

K-means Clustering



Agenda of Today's Session

- What is Clustering?
- Types of Clustering
- What is K- Means Clustering?
- How does a K-Means Algorithm works?
- K-Means with Python





What is Clustering?

What is Clustering?

“Clustering is the process of dividing the dataset into groups, consisting of similar data-points”

- Points in the same group are as similar as possible
- Points in different group are as dissimilar as possible



What is Clustering?



Group of dining
in a restaurant



Items arranged in
a mall

Where is it Used?

The Amazon logo, featuring the word "amazon" in a bold, black, sans-serif font. A yellow curved arrow starts under the letter 'a' and points to the letter 'z'.

Recommendation System

The Netflix logo, consisting of the word "NETFLIX" in a bold, red, sans-serif font.

Recommended Movies

The Flickr logo, featuring two overlapping circles (one blue, one pink) above the word "flickr" in a blue, sans-serif font. The letter 'r' is pink.

Flickr's Photos

How business use Clustering?



Retail Store



Banking



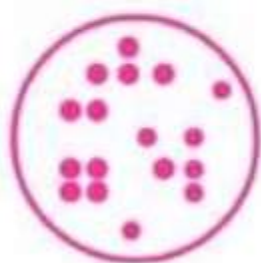
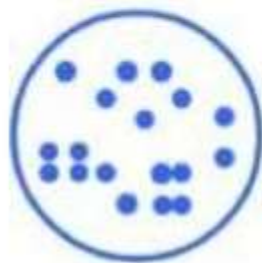
**Insurance
Companies**

Types of Clustering

- Exclusive Clustering
- Overlapping Clustering
- Hierarchical Clustering

Exclusive Clustering

- Hard Clustering
- Data Point / Item belongs exclusively to one cluster
- For Example: K-Means Clustering

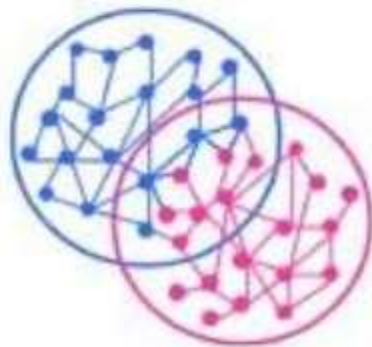


Types of Clustering

- ☐ Exclusive Clustering
- ☒ Overlapping Clustering
- ☐ Hierarchical Clustering

Overlapping Clustering

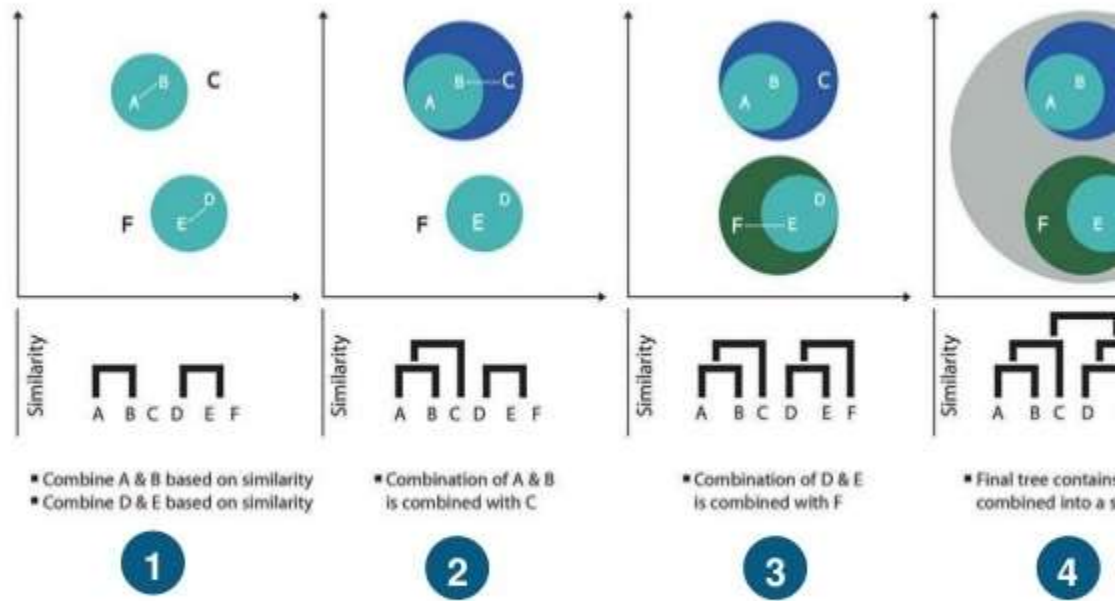
- Soft Cluster
- Data Point/ Item belongs to multiple cluster
- For Example: Fuzzy/ C-Means Clustering



Types of Clustering

- Exclusive Clustering
- Overlapping Clustering
- Hierarchical Clustering

Hierarchical Clustering



What is K-Means Clustering?

“K-Means is a clustering algorithm whose main goal is to group similar elements or data points into K clusters.”

NOTE: ‘K’ in **K-Means** represent the number of clusters

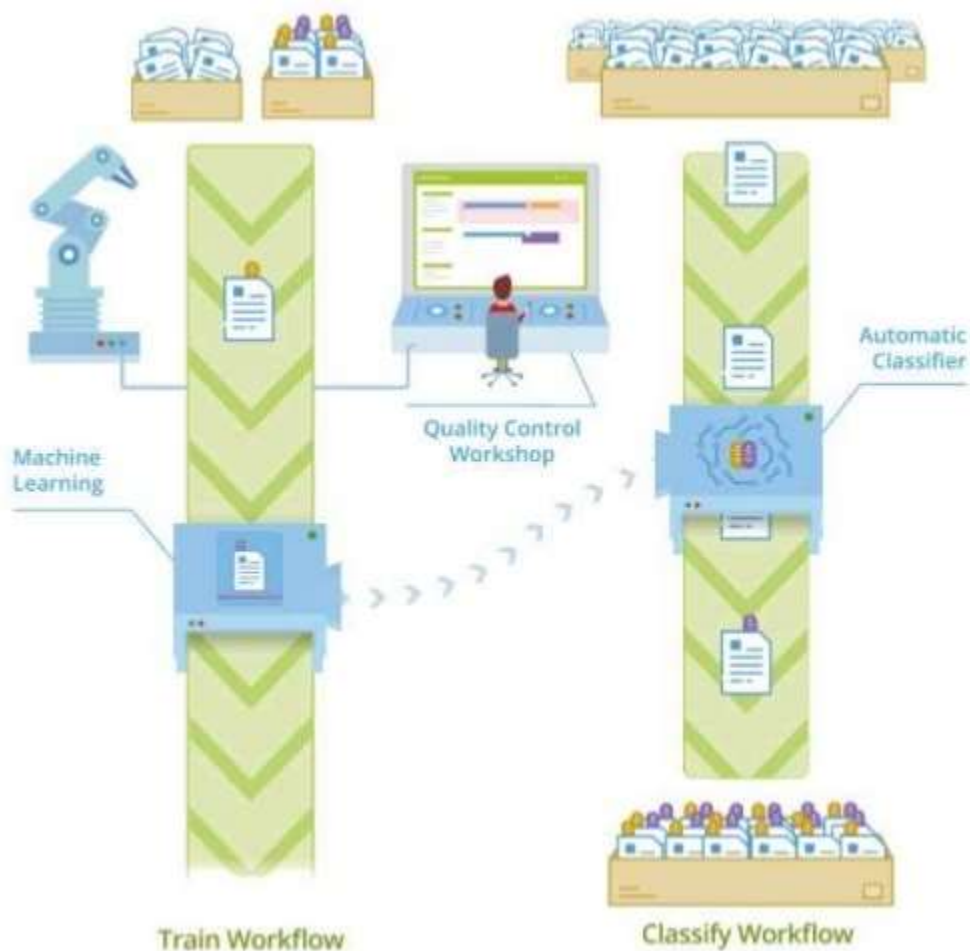


What is K-Means Clustering?

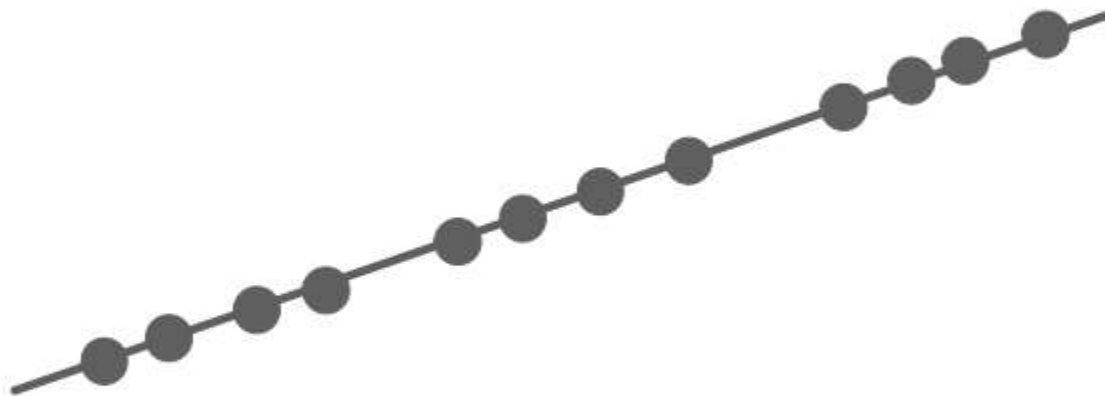


Pile of dirty clothes

Where Can I apply K-Means?



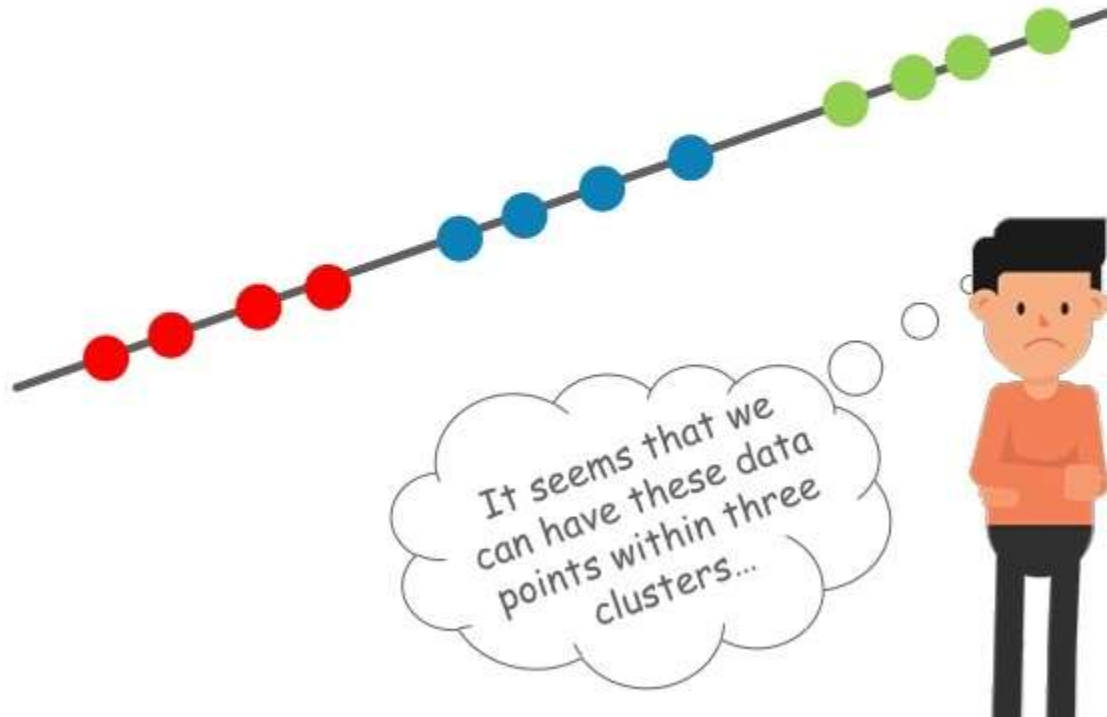
K-Means Algorithm



Number of Clusters = 3

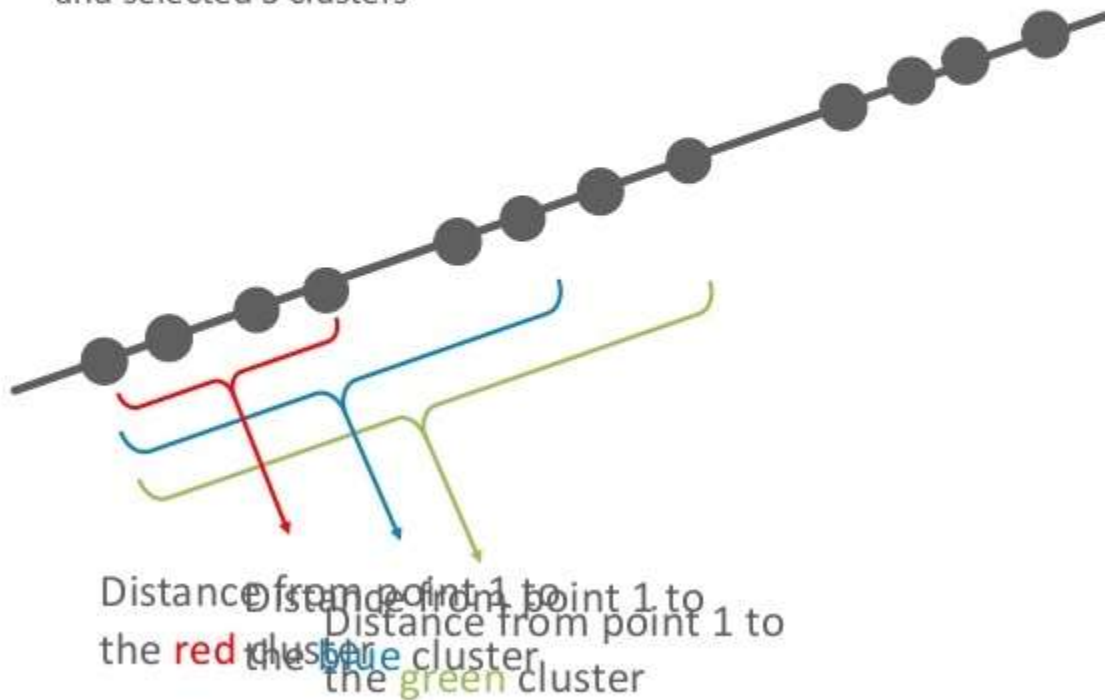
K-Means Algorithm

Number of Clusters, $K = 3$

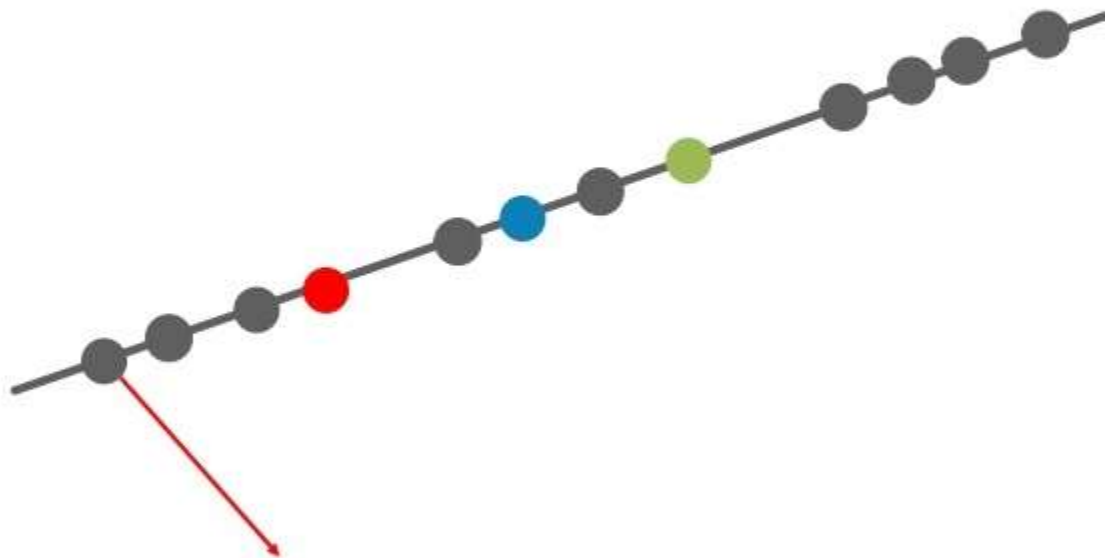


K-Means Algorithm

- **Step 1:** Select the number of clusters to be identified, i.e select a value for $K = 3$ in this case
- **Step 2:** Randomly select 3 distinct data point
- **Step 3:** Measure the distance between the 1st point and selected 3 clusters



K-Means Algorithm



Step 4: Assign the 1st
point to nearest cluster
(red in this case).

K-Means Algorithm

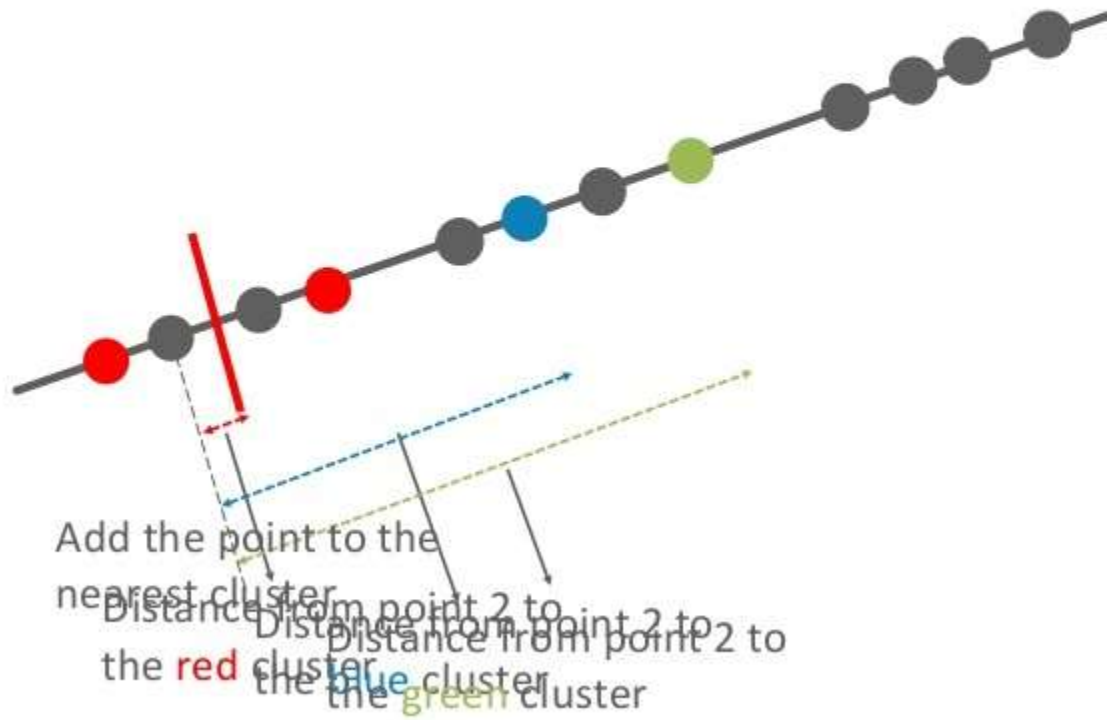


Step 5: Calculate the mean value including the new point for the **red** cluster

K-Means Algorithm

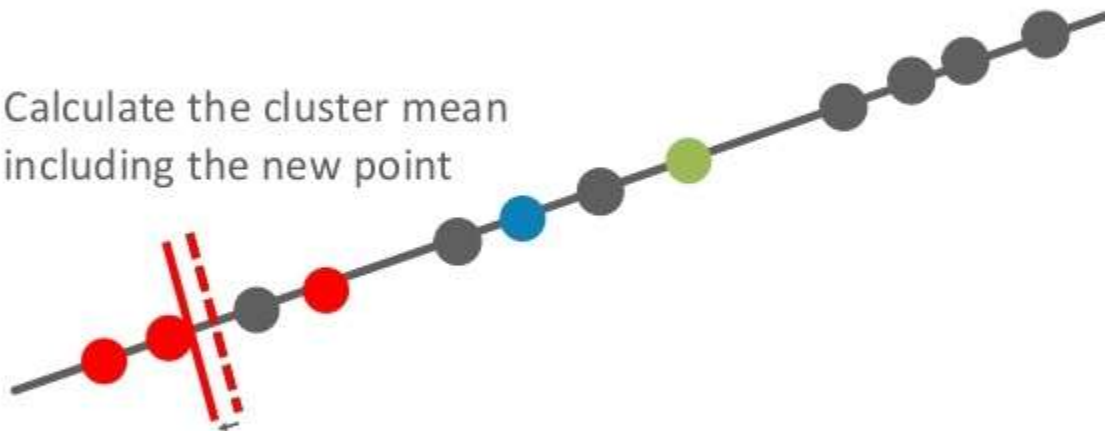
Find to which cluster does point 2 belongs to, how?

- Repeat the same procedure but measure the distance to the **red** mean



K-Means Algorithm

Calculate the cluster mean
including the new point



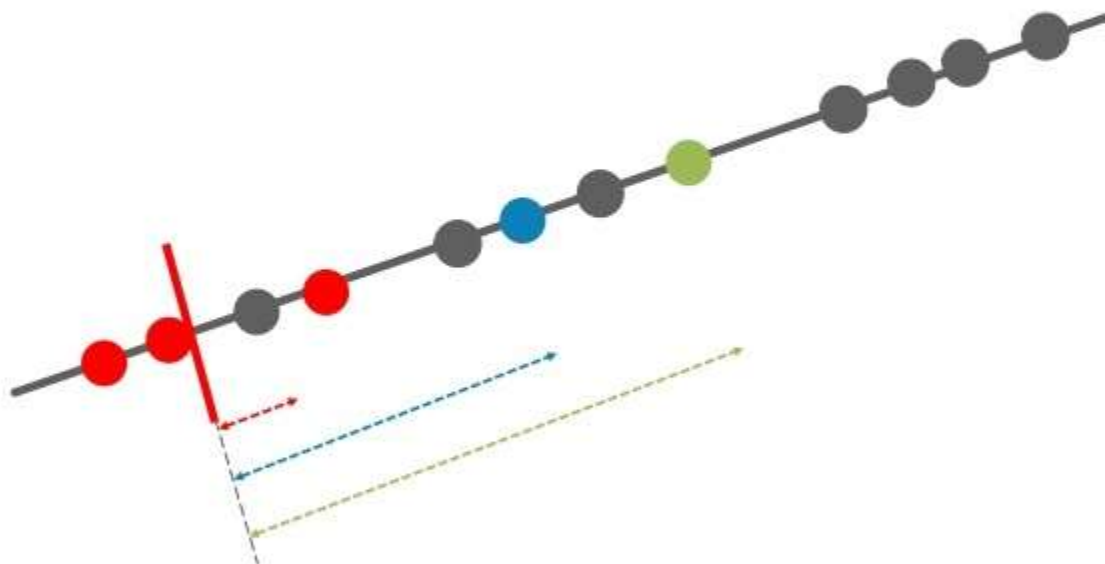
K-Means Algorithm

Find to which cluster does point 3 belongs to, how?

- Repeat the same procedure but measure the distance to the **red** mean



K-Means Algorithm



Measure the distance and add
the 3rd point to the nearest
cluster (red)

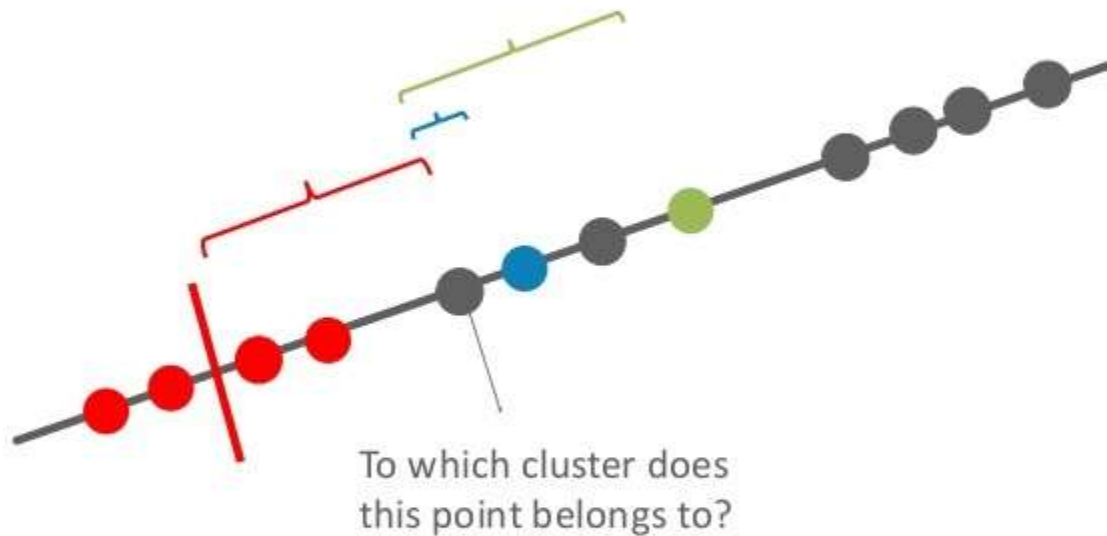
K-Means Algorithm

Calculate the new cluster mean using the new point



K-Means Algorithm

- Measure the distance
- Assign the point to the nearest cluster
- Calculate the cluster mean using the new point



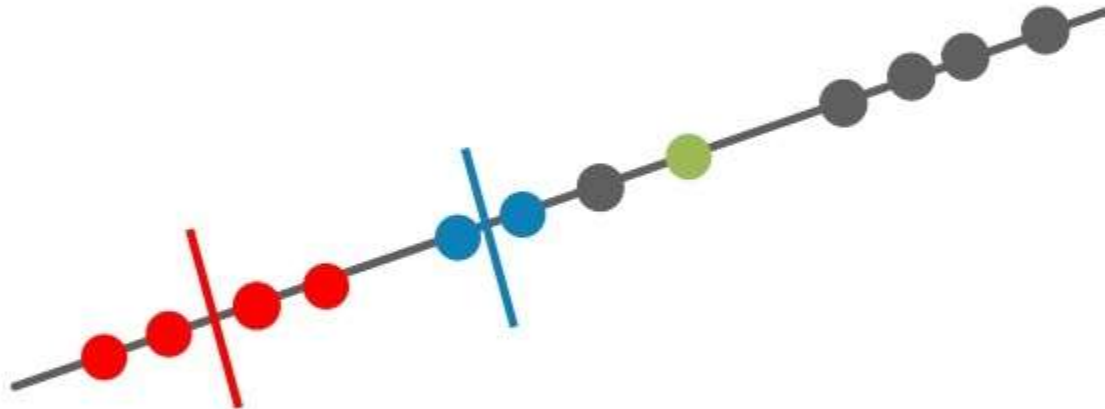
K-Means Algorithm

- Measure the distance
- Assign the point to the nearest cluster
- Calculate the cluster mean using the new point



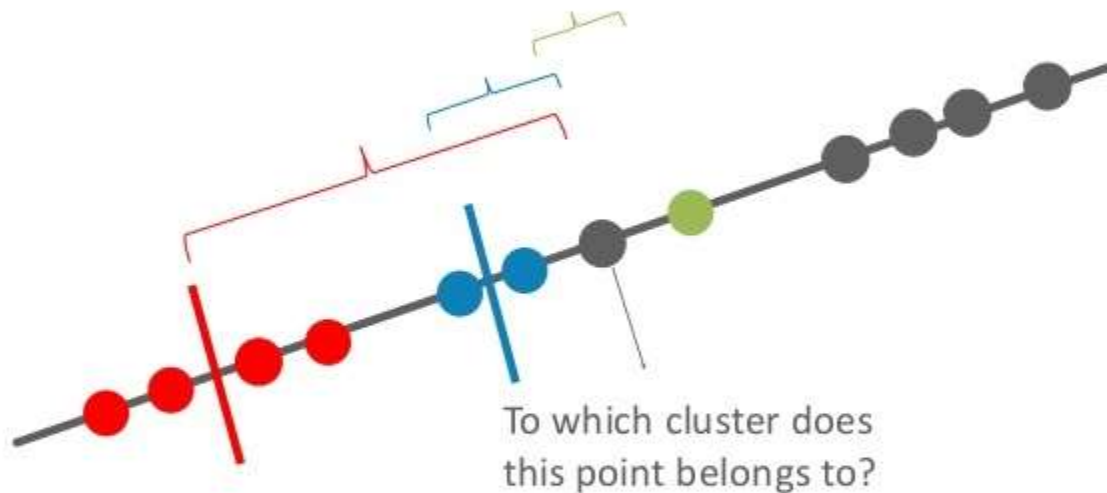
K-Means Algorithm

- Measure the distance
- Assign the point to the nearest cluster
- Calculate the cluster mean using the new point



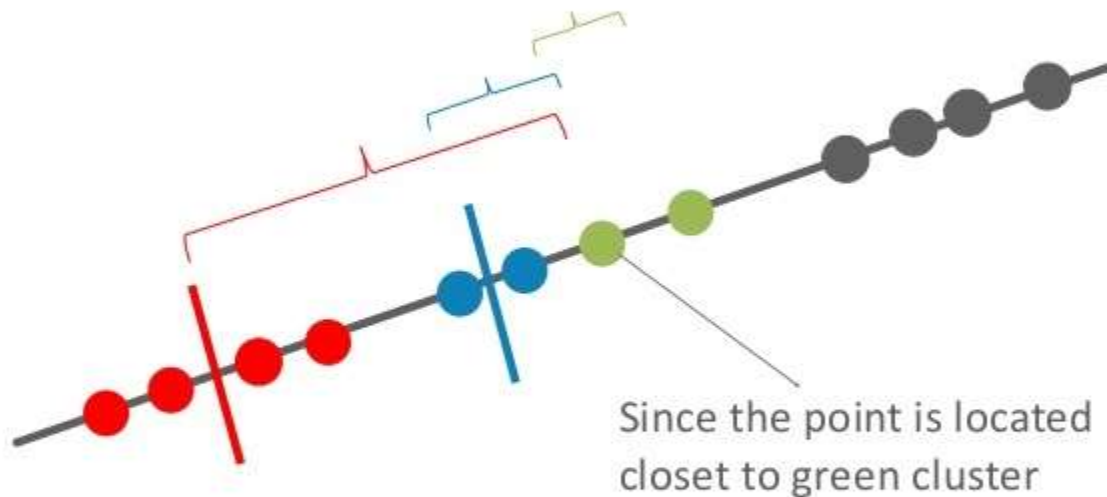
K-Means Algorithm

- Measure the distance from the cluster mean (centroids)
- Assign the point to the nearest cluster
- Calculate the cluster mean using the new point



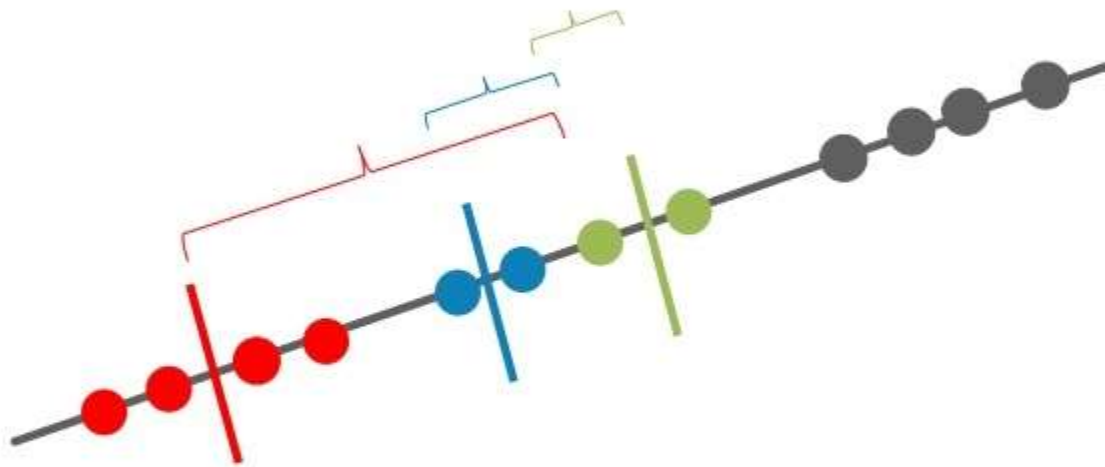
K-Means Algorithm

- Measure the distance from the cluster mean (centroids)
- Assign the point to the nearest cluster
- Calculate the cluster mean using the new point

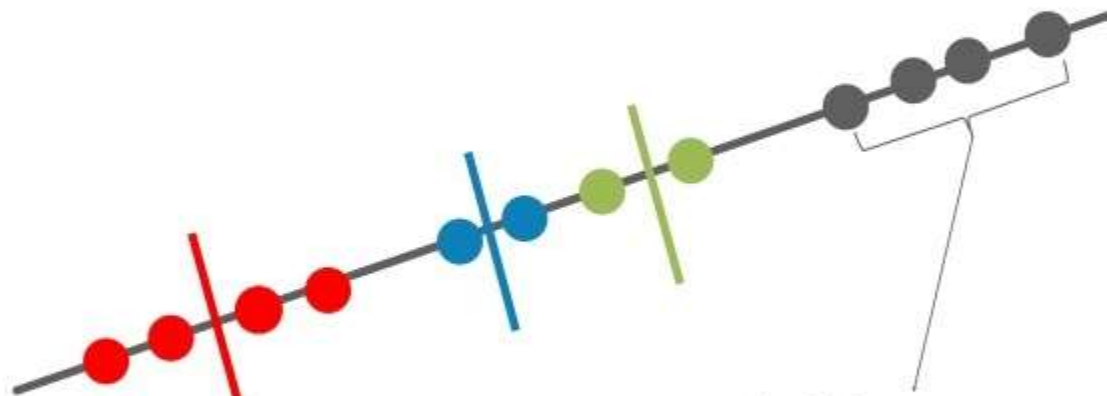


K-Means Algorithm

- Measure the distance from the cluster mean (centroids)
- Assign the point to the nearest cluster
- Calculate the cluster mean using the new point



K-Means Algorithm



Since all of these points are located closest to **green** cluster, so all of them will be assigned to green cluster

K-Means Algorithm

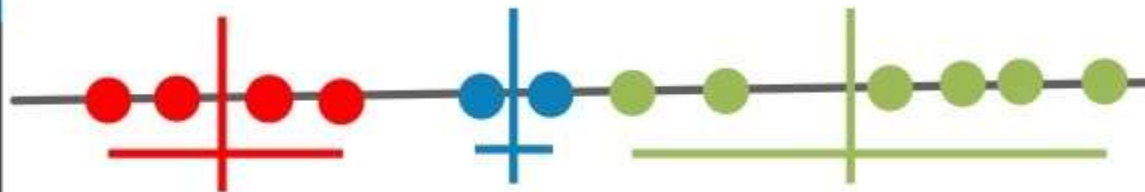


Result from 1st
iteration



Original/Expected Result

K-Means Algorithm

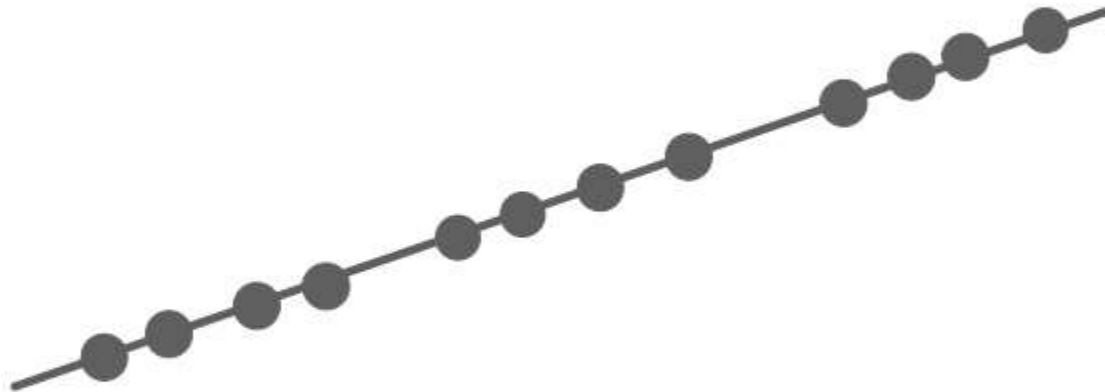


Total variation within the cluster

According to the K-Means Algorithm it iterates over again and again until the data points within each cluster stops changing.

K-Means Algorithm

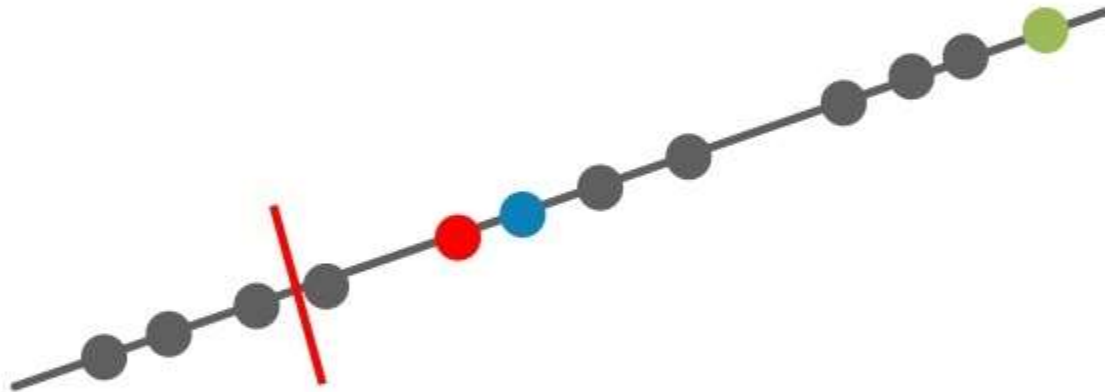
Iteration 2: Again we will start from the beginning. But this time we will be selecting different initial random point (as compared to what we chose in the 1st iteration)



- **Step 1:** Select the number of clusters to be identified, i.e. $K = 3$ in this
- **Step 2:** Randomly select 3 distinct data point
- **Step 3:** Measure the distance between the 1st point and selected 3 clusters

K-Means Algorithm

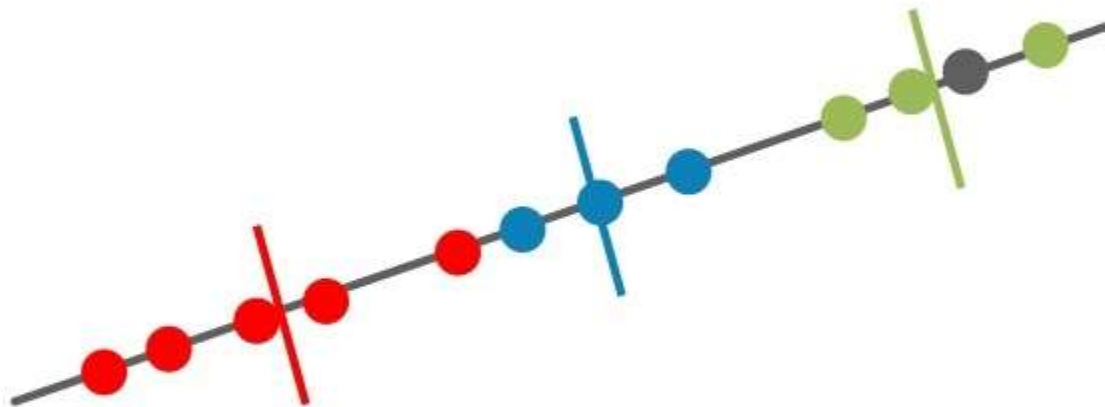
Algorithm picks 3 initial clusters and adds the remaining points to the cluster with the nearest mean, and again recalculating the mean each time a new point is added to the cluster



The diagram illustrates a sequence of data points along a diagonal line. The points are colored red, blue, grey, and green. Vertical lines segment the sequence into groups: a red line after the first red point, a blue line after the blue point, and a green line after the last green point. This represents a partitioning of the data into segments.

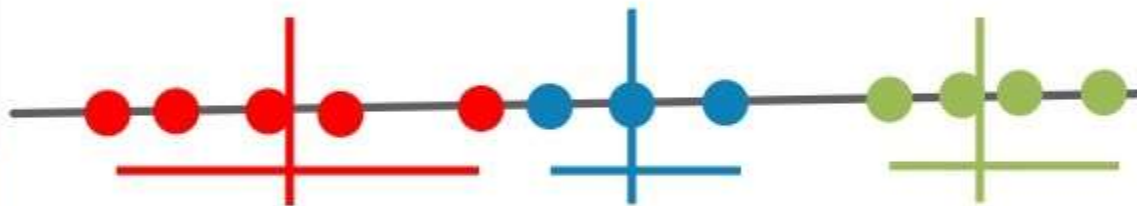
K-Means Algorithm

Algorithm picks 3 initial clusters and adds the remaining points to the cluster with the nearest mean, and again recalculating the mean each time a new point is added to the cluster



Algorithm picks 3 initial clusters and adds the remaining points to the cluster with the nearest mean, and again recalculating the mean each time a new point is added to the cluster

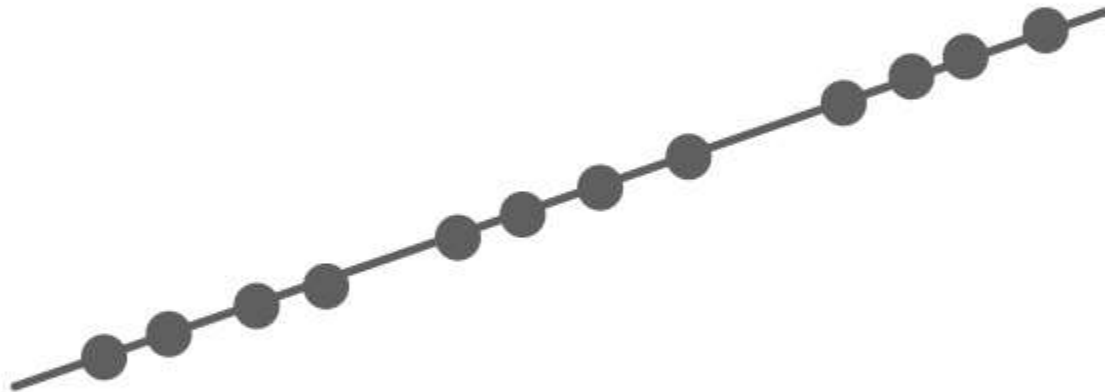
K-Means Algorithm



Total variation within the cluster

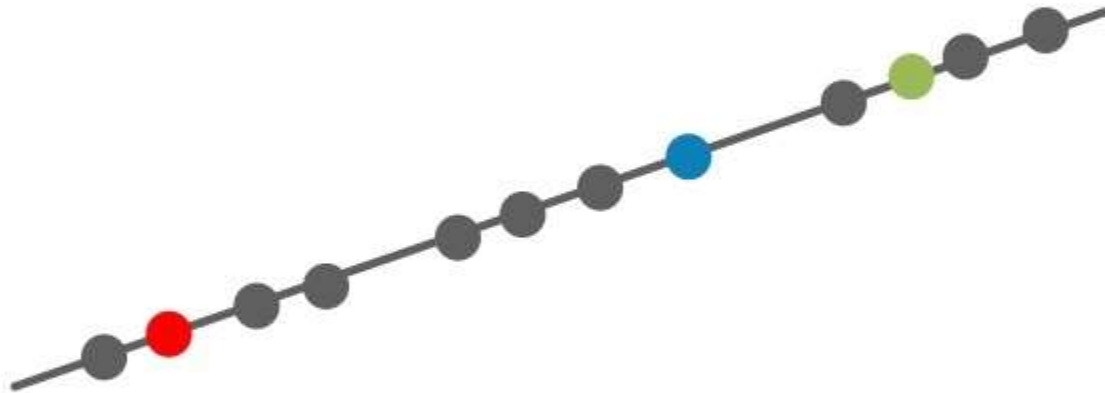
K-Means Algorithm

Iteration 3: Again we will start from the beginning and select different initial random point (as compared to what we chose in the 1st and 2nd iteration)



Pick 3 initial clusters

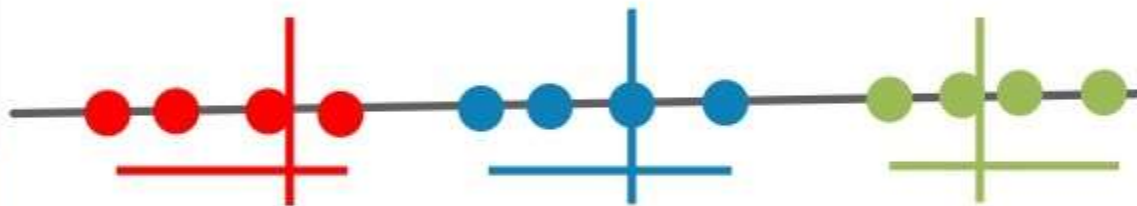
K-Means Algorithm



Cluster the remaining points

K-Means Algorithm

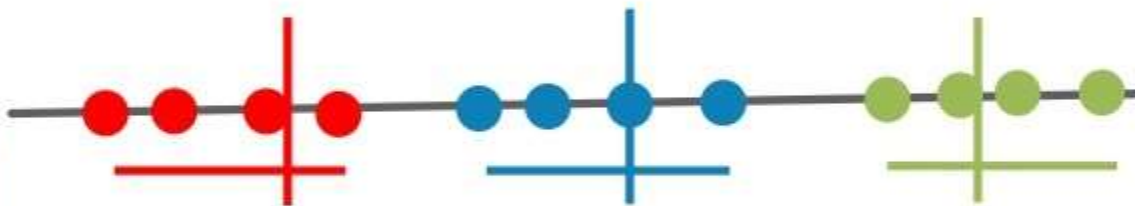
Finally sum the variation within each cluster



Total variation within the cluster

K-Means Algorithm

The algorithm can now compare the result and select the best variance out of it



1st Iteration



2nd Iteration

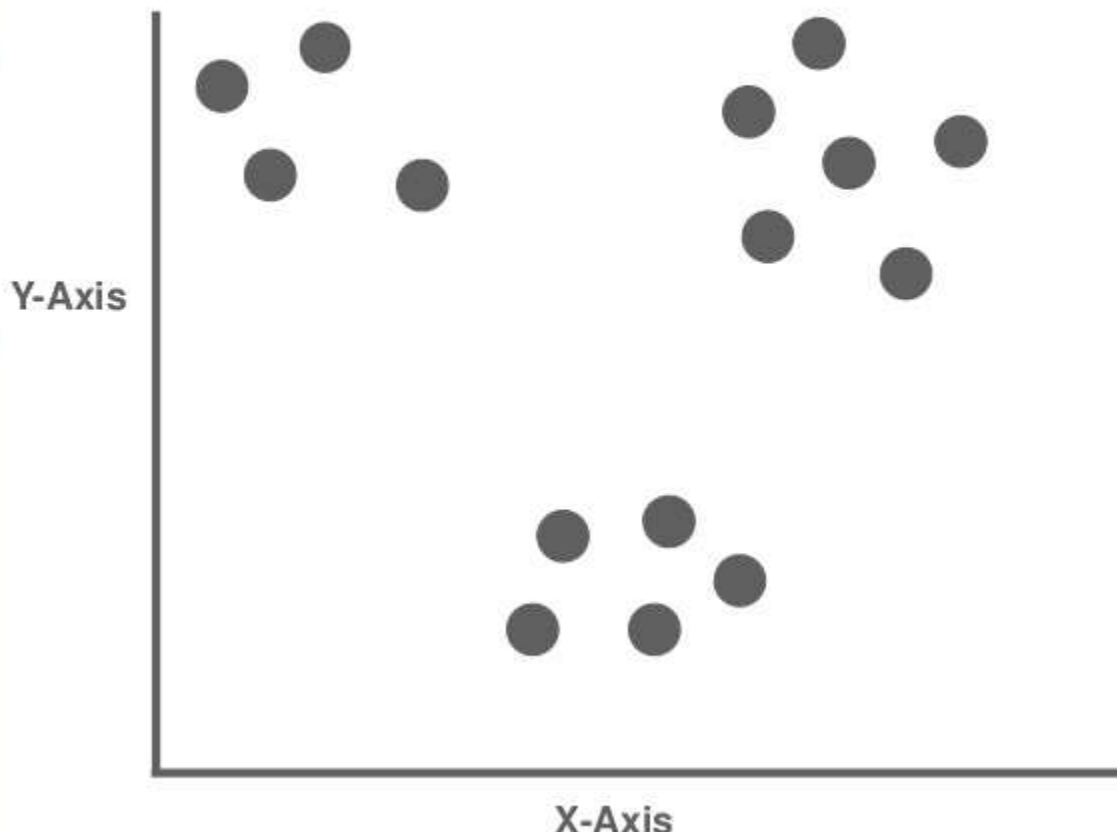


3rd Iteration



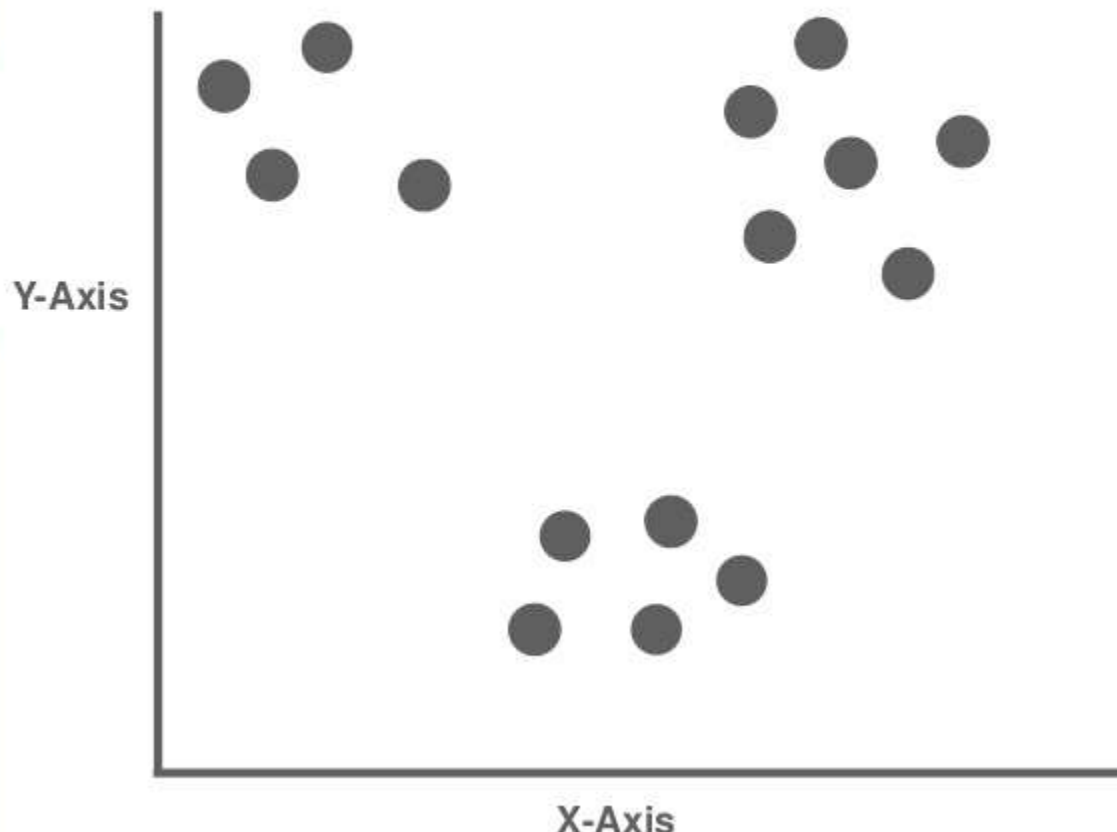
Now what if we have our data plotted on the X and Y axis

K-Means Algorithm



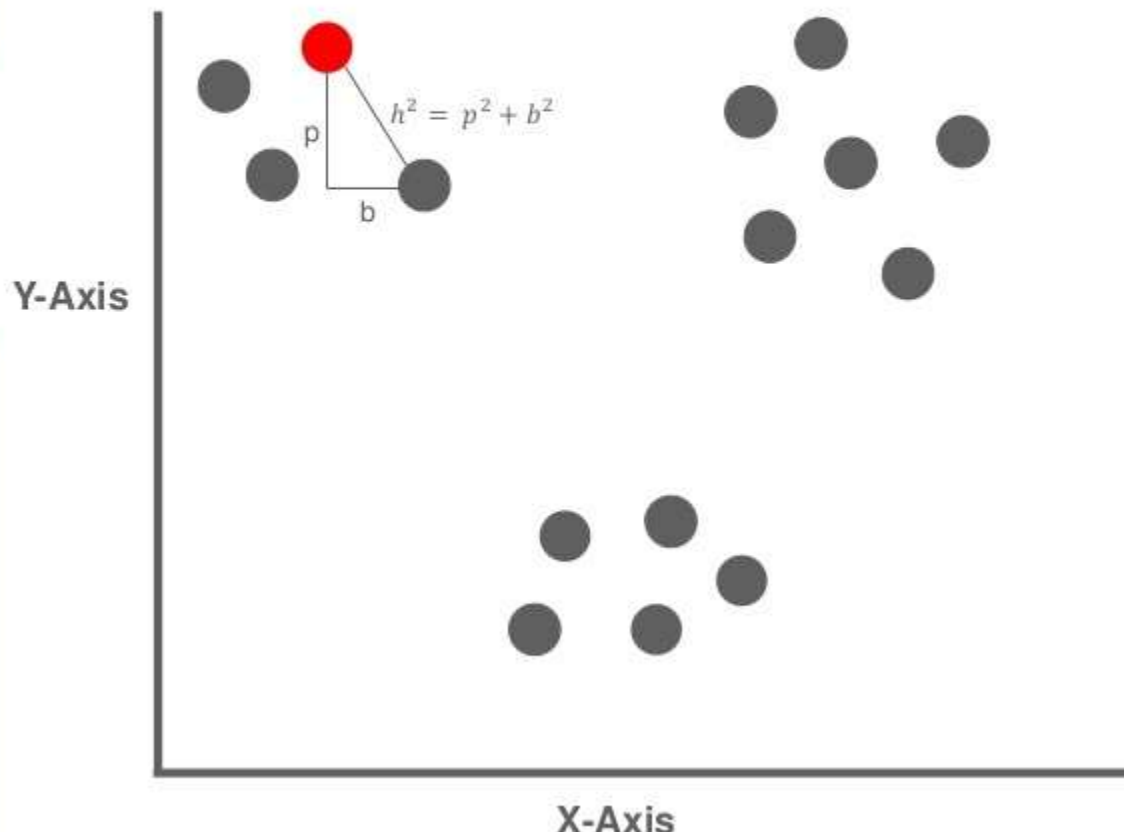
Similarly, pick initial 3 random points..

K-Means Algorithm



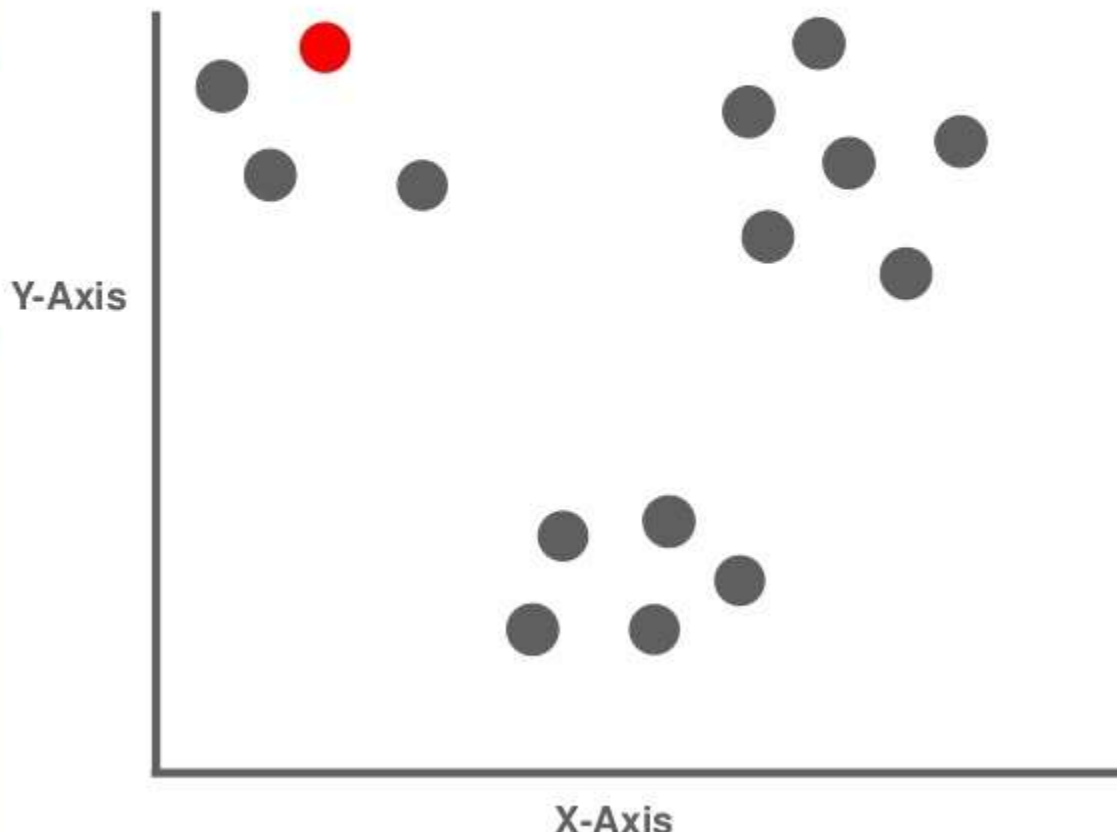
K-Means Algorithm

We will be using the Euclidean distance (in 2D its same as that of a Pythagorean Theorem)



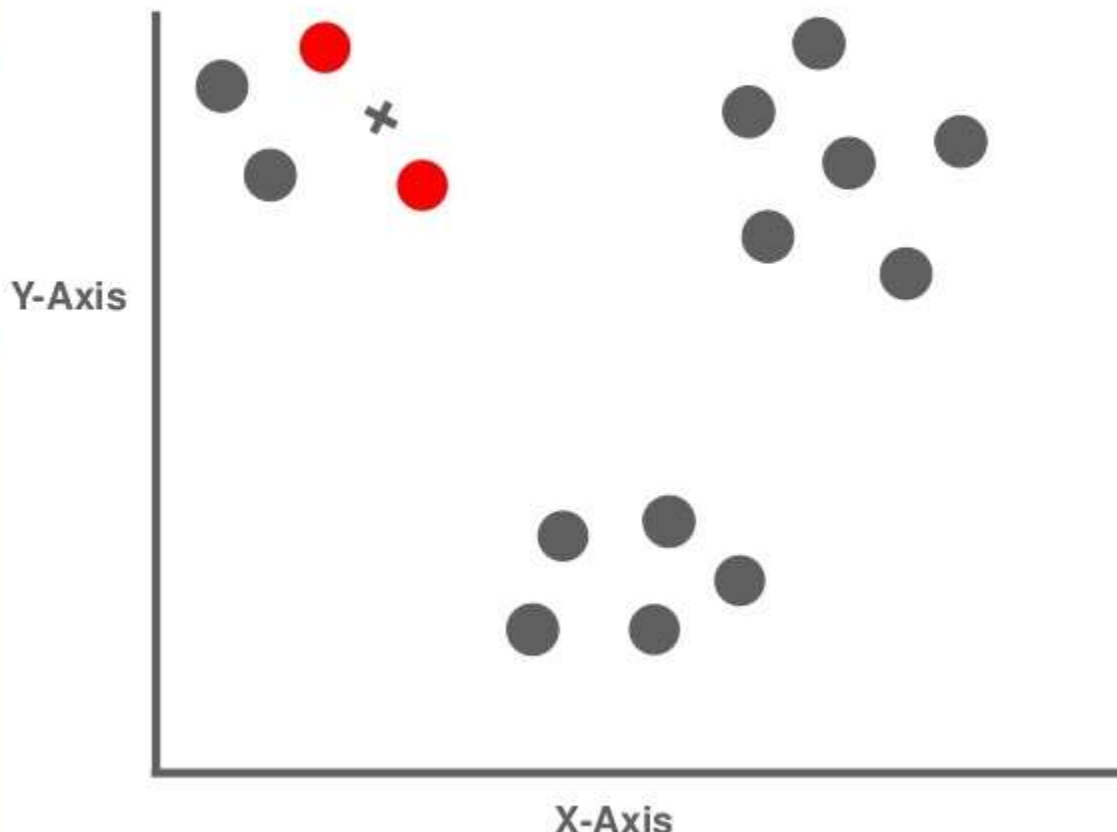
K-Means Algorithm

Again assign the point to the nearest cluster



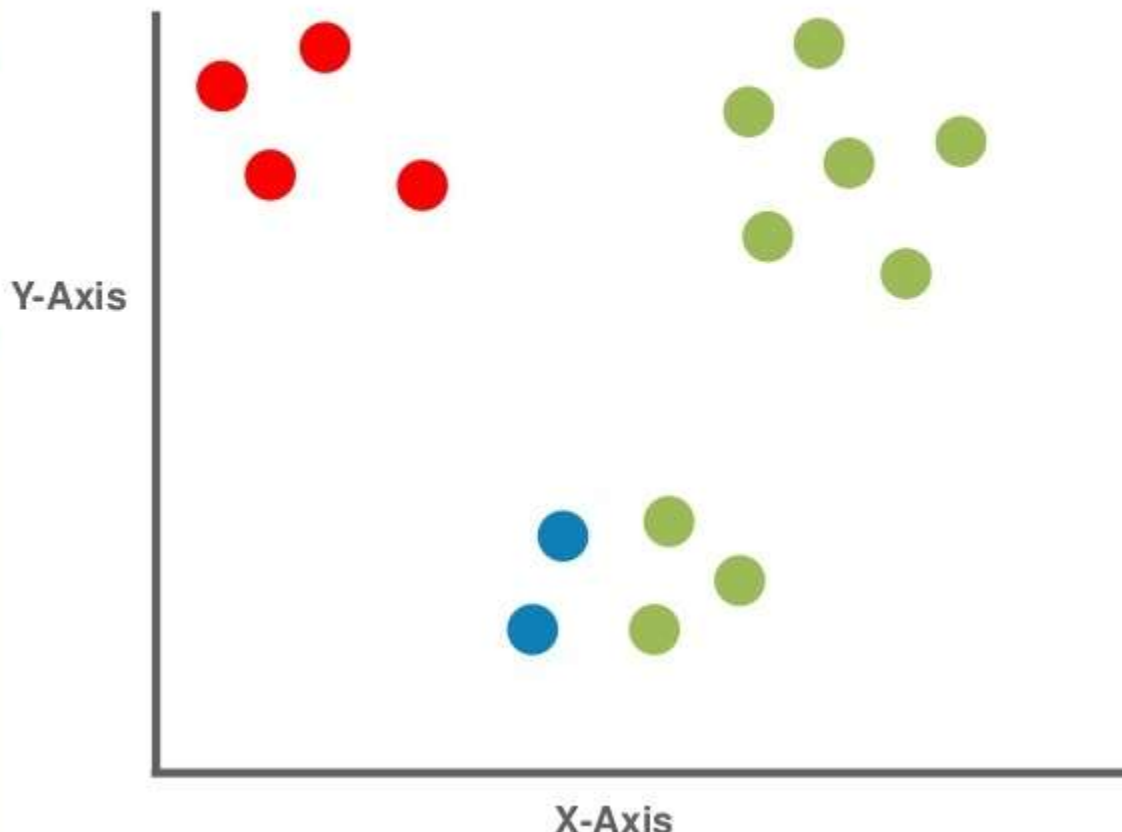
K-Means Algorithm

Finally calculate the centroid (mean of cluster)
including the new point



K-Means Algorithm

Finally in first iteration you get something like this...again
you have to iterate this process to get the final cluster



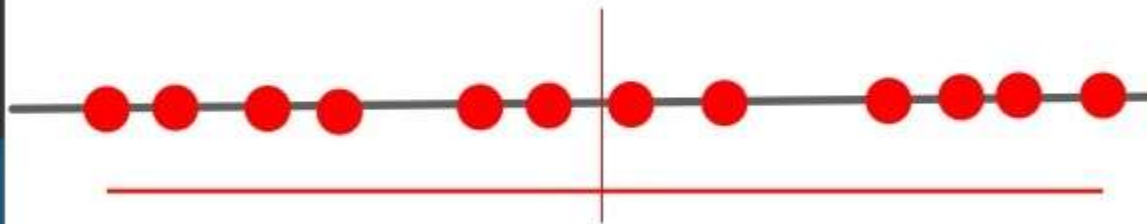
In the previous scenario k value was known to be 3, but this is not always true

How will you find K value



How will you find K value

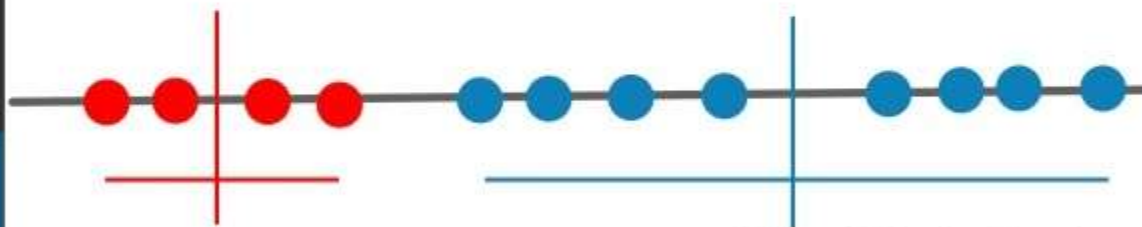
For deciding the value of k , you have to use hill climbing and trail method, starting from $K = 1$



$K=1$ is the worst case scenario, even you cross-verify it with total variation

How will you find K value

Now try with $K = 2$



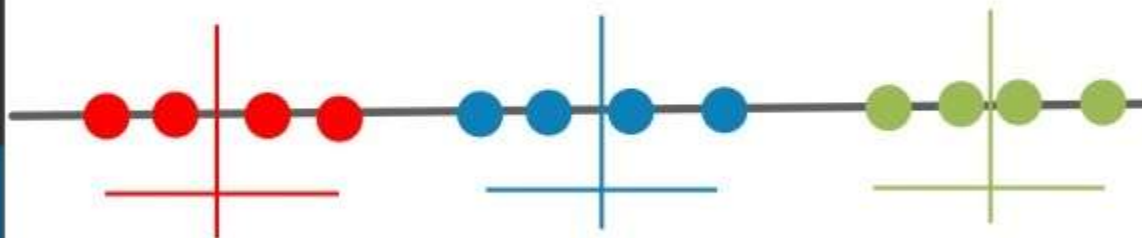
$K=2$ is still better than $K = 1$ (Total Variation)

$K = 1$

$K = 2$

How will you find K value

Now try with $K = 3$



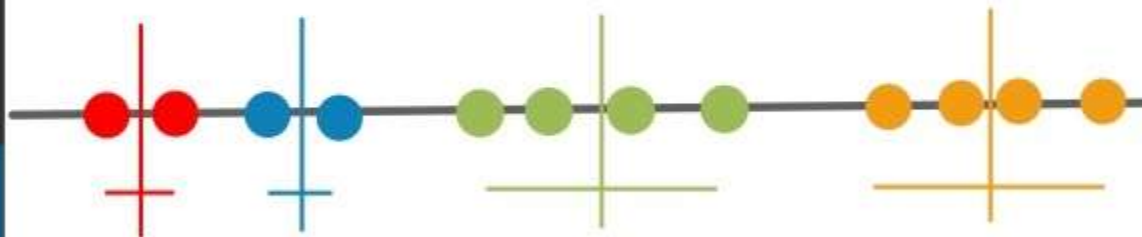
$K=3$ is even better than $K=2$ (Total Variation)

$K = 1$

$K = 2$

How will you find K value

Now try with $K = 4$



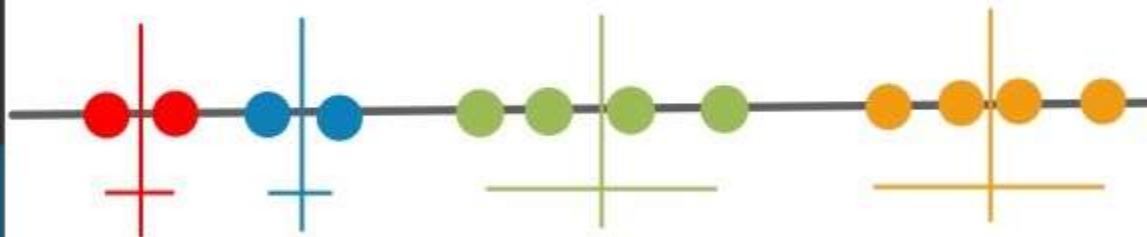
Total variation in $K=4$ is less than $K=3$



How will you find K value

Now try with $K = 4$

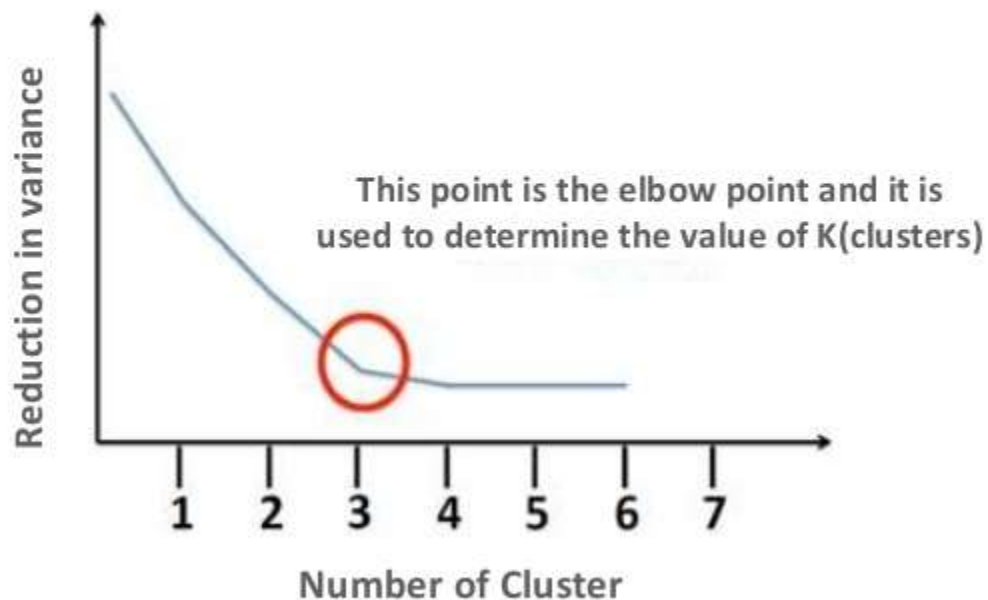
Each time you increase the cluster the variation decreases, no. of clusters = no. of data points then in that case the variation = 0



Total variation in $K=4$ is less than $K=3$



How will you find K value





**Let's learn to
code**

K-Means Algorithm

Summarizing the K-Means Algorithm

```
randomly chose k examples as initial centroids
while true:
    create k clusters by assigning each
        example to closest centroid
    compute k new centroids by averaging
        examples in each cluster
    if centroids don't change:
        break
```

600,000 +
SATISFIED LEARNERS
Thank you! 😊

For more information please visit our website
www.edureka.co

