('Number of clusters')
13 plt.ylabel('WCSS')
14 plt.show()
15 print(wcss)
```

### Training the K-Means model on the dataset

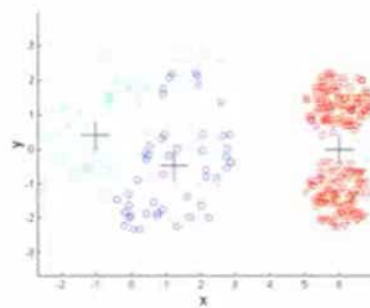
Type here to search

17:35 14-05-2021 ENG

## Limitations of K-means: Differing Density

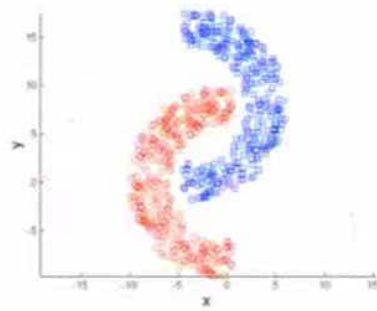


Original Points

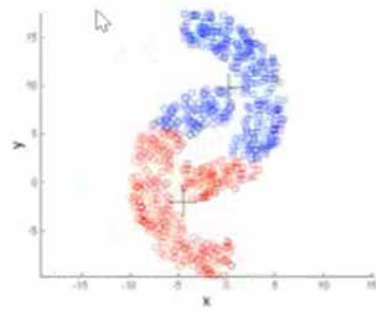


K-means (3 Clusters)

## Limitations of K-means: Non-globular Shapes



Original Points



K-means (2 Clusters)

Home Page - Select or create a... x Edge detection - Wikipedia x Home Page - Select or create a... x k\_means\_clustering - Jupyter No... x

localhost:8890/notebooks/k\_means\_clustering.ipynb

File Edit View Insert Cell Kernel Widgets Help Python 3

## K-Means Clustering

### Importing the libraries

```
1 import numpy as np
2 import matplotlib.pyplot as plt
3 import pandas as pd
4
```

### Importing the dataset

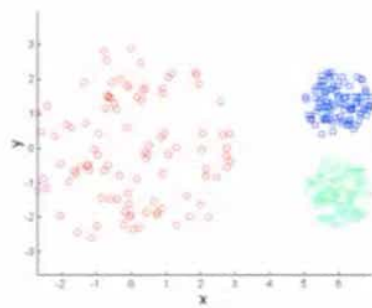
```
1 dataset = pd.read_csv('customer_gp.csv')
2 X = dataset.iloc[:, :].values
3 x
```

### Using the elbow method to find the optimal number of clusters

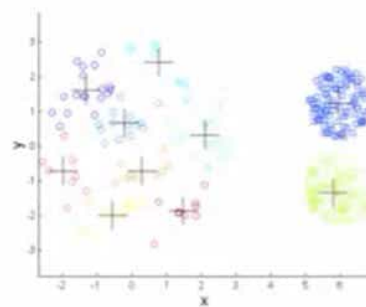
Type here to search

17:34 14-05-2021 ENG

## Overcoming K-means Limitations

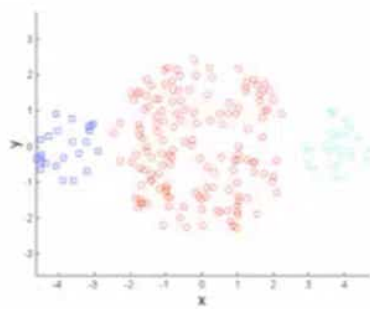


Original Points

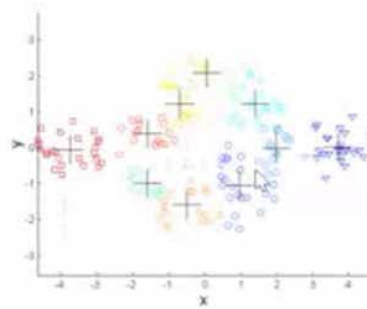


K-means Clusters

## Overcoming K-means Limitations



Original Points



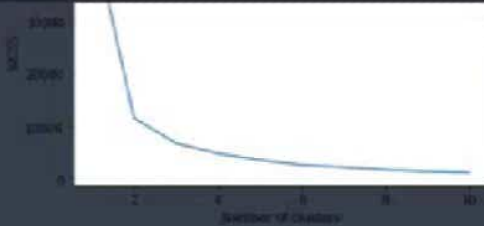
K-means Clusters

One solution is to use many clusters.  
Find parts of clusters, but need to put together.

Home Page - Select or create a ... x Edge detection - Wikipedia x Home Page - Select or create a ... x k\_means\_clustering - Jupyter No x +

localhost:8890/notebooks/k\_means\_clustering.ipynb

File Edit View Insert Cell Kernel Widgets Help Trusted Python 3



Training the K-Means model on the dataset

```
1 kmeans = KMeans(n_clusters = 2, init = 'k-means++', random_state = 42)
2 y_kmeans = kmeans.fit_predict(X)
3 print(y_kmeans)
```

Visualising the clusters

```
1 x
2 y_kmeans
```

Type here to search

17:39 14-05-2021 ENG



Home Page - Select or create a ... x

Edge detection - Wikipedia x

Home Page - Select or create a ... x

k\_means\_clustering - Jupyter No x

+

localhost:8890/notebooks/k\_means\_clustering.ipynb

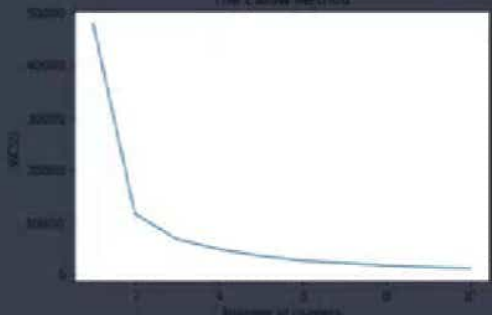
File Edit View Insert Cell Kernel Widgets Help

Python 3

```
(2) plt.ylabel('WCSS')
(3) plt.show()
(4) #print(wcss)
```

[47860.82, 31617.521794871793, 6884.461115942, 4860.863895238096, 3608.5829137528136, 2656.418658793651, 2236.379365979365, 1744.984362840866, 1406.0984526984129, 1229.862380523811]

The Elbow Method



Number of Clusters	WCSS
1	47860.82
2	31617.521794871793
3	6884.461115942
4	4860.863895238096
5	3608.5829137528136
6	2656.418658793651
7	2236.379365979365
8	1744.984362840866
9	1406.0984526984129
10	1229.862380523811

Training the K-Means model on the dataset

Type here to search

17:38

14-05-2021

ENG

