



# Oriental Institute of Science and Technology, Bhopal

## Department of Computer Science and Engineering

# Progress Seminar - I

on

## FULL STACK SCRUTINY

(Minor Project II)

Session 2022



**Guided By :** Prof. Goldi Jarbais

**Presented by:**

1. Anshul Verma (0105CS191023) [Team Leader]
2. Ayush Waghmare (0105CS191028)
3. Harshit Shrivastava (0105CS191048)
4. Manish Nathrani (0105CS191062)

# OVERVIEW

- Introduction and Detail Idea
- Design Process
- Flowchart
- UML Diagrams
- Description of Modules
- Module 1
- References

# INTRODUCTION & DETAIL IDEA OF PROJECT

All the three things at the same platform.

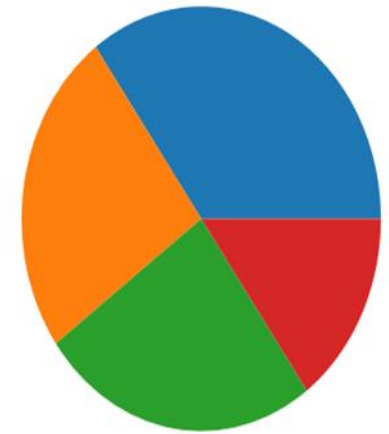
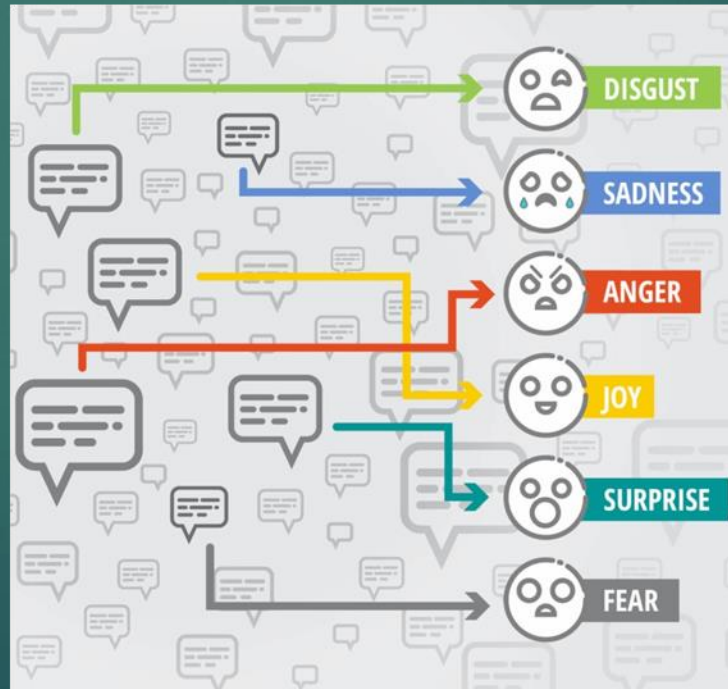
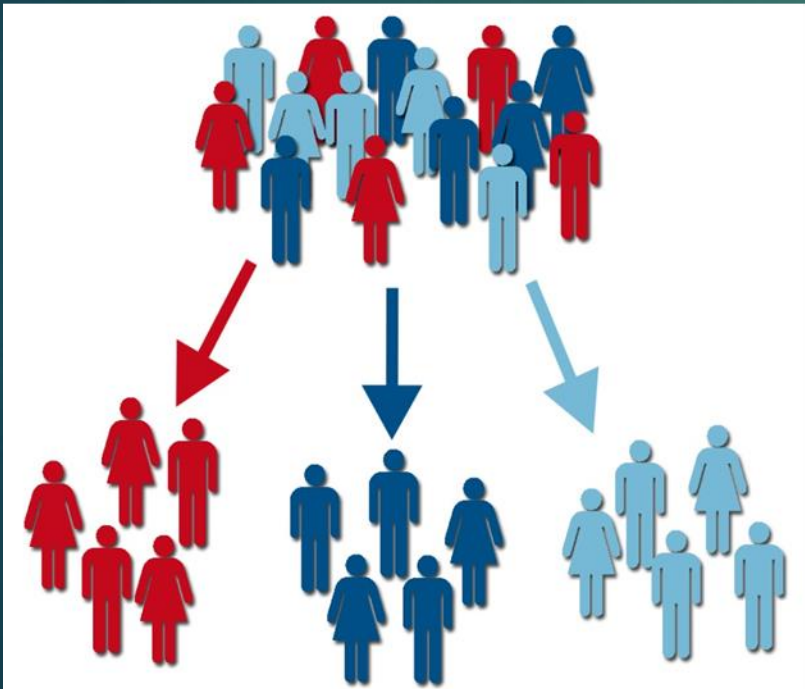
Customer  
Segmentation

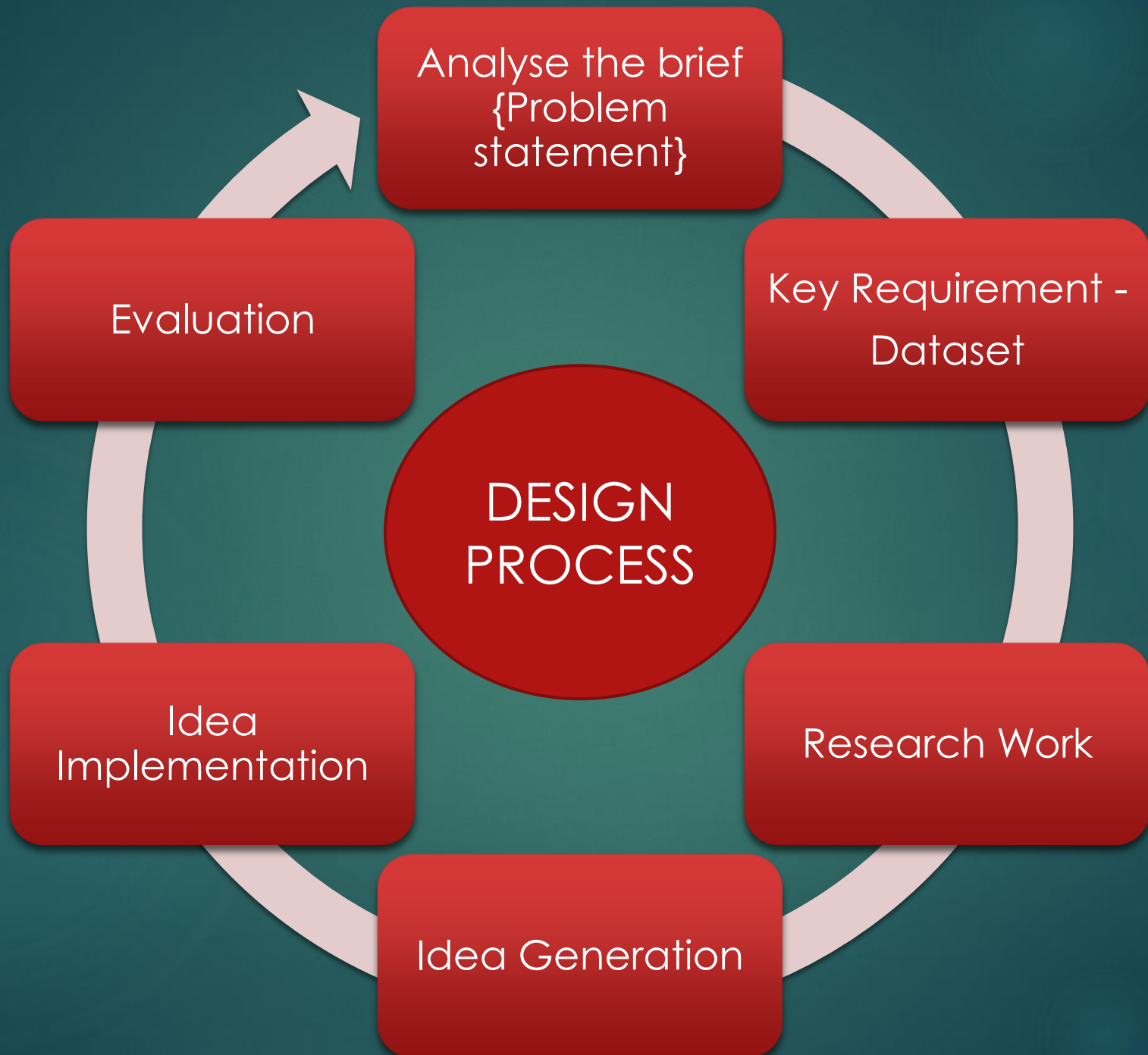


Sentiment Analysis



Exploratory Data  
Analysis





# FLOWCHART

DATA PREPERATION  
STEP 1



DEVELOPEMENT  
STEP 2



DEPLOYMENT  
STEP 3

Collect Data



Clean Data

Design Model



Train Model



Test Model



Optimize Model

Deploy Model



Retrain Model



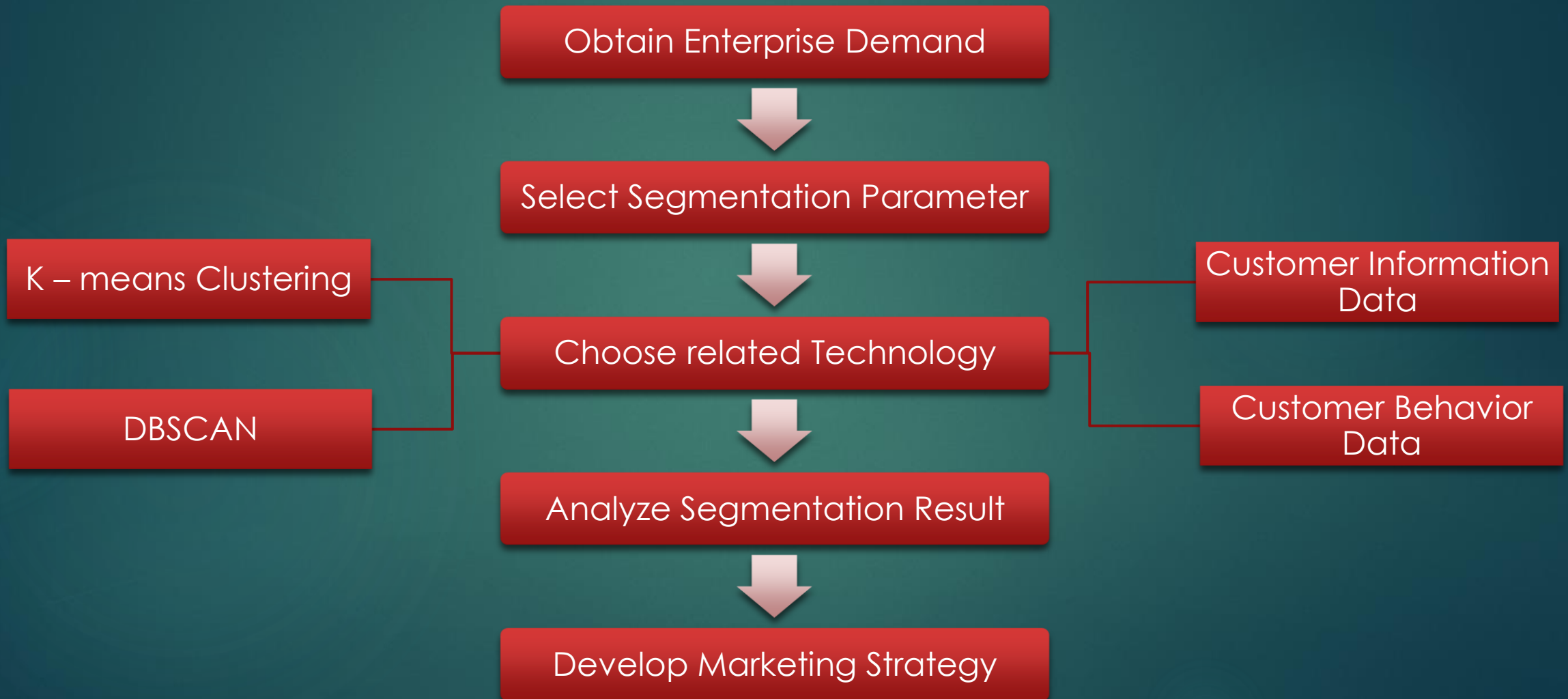
Result

# UML DIAGRAMS

- Customer Segmentation
- Sentiment Analysis



# CUSTOMER SEGMENTATION



# SENTIMENT ANALYSIS

A typical Sentiment  
Analysis Model

Reviews



Data Preparation



Review Analysis



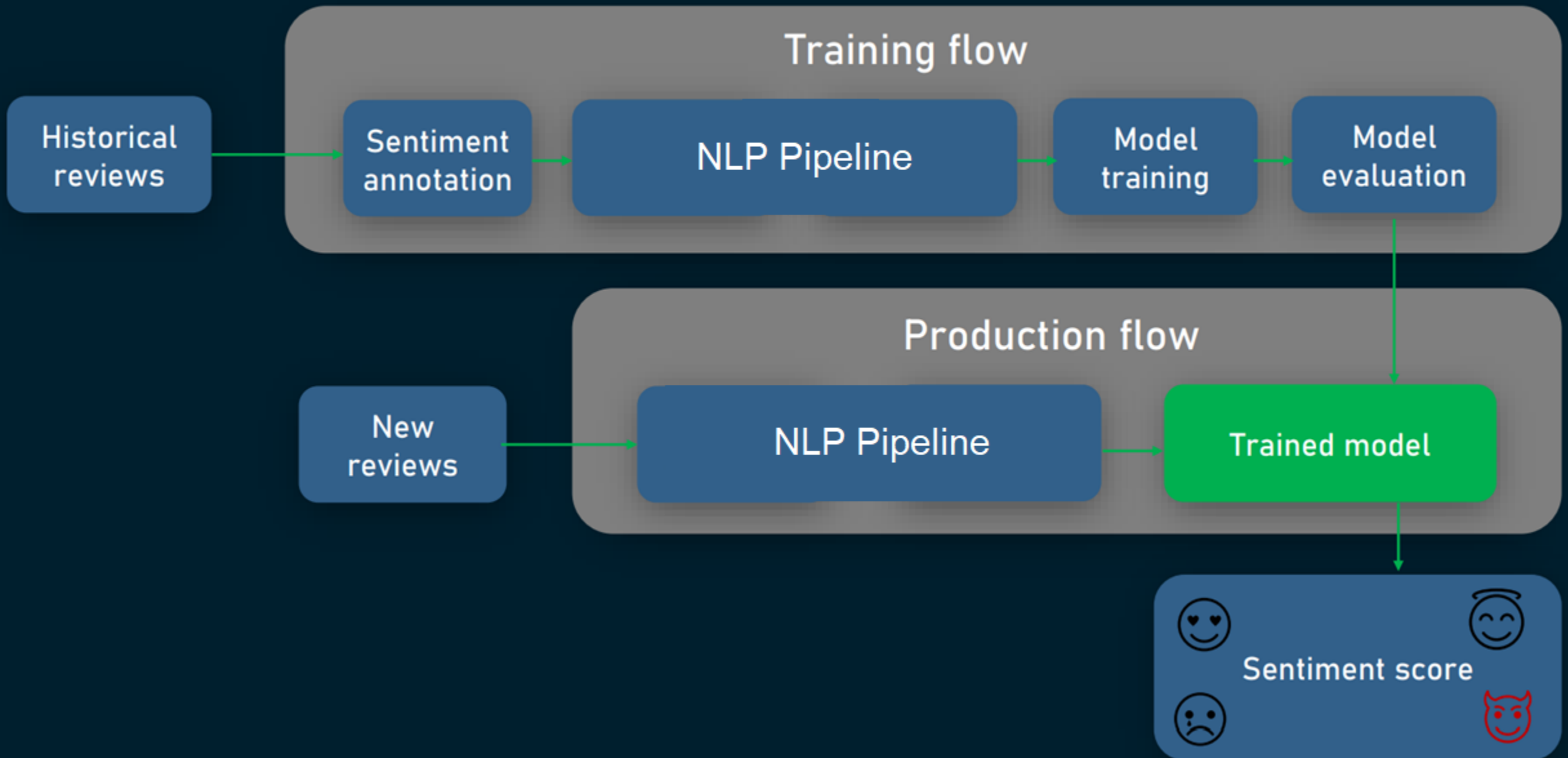
Sentiment Classification



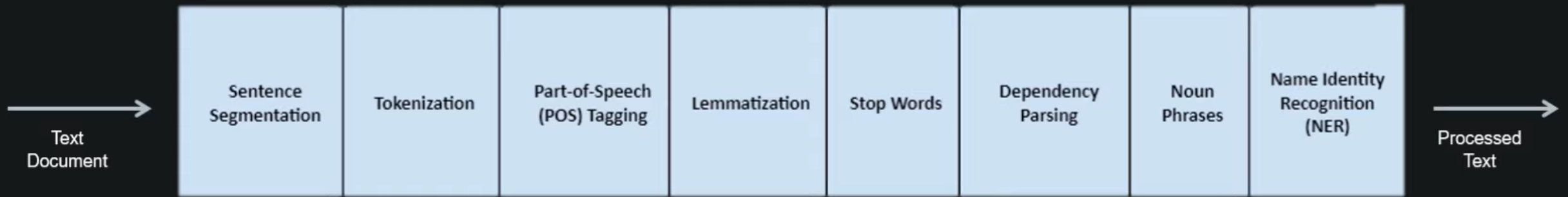
Result



# SENTIMENT ANALYSIS WITH MACHINE LEARNING



# NLP PIPELINE



A pipeline is just a way to design a program where the output of one module feeds to the input of the next. For example, Linux shells feature a pipeline where the output of a command can be fed to the next using the pipe character

# DESCRIPTION OF MODULES

M 1

Customer Segmentation

M 2

Sentiment Analysis

M 3

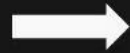
Exploratory Data Analysis

M 4

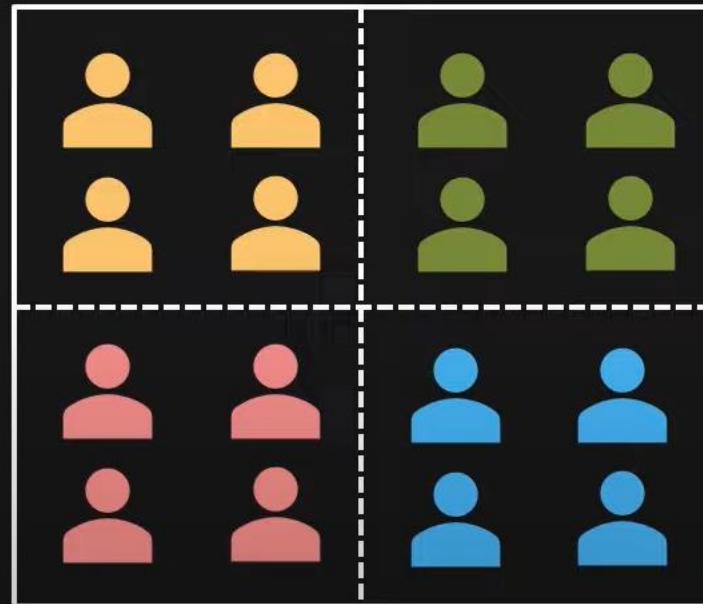
Data Visual Representation & Hypothesis Generation

# MODULE 1 – CUSTOMER SEGMENTATION

Determine customers to sell products



Applying clustering algorithm



Sell product to targeted customers



# CLUSTERING OPTIONS

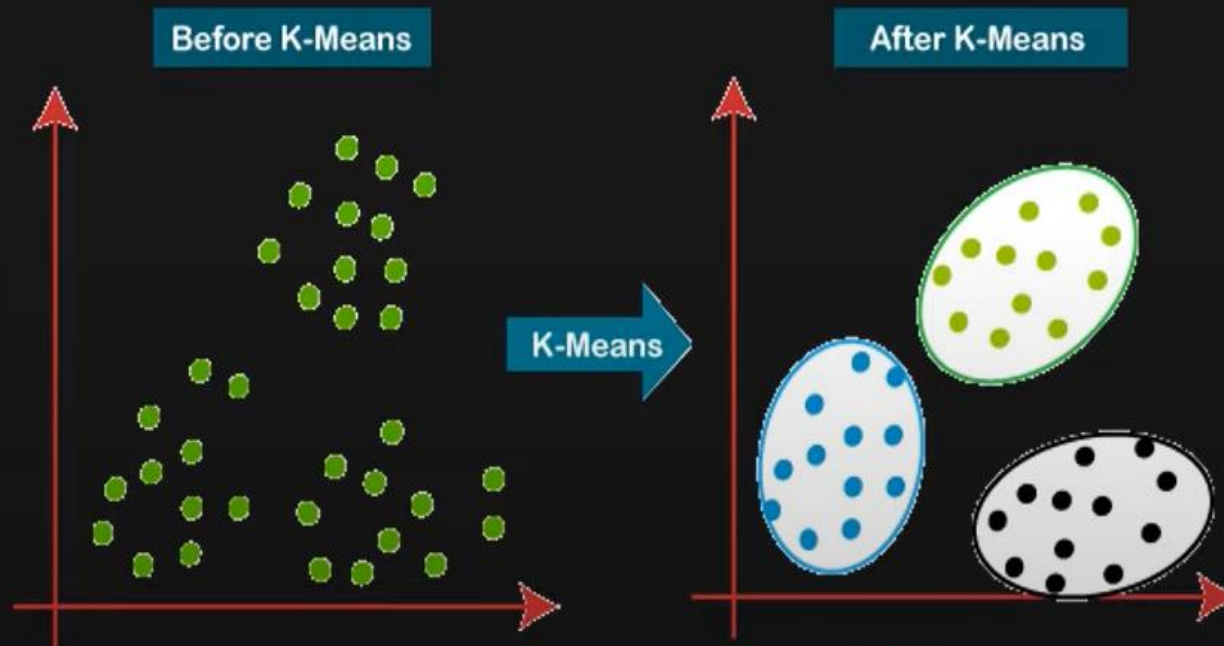


K – means  
Clustering

DBSCAN

# K-means Algorithm

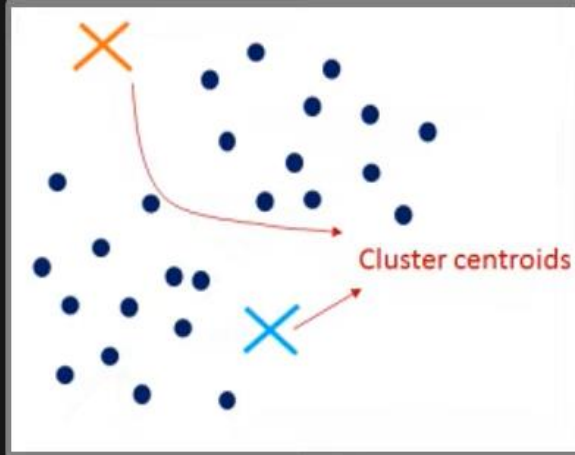
It is an iterative algorithm that divides the unlabelled dataset into  $k$  different clusters in such a way that each dataset belongs only one group that has similar properties.



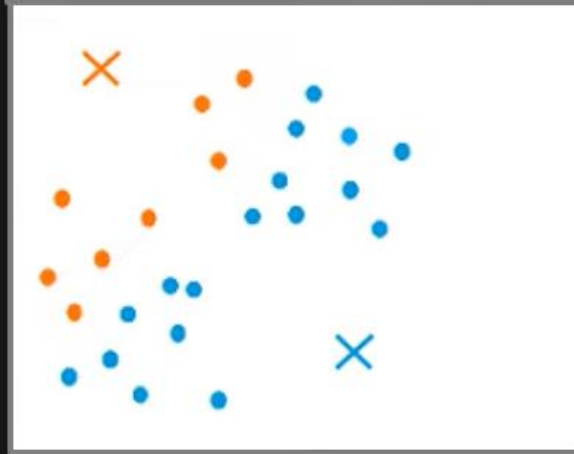


# K-means Algorithm

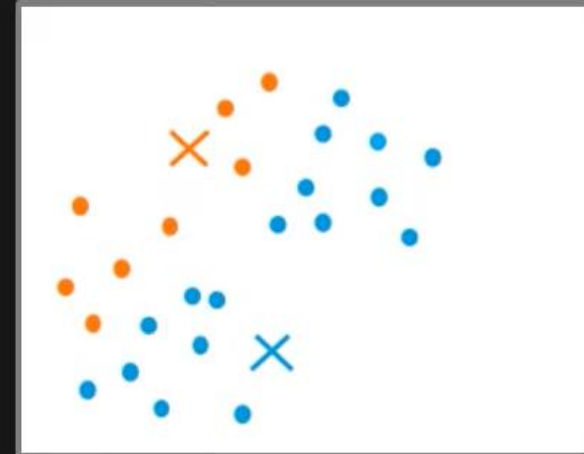
Initialization



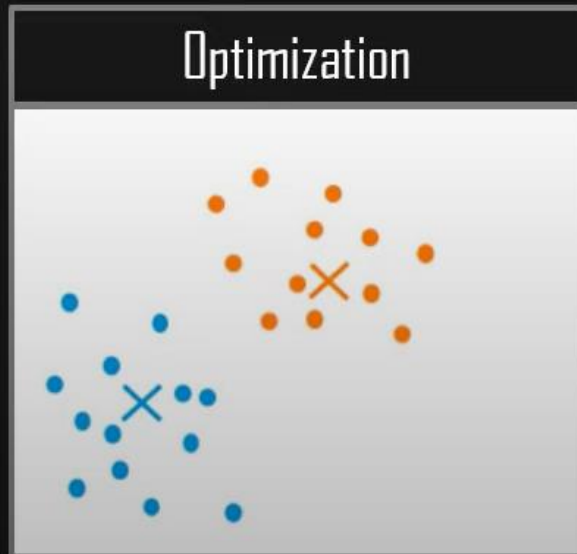
Cluster Assignment



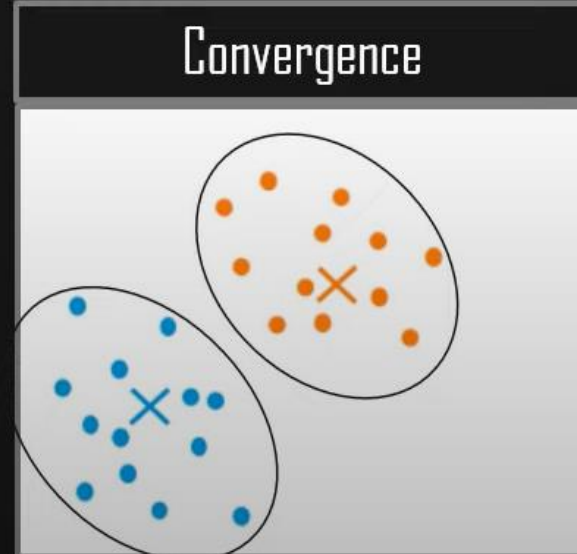
Move Centroid



Optimization



Convergence

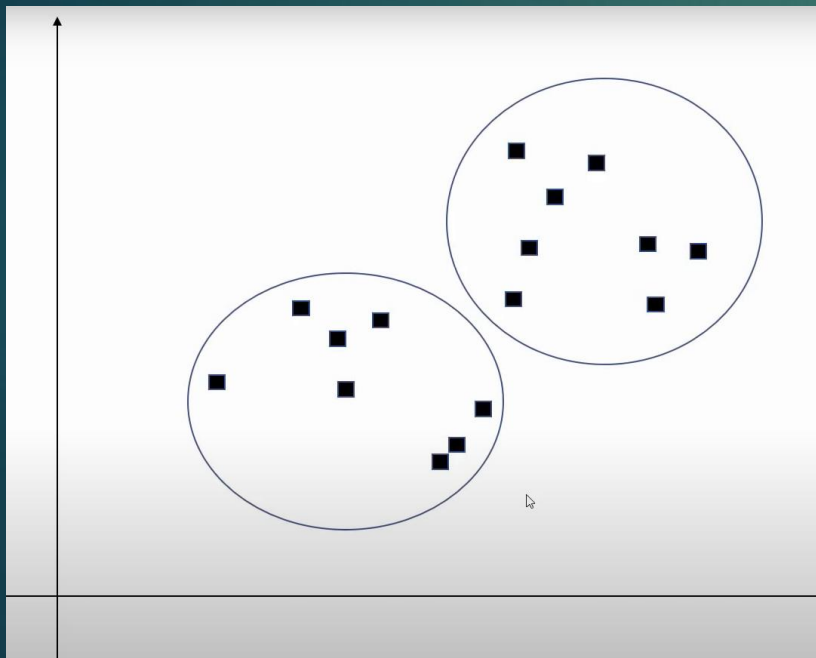


But there is a Catch...

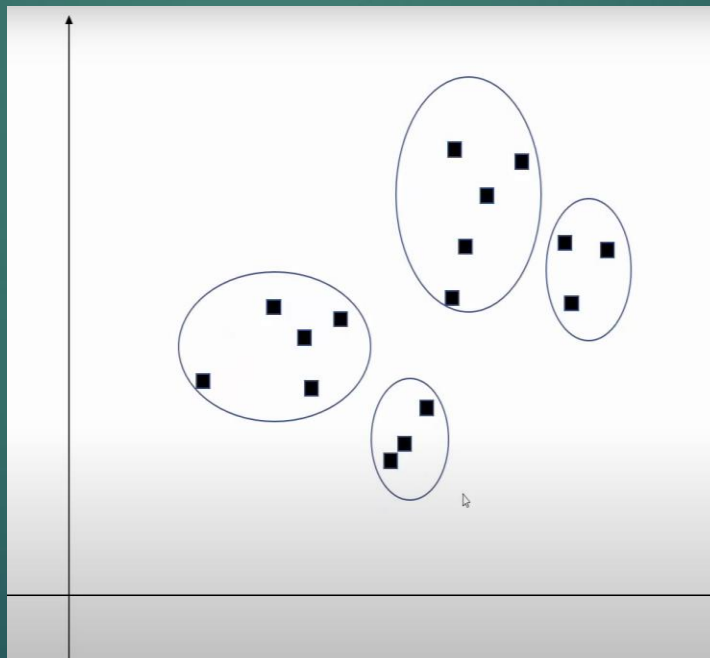
How to determine the Number of  
Clusters ( $K$ )?

# The Problem is...

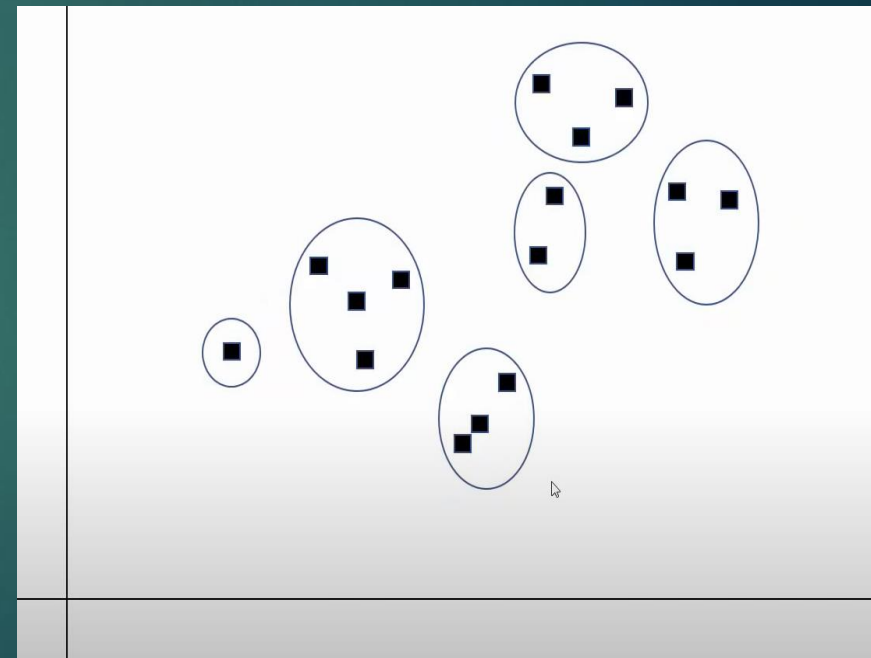
$K = 2$



$K = 4$



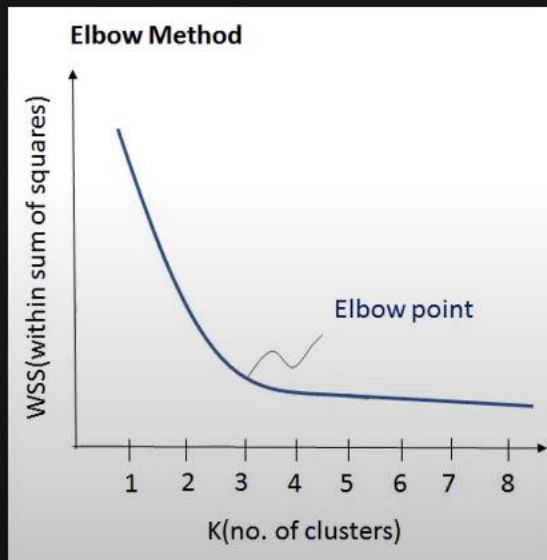
$K = 6$



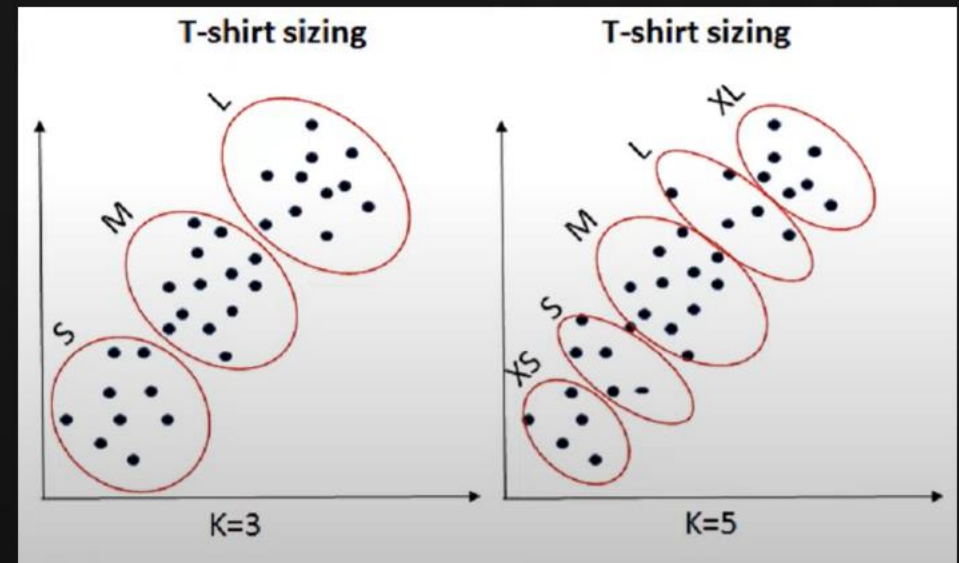
# Solution is...

Ways to choose the optimum K value

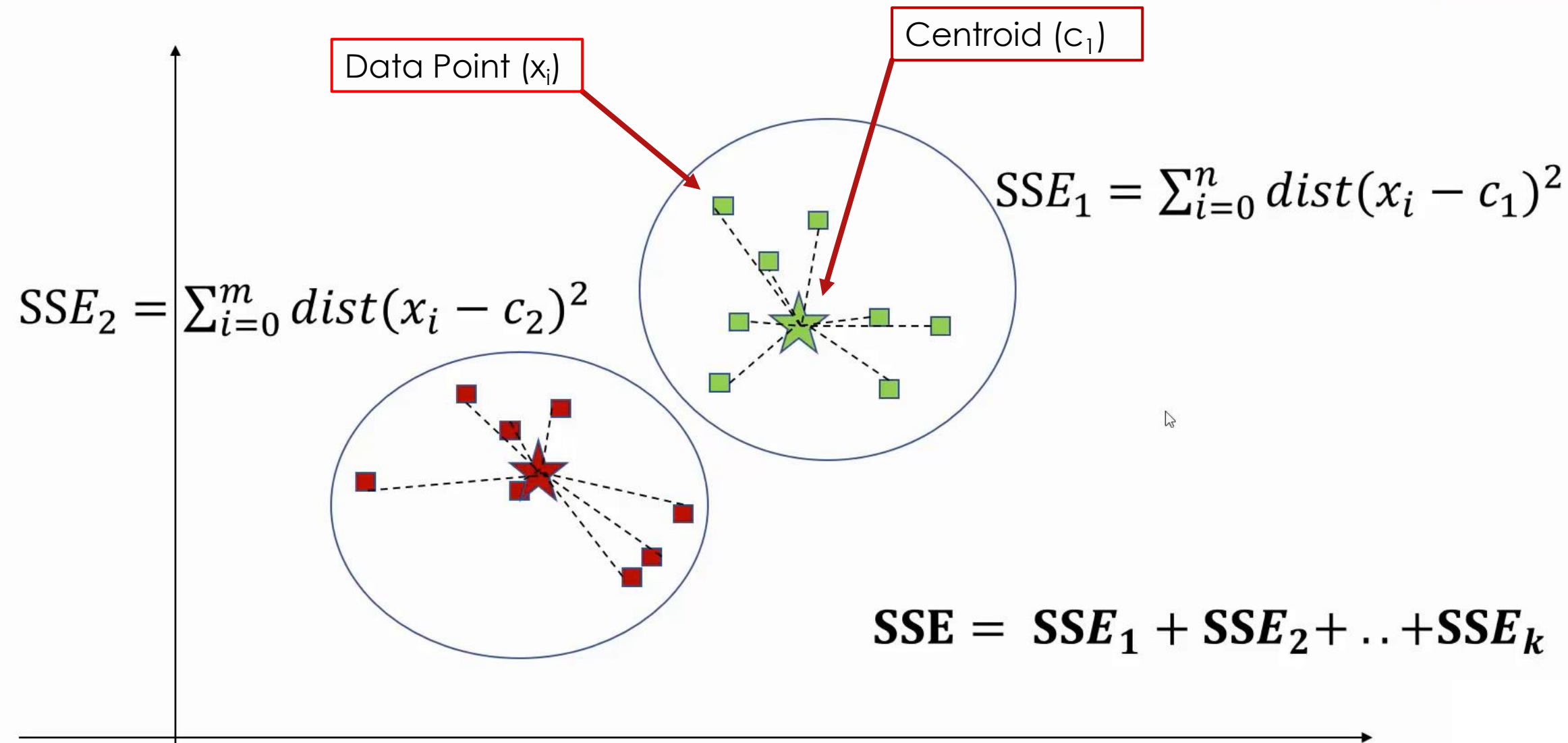
Elbow Method

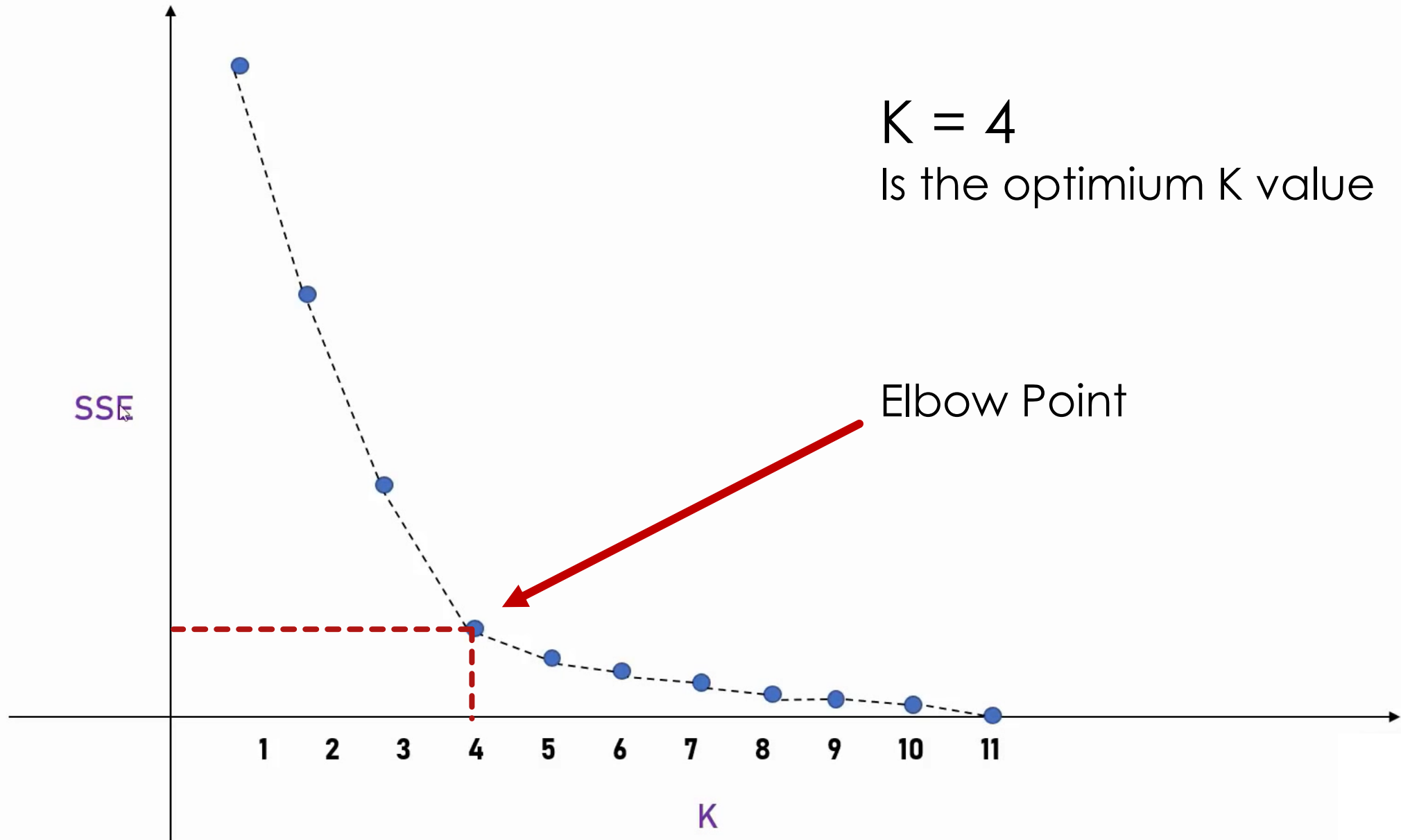


Purpose Based



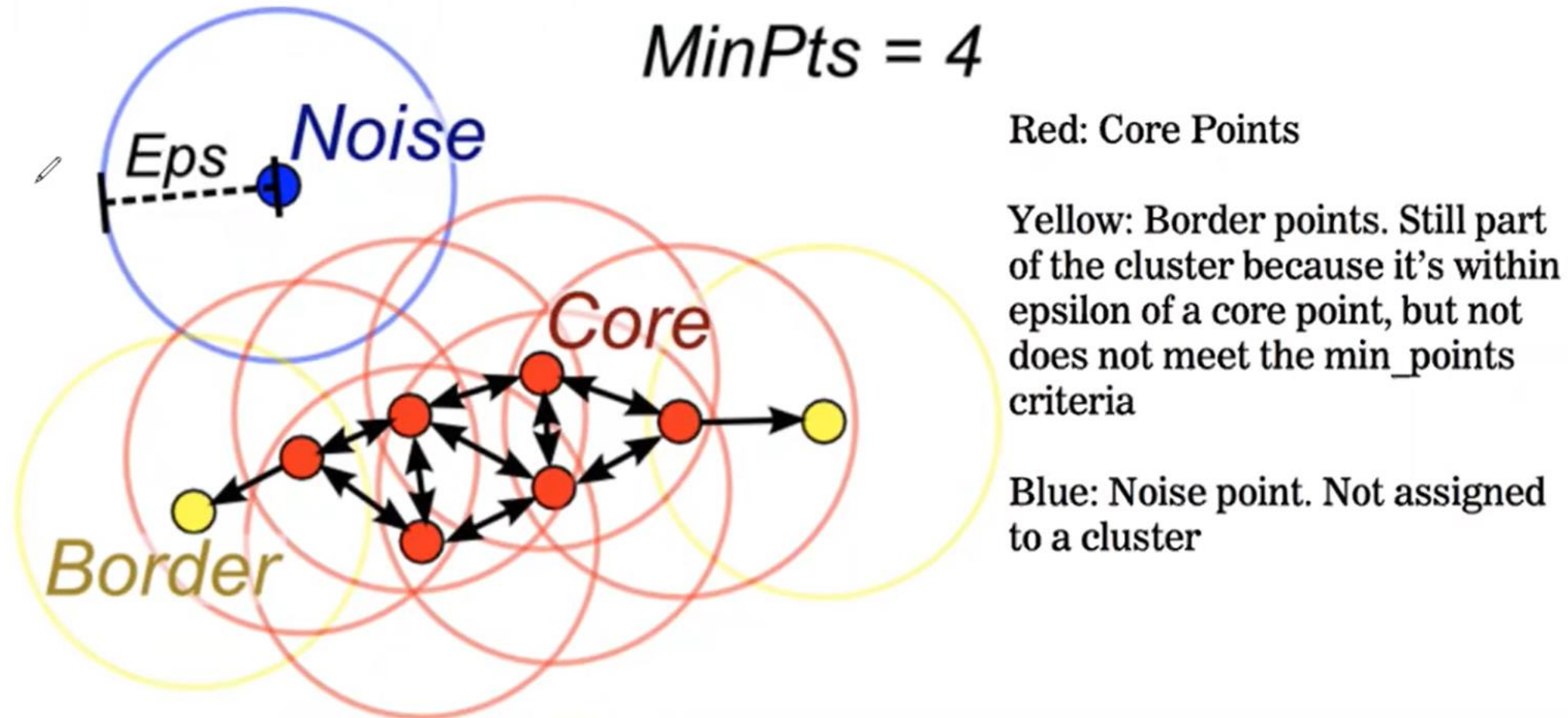
SSE = Sum of Squared Errors



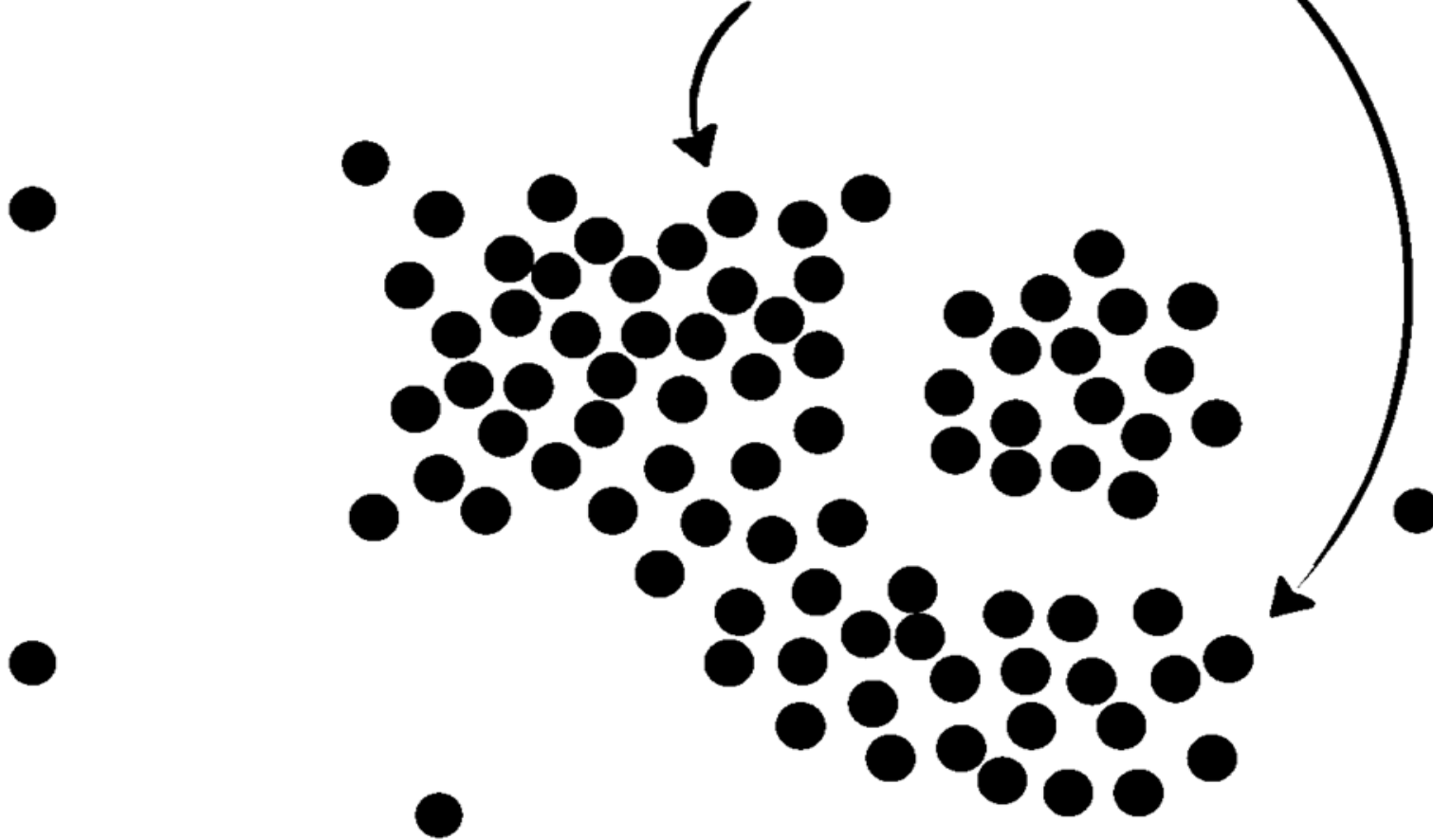




# Density-Based Spatial Clustering of Applications with Noise (DBSCAN)

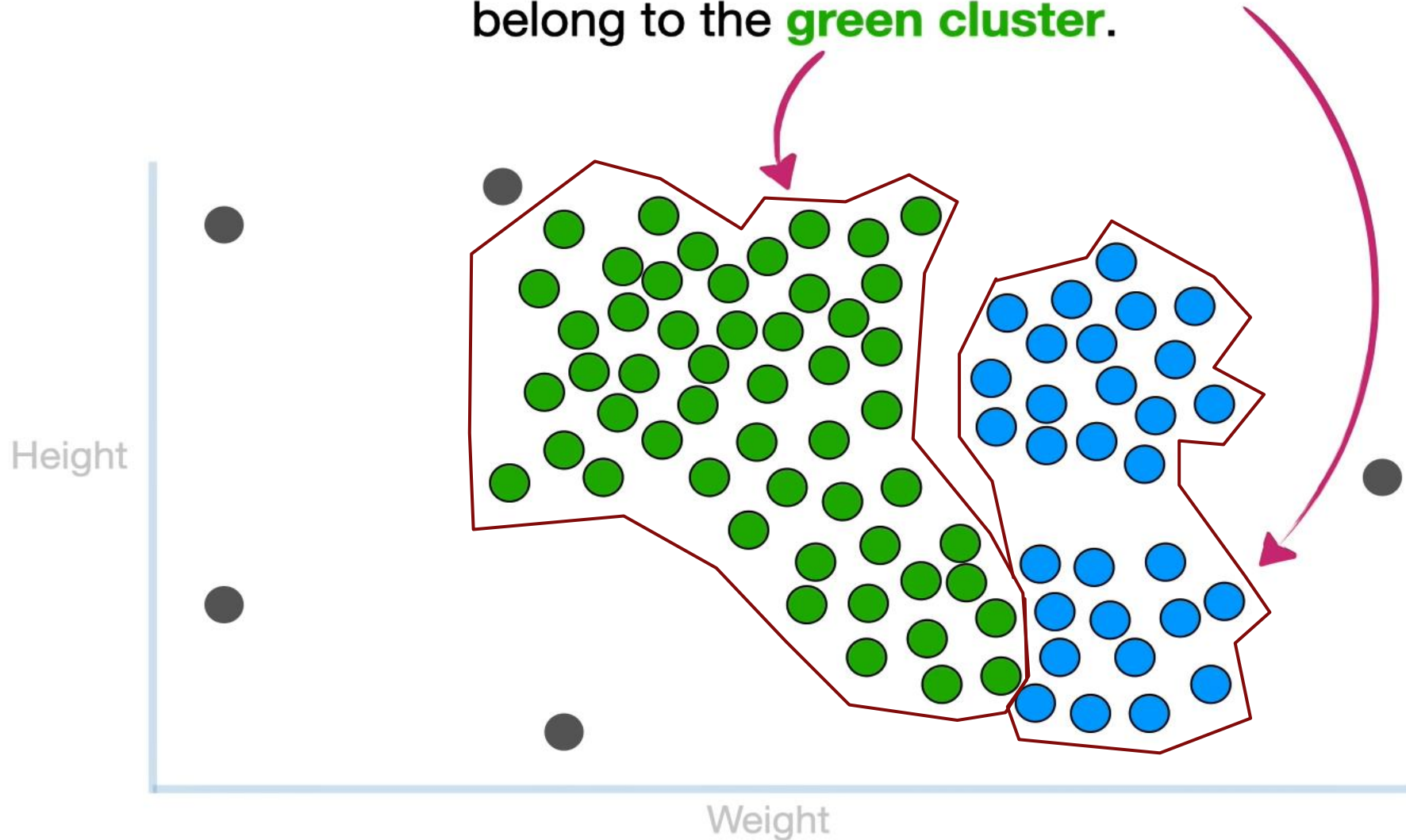


...where these points are assigned to the **blue cluster** even though they look like they belong to the **green cluster**.



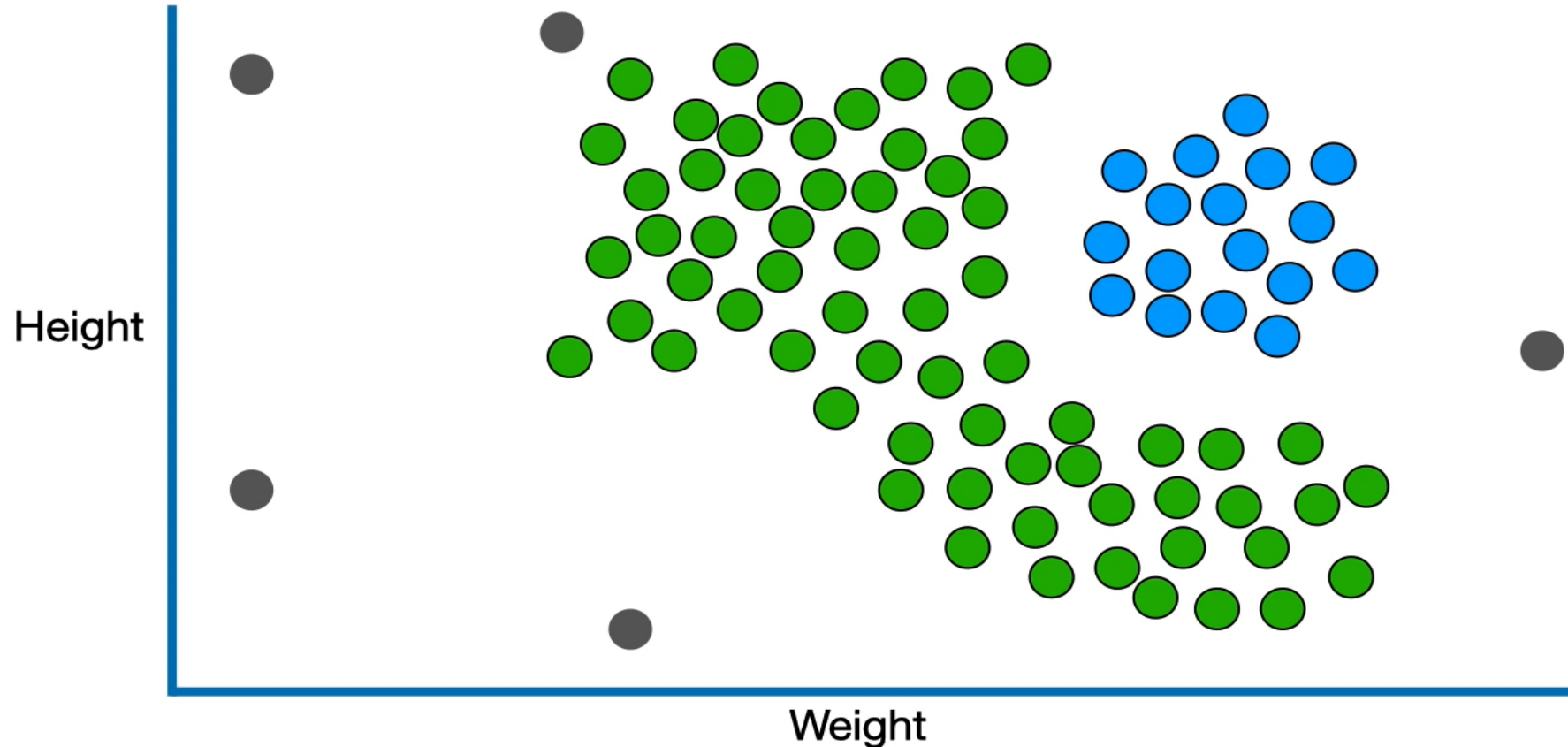
# Clustering using K – means Clustering

...where these points are assigned to the **blue cluster** even though they look like they belong to the **green cluster**.



# Clustering using DBSCAN Clustering

So let's go back to the original **2**-dimensional graph and see how **DBSCAN** tries to mimic what we can easily do by eye.



# Conclusion

- It has a capacity of analyzing the needs of the Customer.
- It has a target of reaching the Products & services for the particular group of Customers.
- For Customer Segmentation we are using DBSCAN over K-Means Algorithm as DBSCAN is more efficient & K-means is a Show stopper.

# References

- <https://technologyadvice.com/blog/marketing/customer-segmentation-methods/>
- <https://towardsdatascience.com/sentiment-analysis-concept-analysis-and-applications-6c94d6f58c17#:~:text=Sentiment%20analysis%20is%20contextual%20mining,service%20while%20monitoring%20online%20conversations.>
- <https://towardsdatascience.com/exploratory-data-analysis-8fc1cb2ofd15>
- [https://thesai.org/Downloads/Volume10No2/Paper\\_48-A\\_Study\\_on\\_Sentiment\\_Analysis\\_Techniques.pdf](https://thesai.org/Downloads/Volume10No2/Paper_48-A_Study_on_Sentiment_Analysis_Techniques.pdf)
- [https://www.researchgate.net/publication/313737530\\_Review\\_on\\_Customer\\_Segmentation\\_Technique\\_on\\_Ecommerce#:~:text=This%20paper%20will%20review%20customer%20segmentation%20using%20data%2C,and%20survey%20data%20were%20as%20the%20external%20data.](https://www.researchgate.net/publication/313737530_Review_on_Customer_Segmentation_Technique_on_Ecommerce#:~:text=This%20paper%20will%20review%20customer%20segmentation%20using%20data%2C,and%20survey%20data%20were%20as%20the%20external%20data.)





Thank You!

# Technologies

- Python Libraries –

- matplotlib.
- seaborn.
- plotly
- sklearn.
- Pandas.
- Numpy



- Machine Learning Algorithms –

- Customer Segmentation :

- DBSCAN ( Density Based Spatial Clustering of Application with Noise )
- K-means Clustering

- Sentiment Analysis :

- Natural Language Processing ( NLP )
- Using different classifiers
  - SVM (Support Vector Machine)
  - Random Forest Classifier
  - Multinomial Naive Bayes
  - Naive Bayes
  - Combination of SVM and Multinomial Naive Bayes



## Work Flow



Customer  
Data



Data pre processing

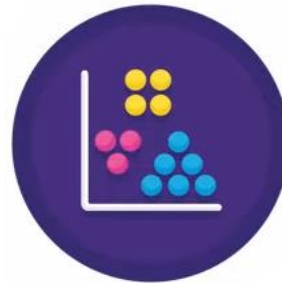


Data Analysis

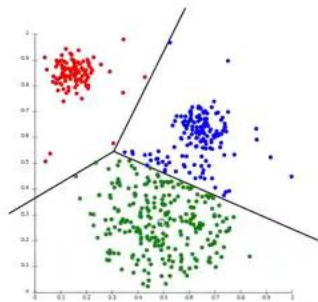


$$WCSS = \sum_{i \in n} (X_i - Y_i)^2$$

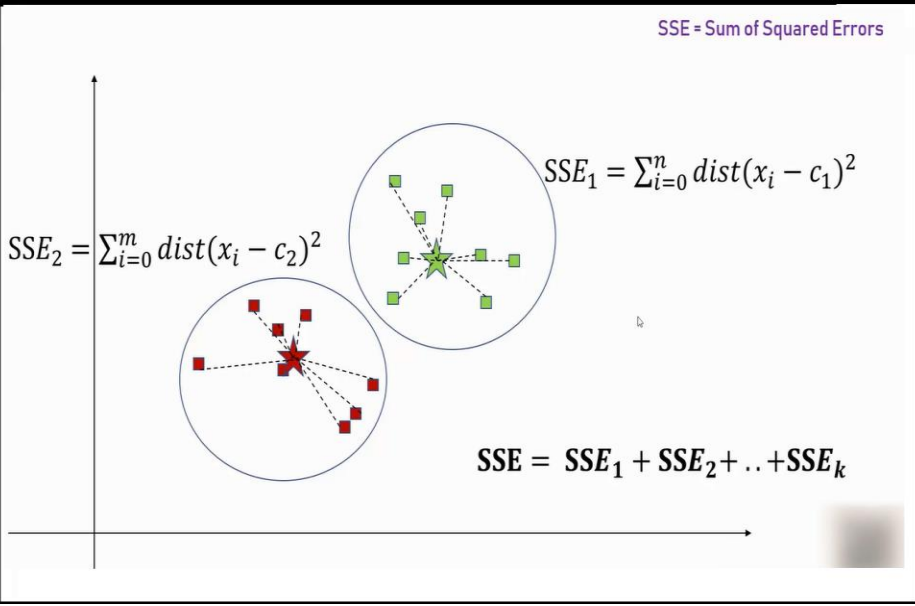
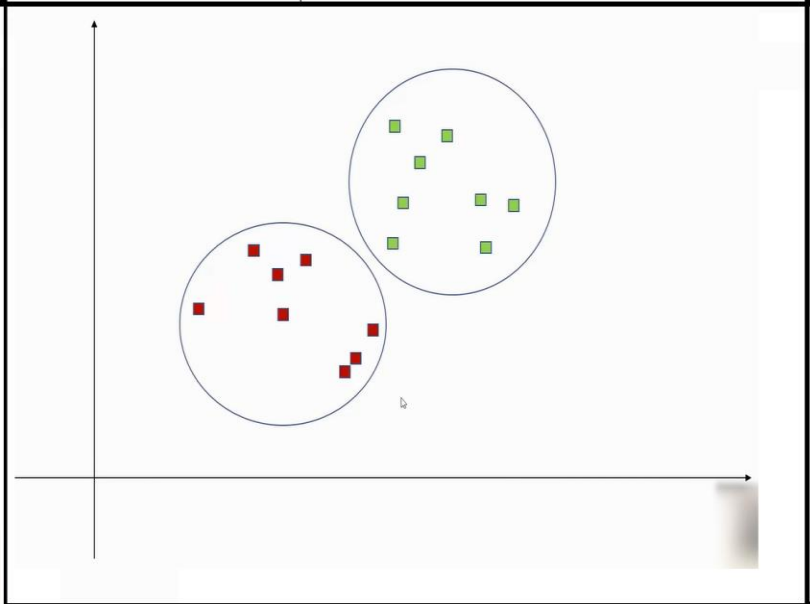
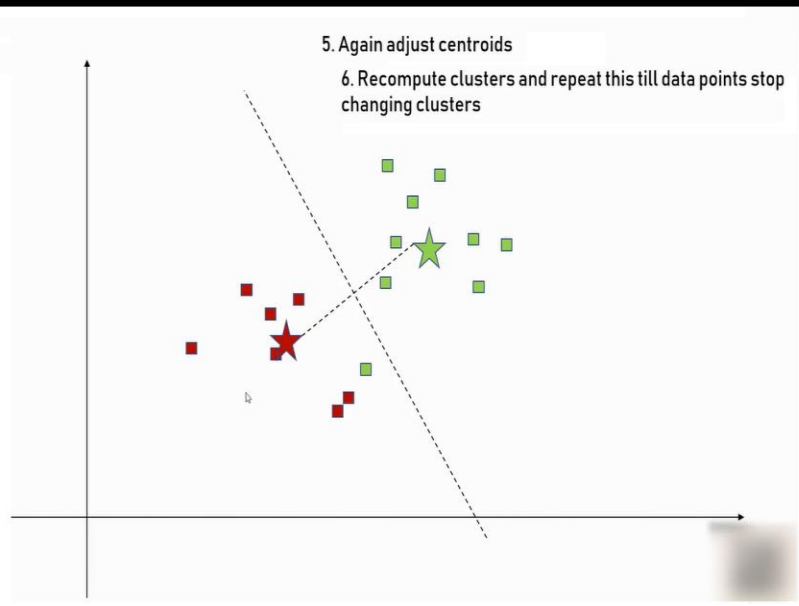
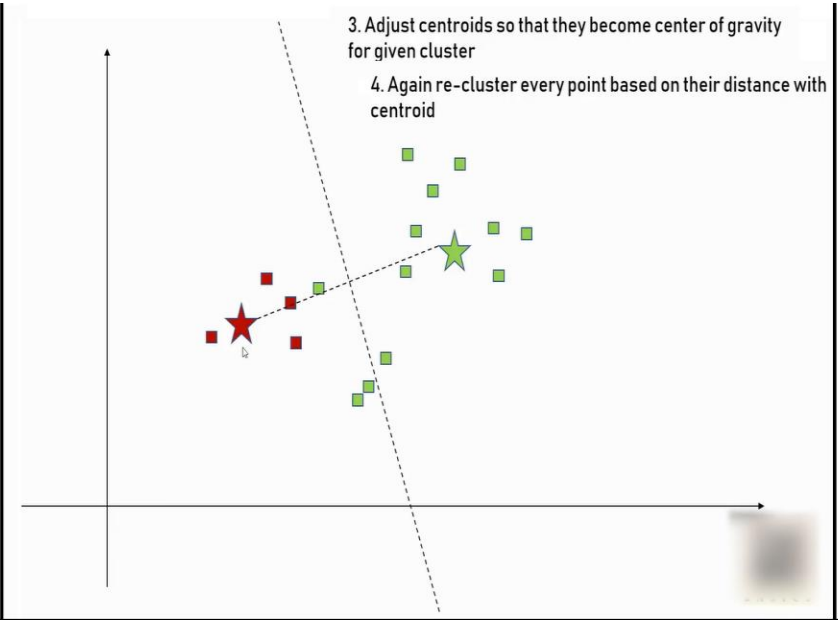
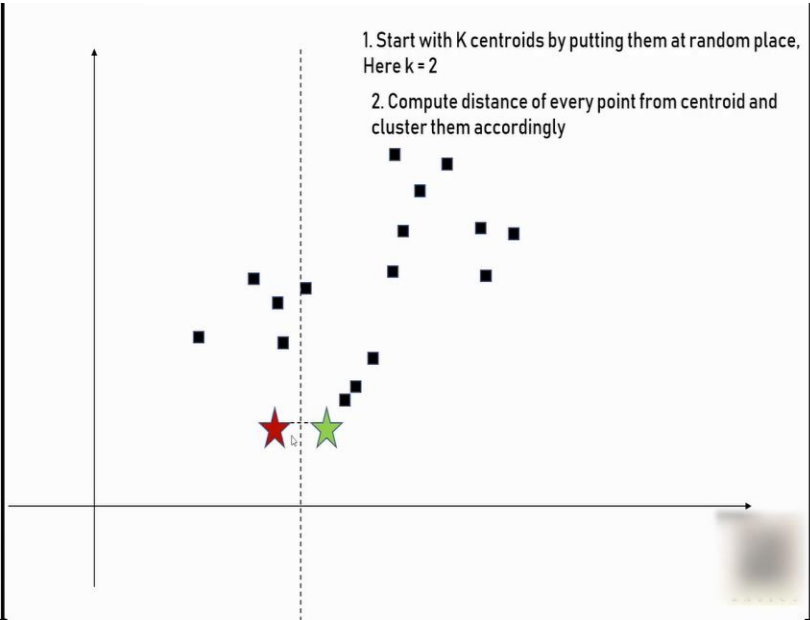
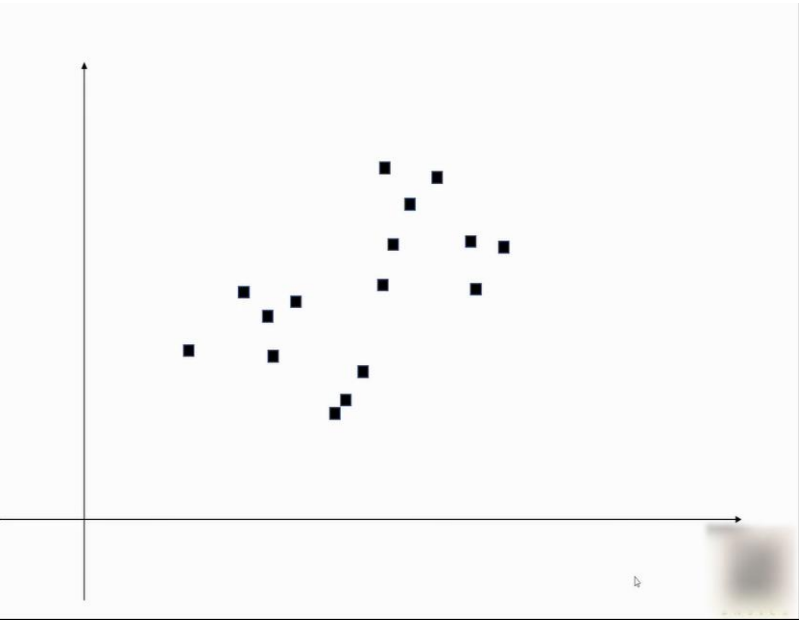
Optimum number of Clusters

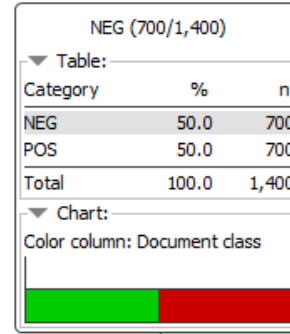


K-Means Clustering



Visualizing the Clusters



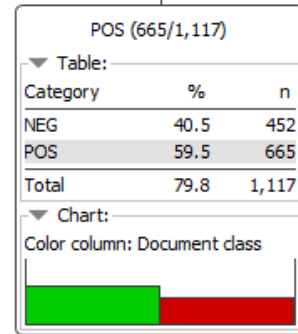


*baa*



$\leq 0.5$

$> 0.5$

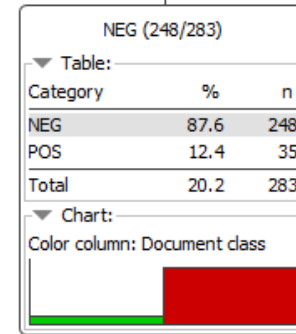
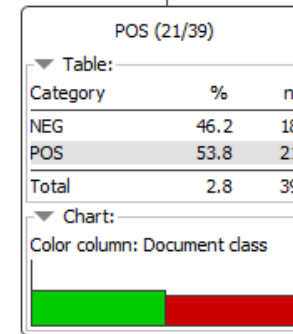
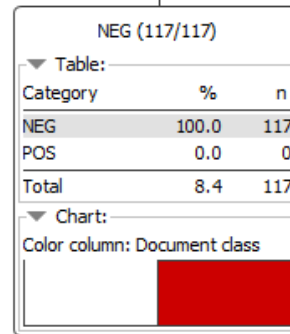
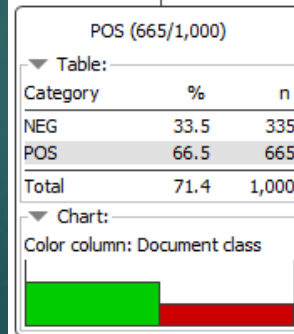


*wast*



$\leq 0.5$

$> 0.5$



*film*



$\leq 0.5$

$> 0.5$

