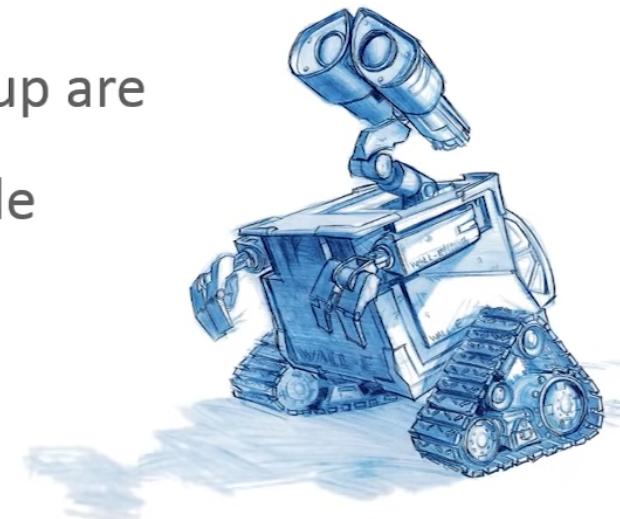


# What is Clustering?

“**Clustering** is the process of dividing the datasets 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 diners  
in a restaurant

Items arranged in  
a mall



# Where is it Used?



Recommendation System



Recommended Movies



Flickr's Photos

# How business use Clustering?



Banking



Retail Store



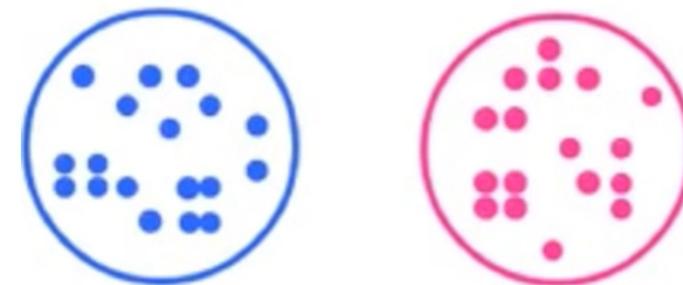
Insurance  
Companies

# 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

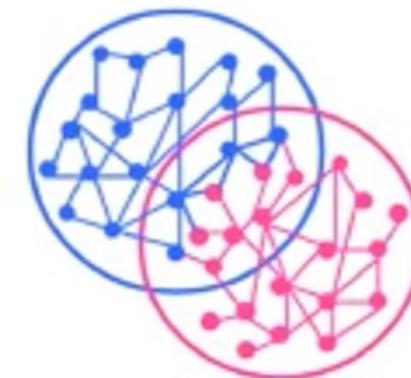


# Types of Clustering

- 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**

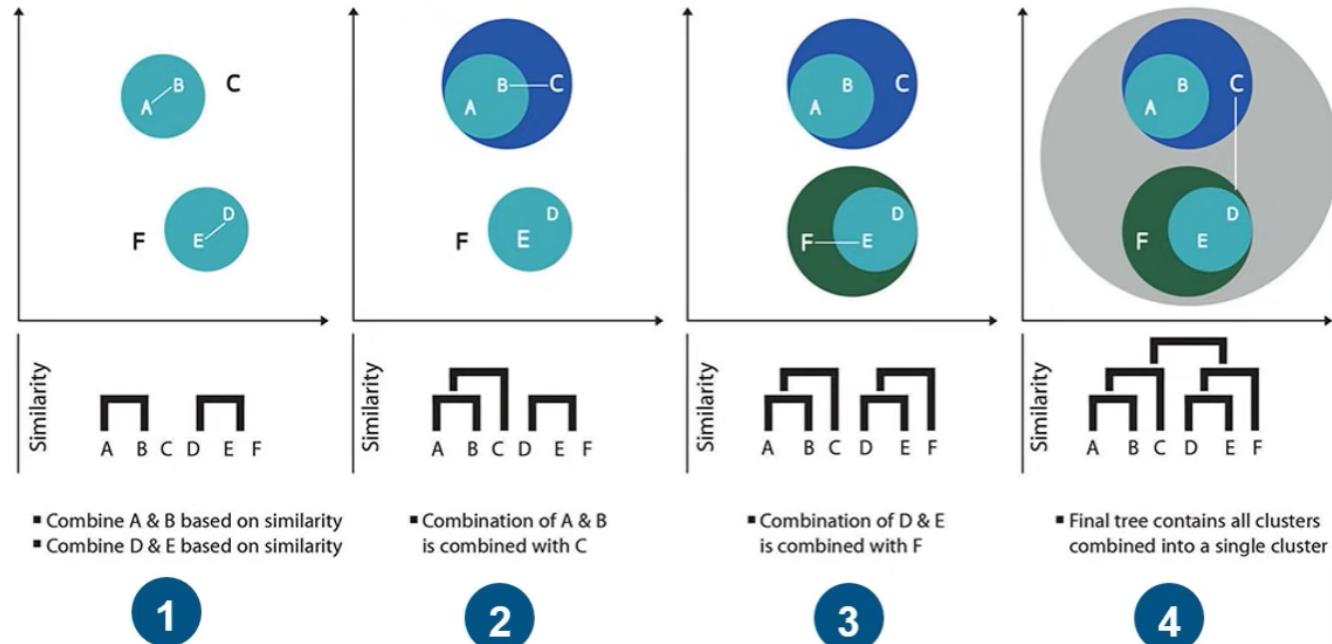


**Overlapping Clustering**



**Hierarchical Clustering**

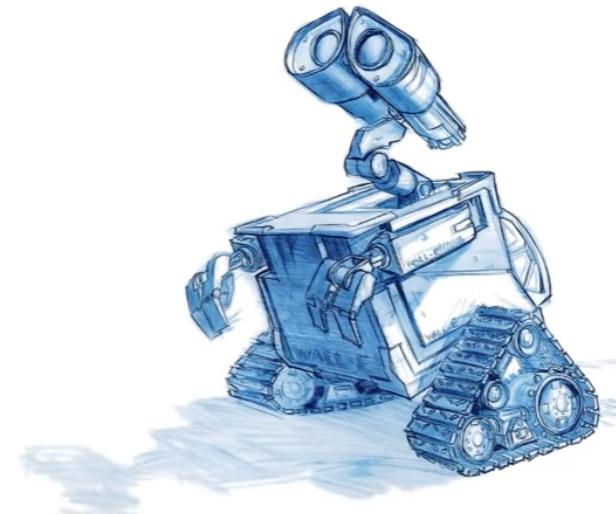
## Hierarchical Clustering



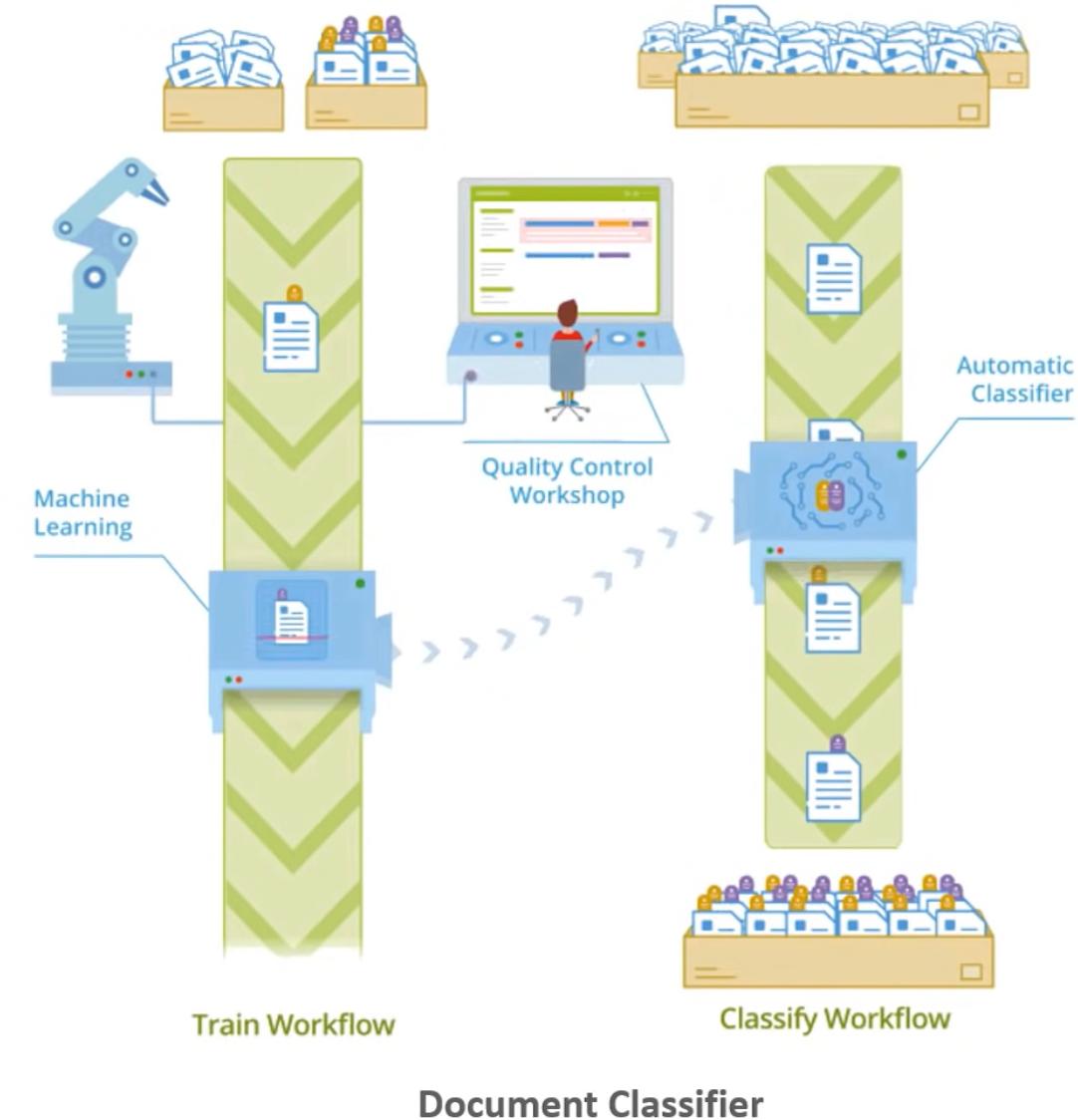
# What is K-Means Clustering?

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

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



# Where Can I apply K-Means?



# 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:** Move the cluster mean towards the center of the cluster  
and repeat step 2



# K-Means Algorithm

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



# K-Means Algorithm

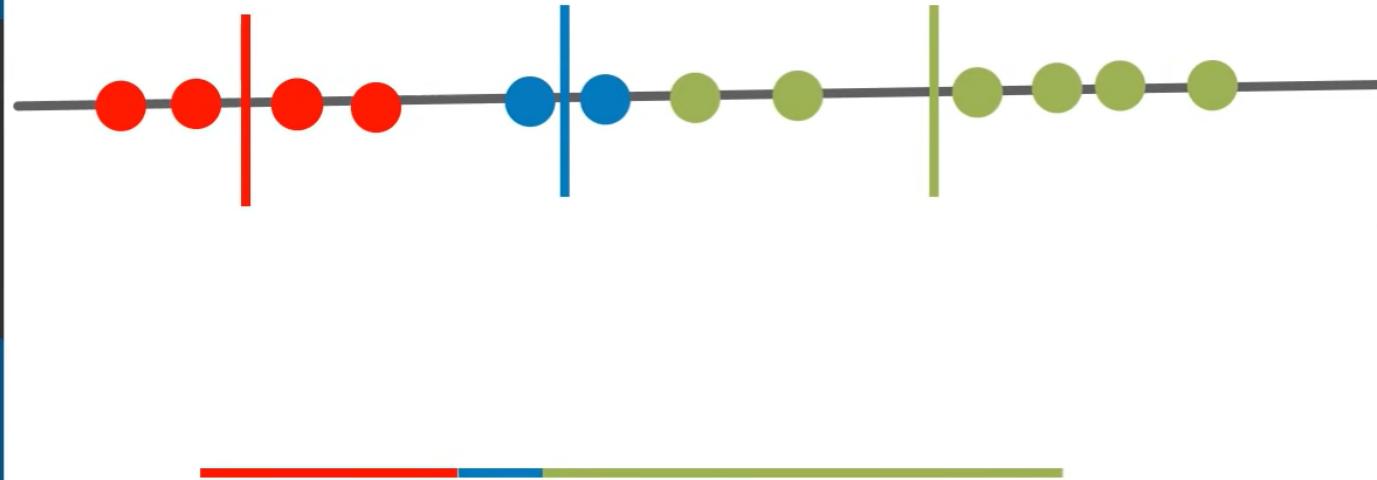


Result from 1<sup>st</sup>  
iteration



Original/Expected Result

# K-Means Algorithm

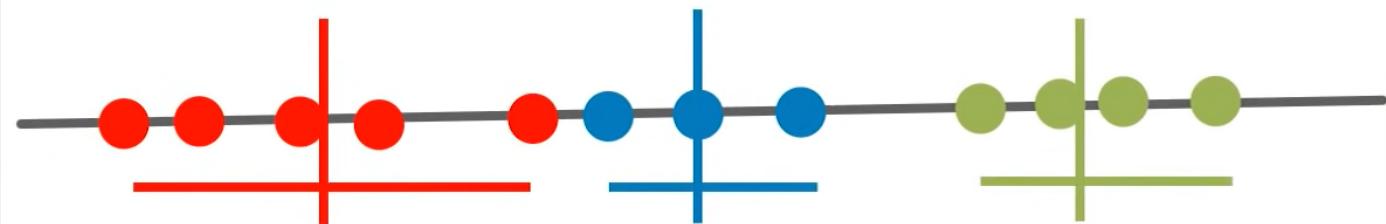


Total variation within the cluster

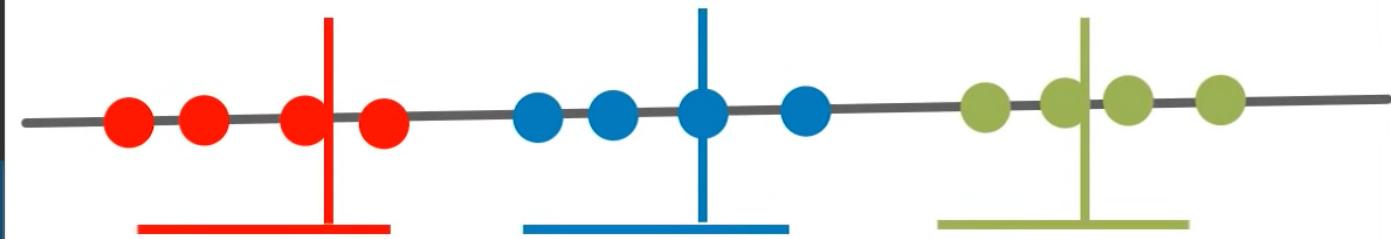
According to the K-Means Algorithm it iterates over again and again unless and until the data points within each cluster stops changing

# 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

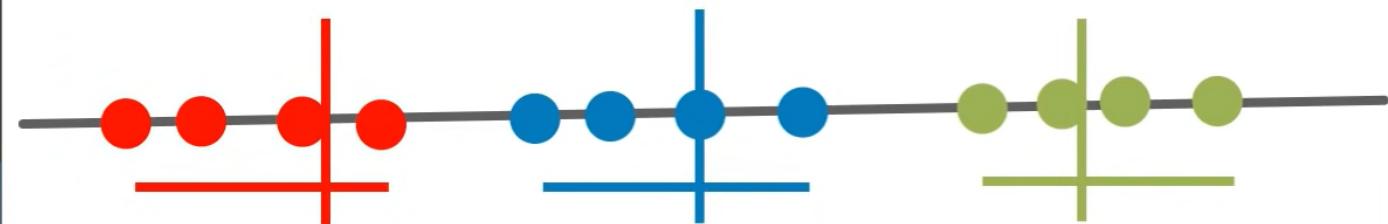


# K-Means Algorithm



# K-Means Algorithm

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



1<sup>st</sup> Iteration

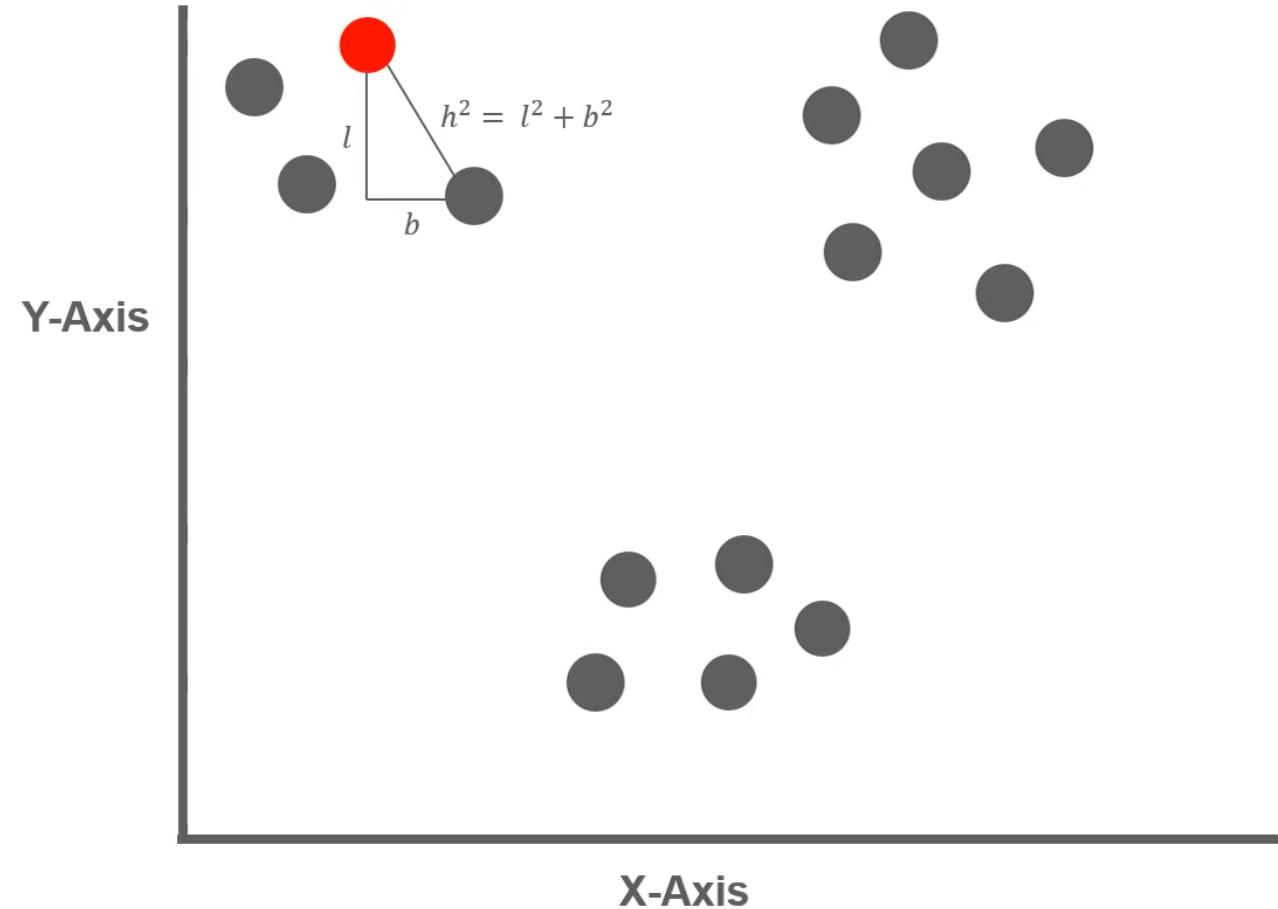
2<sup>nd</sup> Iteration

3<sup>rd</sup> Iteration



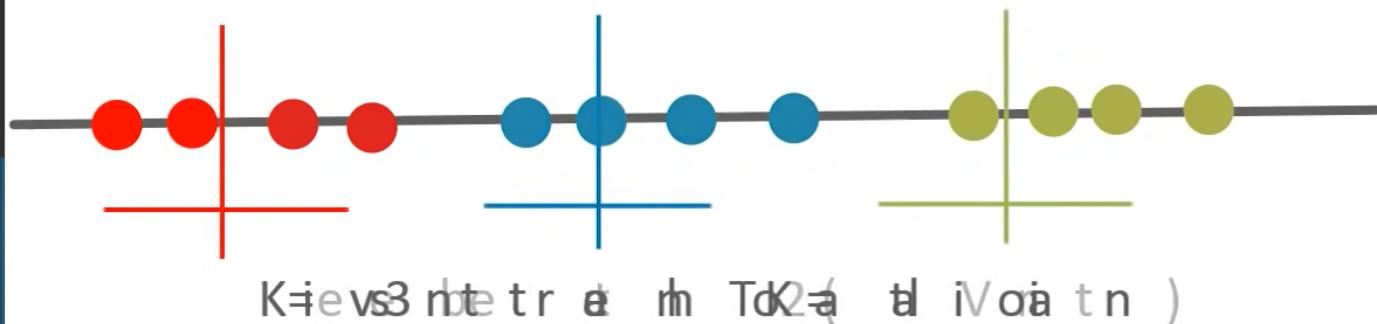
# K-Means Algorithm

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



# How will you find K value

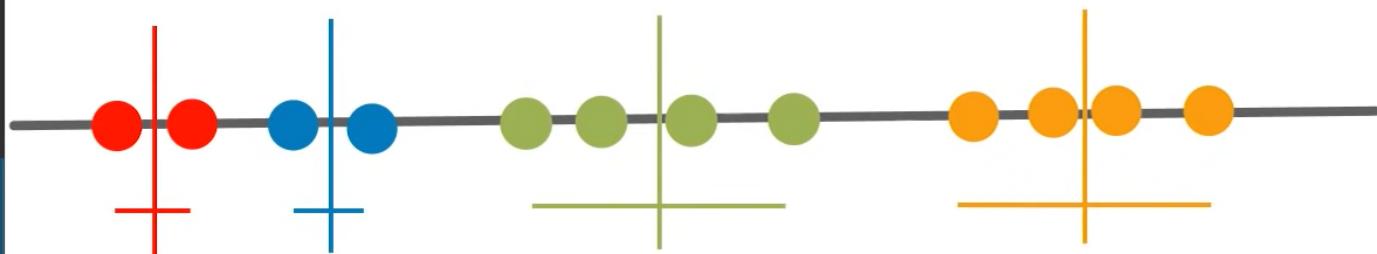
Now try with K = 3



# 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

