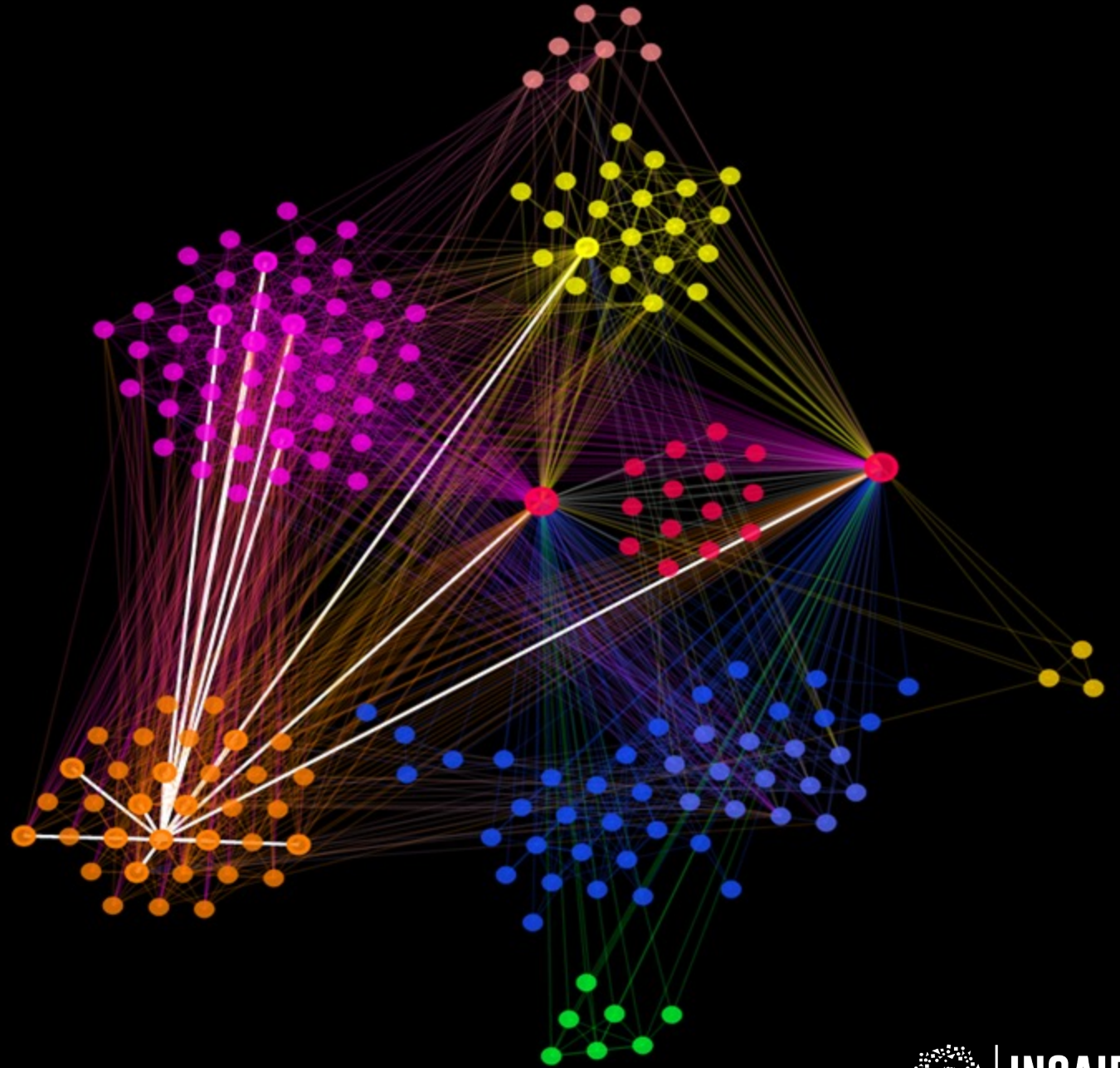


K – Means Clustering



Course Overview

Term	CDF	GCD	GCDAI	PGPDSAI
Term 1	Data Analytics with Python	Data Analytics with Python	Data Analytics with Python	Data Analytics with Python
Term 2	Data Visualization Techniques	Data Visualization Techniques	Data Visualization Techniques	Data Visualization Techniques
Term 3	EDA & Data Storytelling	EDA & Data Storytelling	EDA & Data Storytelling	EDA & Data Storytelling
		Minor Project	Minor Project	Minor Project
Term 4		Machine Learning Foundation	Machine Learning Foundation	Machine Learning Foundation
Term 5		Machine Learning Intermediate	Machine Learning Intermediate	Machine Learning Intermediate
Term 6		Machine Learning Advanced (Mandatory)	Machine Learning Advanced (Mandatory)	Machine Learning Advanced (Mandatory)
		Data Visualization with Tableau (Elective - I)	Data Visualization with Tableau (Elective - I)	Data Visualization with Tableau (Elective - I)
		Data Analytics with R (Elective - II)	Data Analytics with R (Elective - II)	Data Analytics with R (Elective - II)
		Capstone Project	Capstone Project	Capstone Project
Term 7		Bonus: Industrial ML (ML – 4 & 5)	Basics of AI, TensorFlow, and Keras	Basics of AI, TensorFlow, and Keras
Term 8			Deep Learning Foundation	Deep Learning Foundation
Term 9			NPL – I/CV – I	CV – I
Term 10			NLP – II/CV – II	NLP – I
			Capstone Project	Capstone Project
Term 11				CV – II
Term 12				NLP – II
				NLP – III + CV – III
				AutoVision & AutoNLP
				Building AI product

You are here...

Term Context

- K – Nearest Neighbor
- **K-means Clustering**
- Ensemble Learning
- Optimization

← You are here...

Agenda

1. What is Clustering?

2. K-Means Clustering

3. When to use K-Means Clustering?

4. What is K?

5. Euclidean Distance

6. Inertia

7. Stopping Criteria

8. K-Means Clustering Example

9. Elbow Method

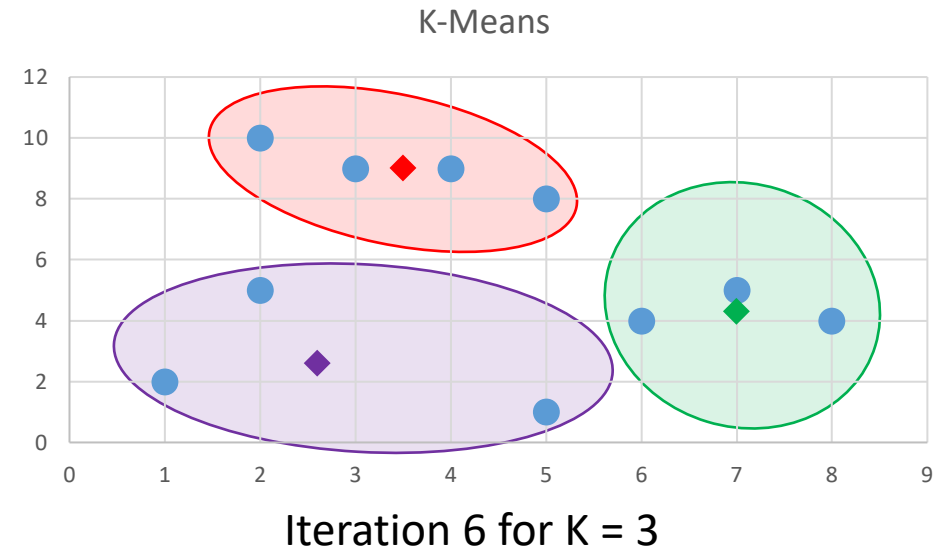
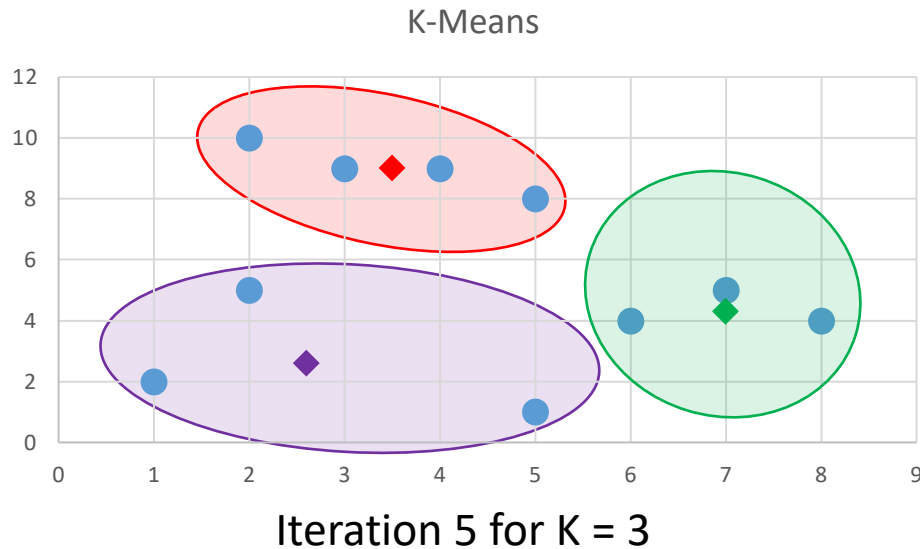
10. Advantages

11. Limitations

12. Applications

Stopping Criteria

- There are essentially three stopping criteria that can be adopted to stop the K-means algorithm:
 - Centroids of newly formed clusters do not change.
 - Points remain in the same cluster after many iterations.
 - Maximum number of iterations are reached.



Agenda

1. What is Clustering?

2. K-Means Clustering

3. When to use K-Means Clustering?

4. What is K?

5. Euclidean Distance

6. Inertia

7. Stopping Criteria

8. K-Means Clustering Example

9. Elbow Method

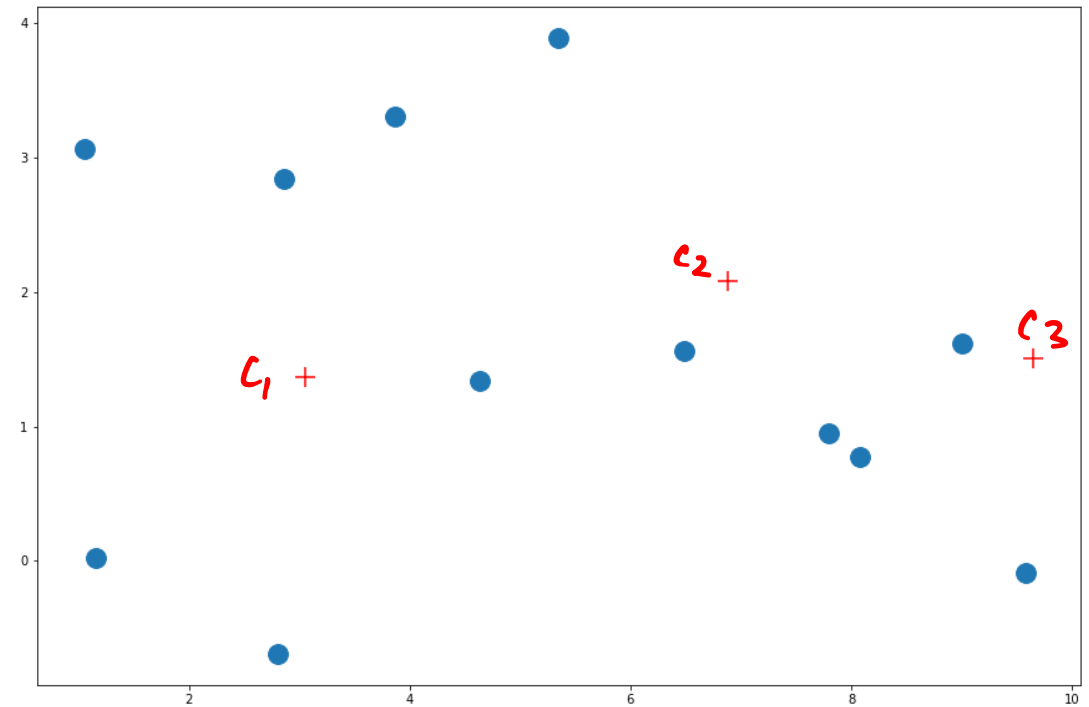
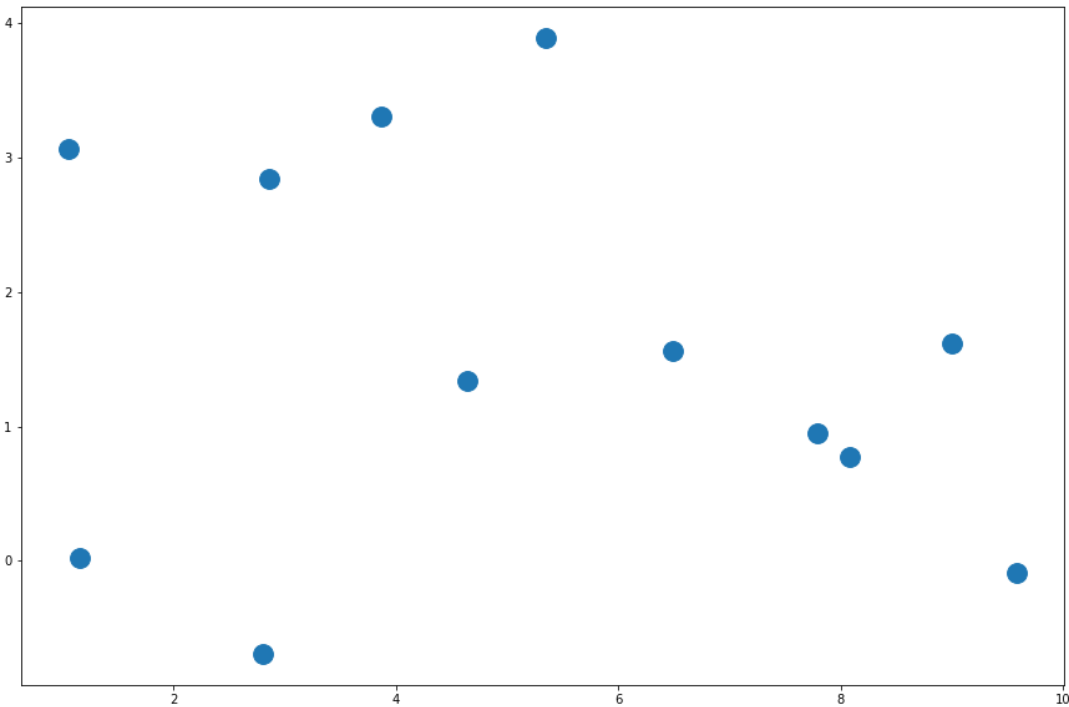
10. Advantages

11. Limitations

12. Applications

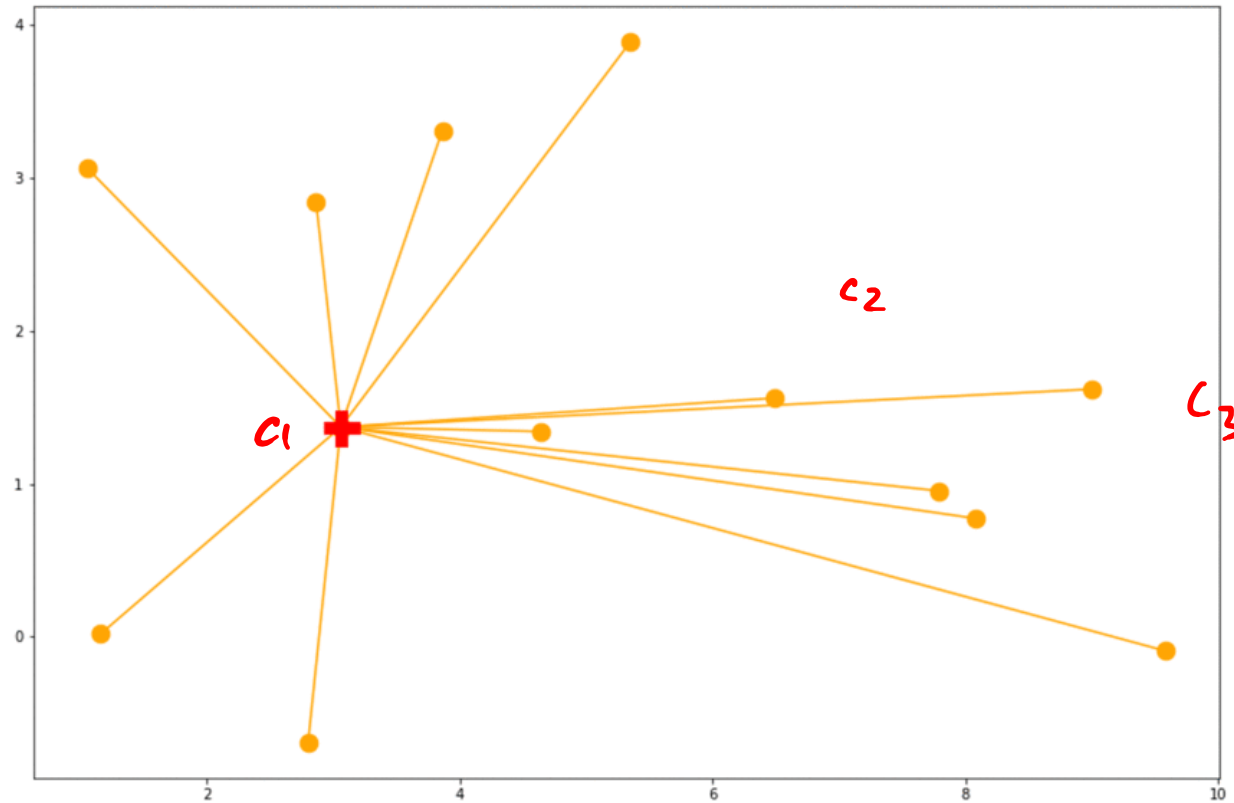
K-Means Clustering (Example) – Iteration 1

- Let's apply K-Means clustering with $K=3$ to the following set of points. We will limit the number of iterations to 5.
- Step 1: Assume $K=3$ random points anywhere in the graph as cluster representatives.



K-Means Clustering (Example) – Iteration 1

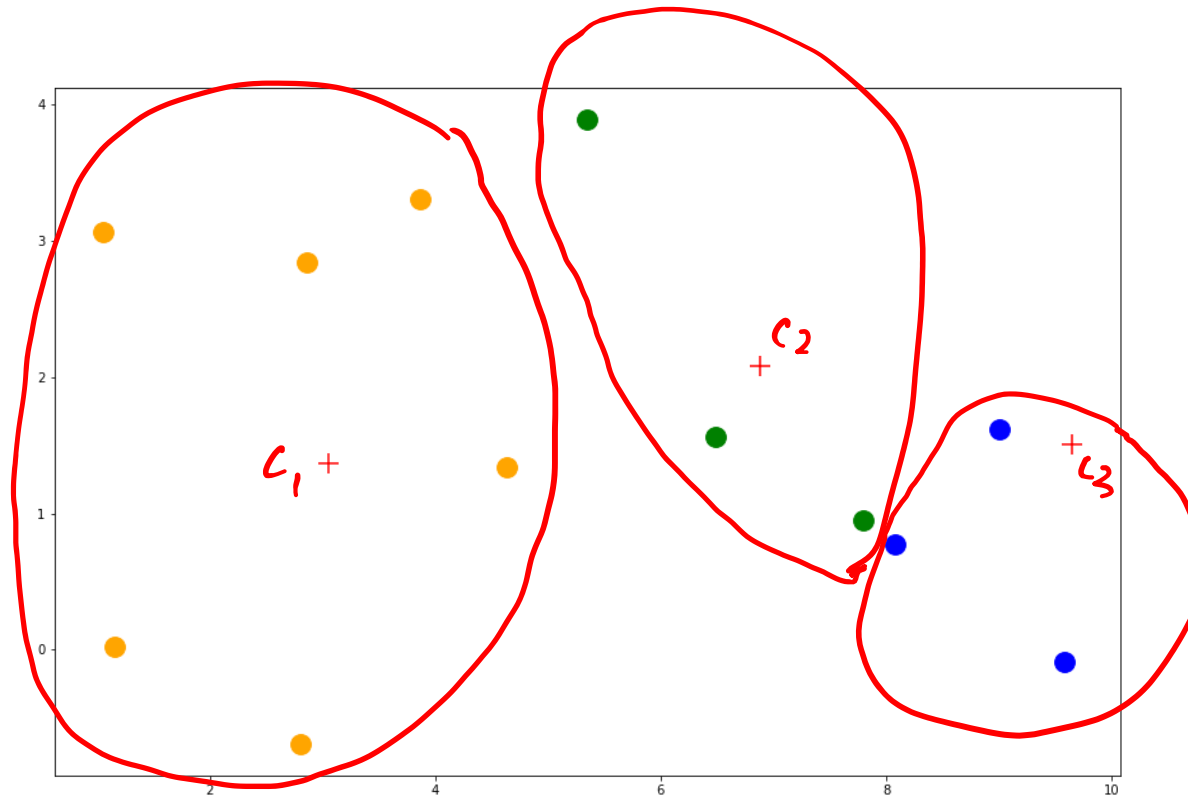
- Step 2: Calculate the distance of each cluster representatives with every point in the data



Distance of cluster representative 1 with all the data points

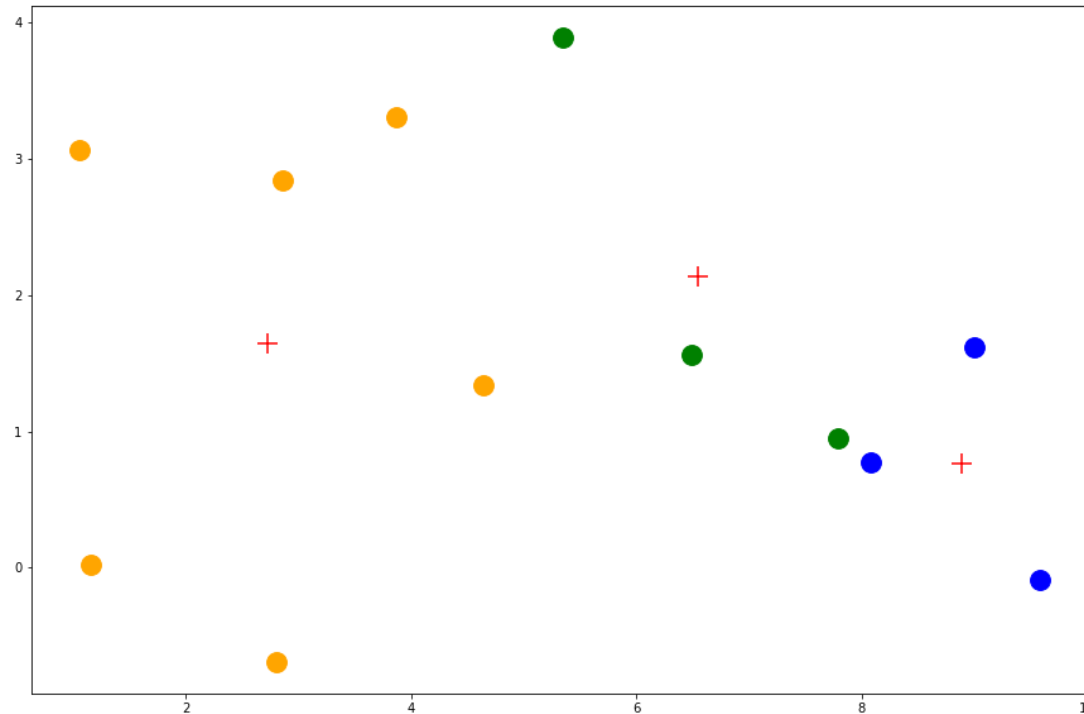
K-Means Clustering (Example) – Iteration 1

- Step 3: Assign the nearest data points to the chosen cluster representatives.



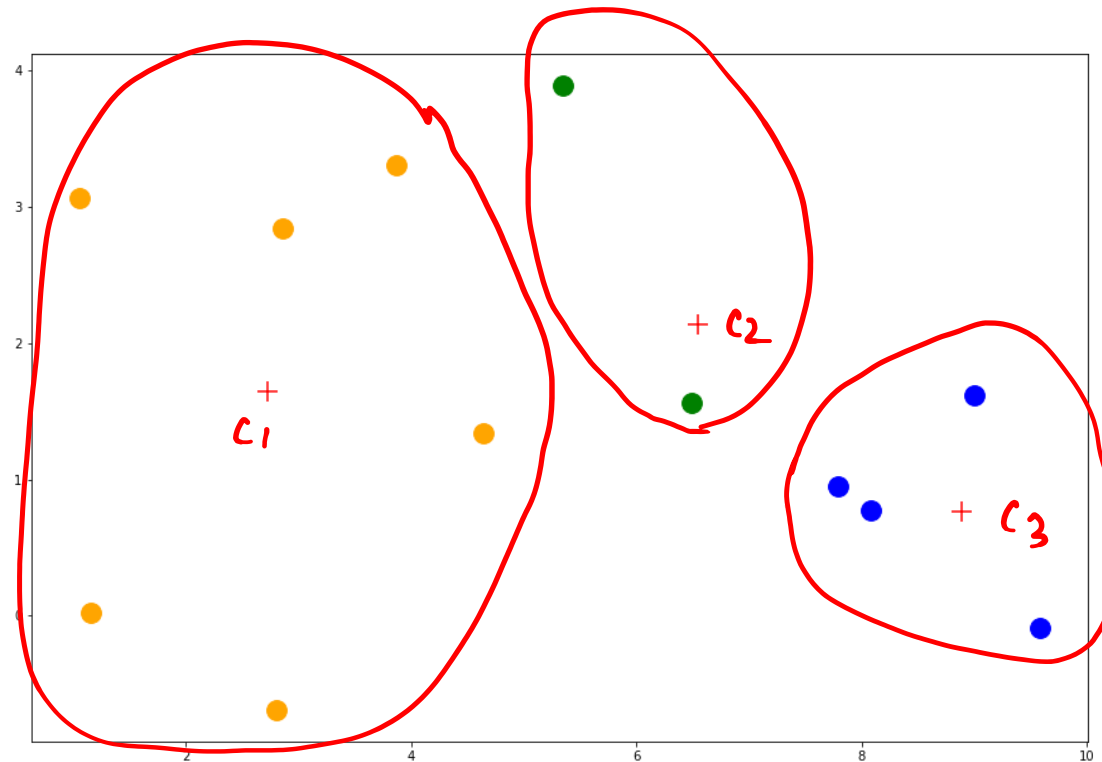
K-Means Clustering (Example) – Iteration 1

- Step 4: Calculate the new center representatives by taking the average of the points in each cluster.
- The total inertia of these clusters is 20.65.



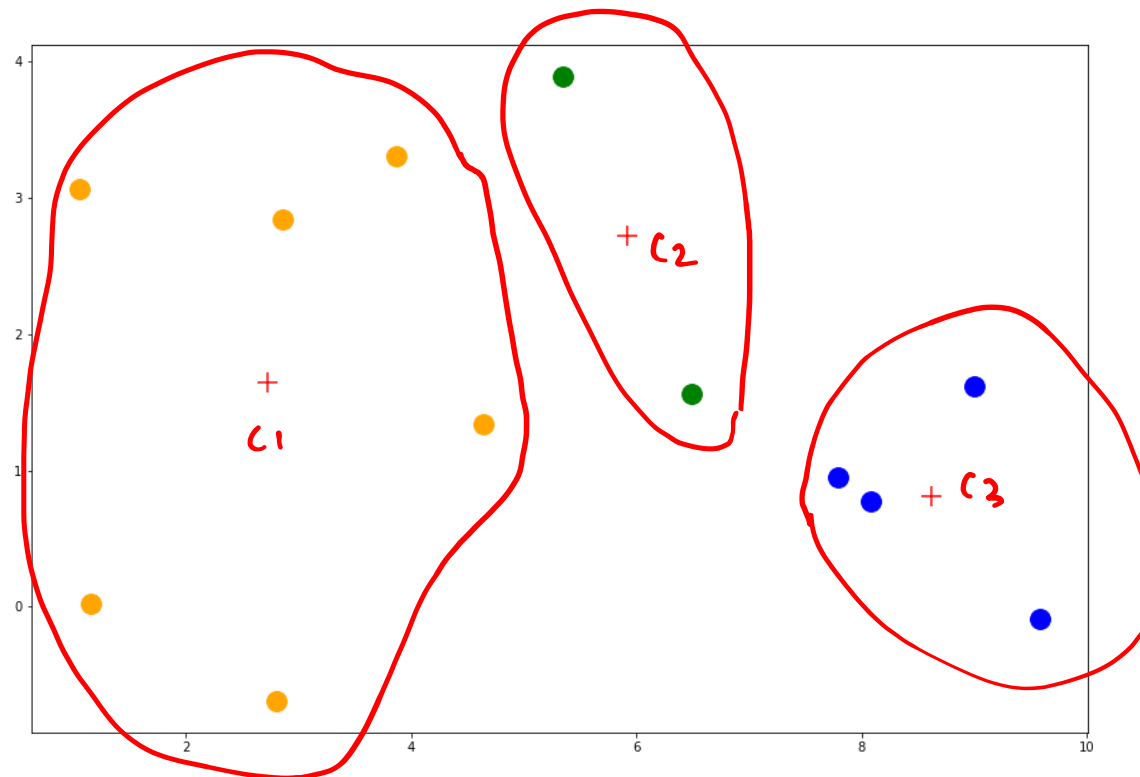
K-Means Clustering (Example) – Iteration 2

- We will repeat steps 2 and 3 with the new cluster representatives. The new clusters are:



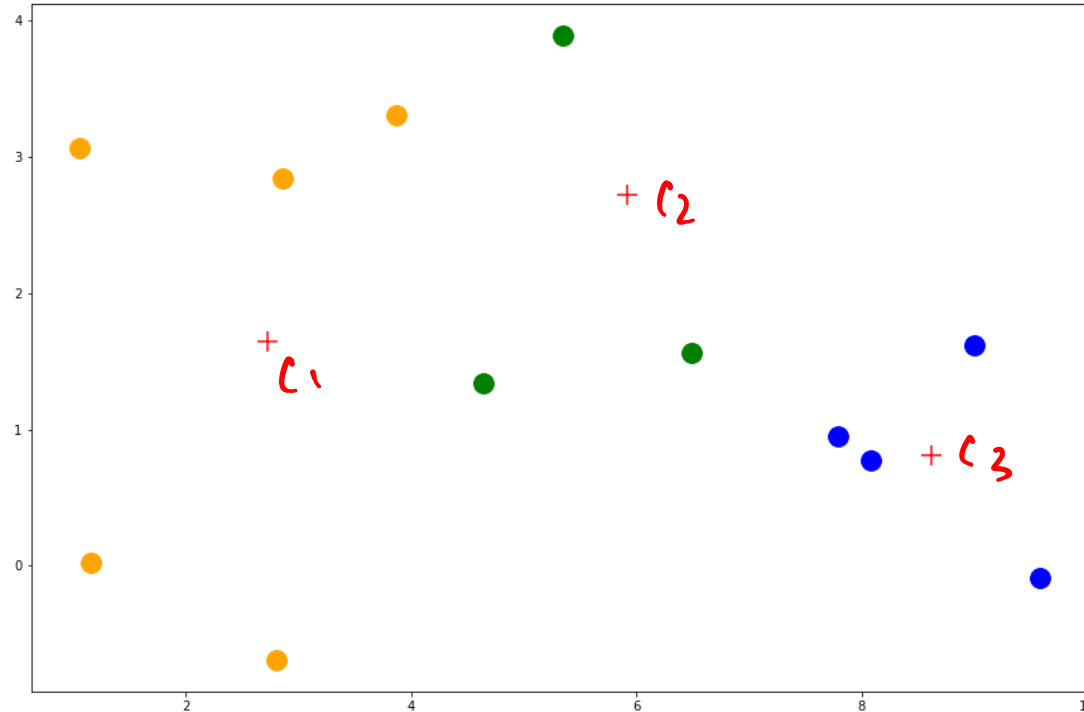
K-Means Clustering (Example) – Iteration 2

- We will re-calculate the new cluster representatives. The total **inertia** for these clusters are **18.5**.



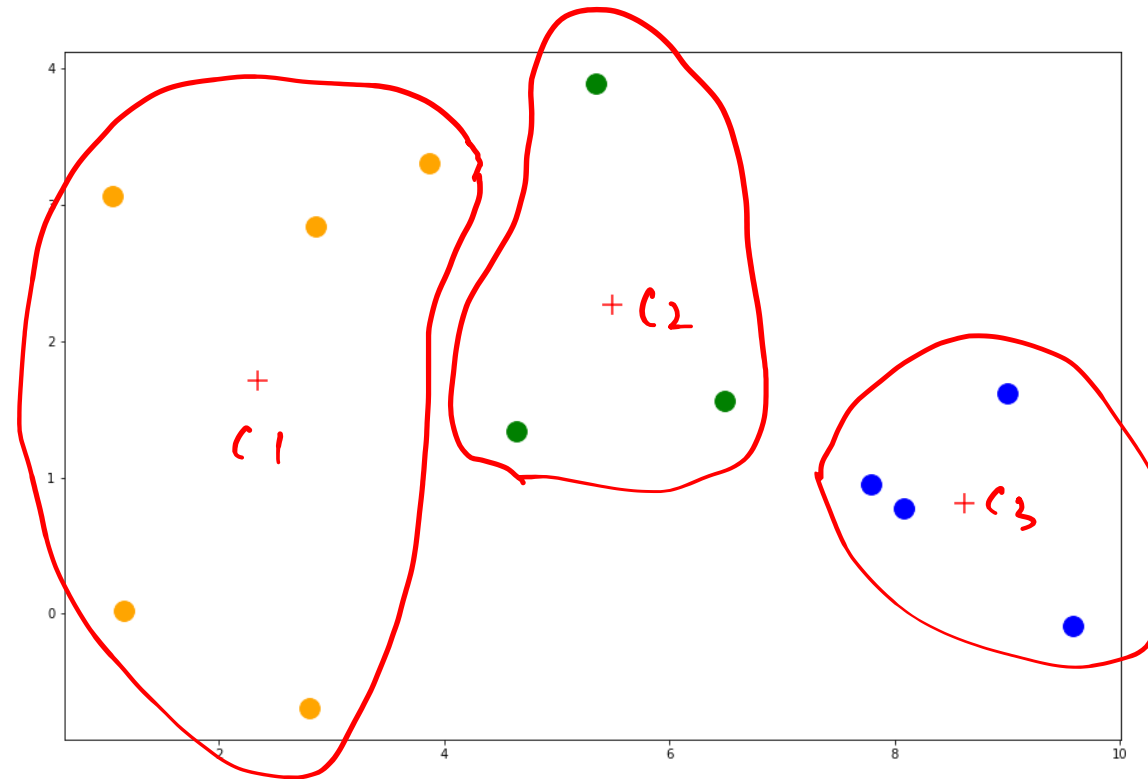
K-Means Clustering (Example) – Iteration 3

- We will repeat steps 2 and 3 with the new cluster representatives. The new clusters are:



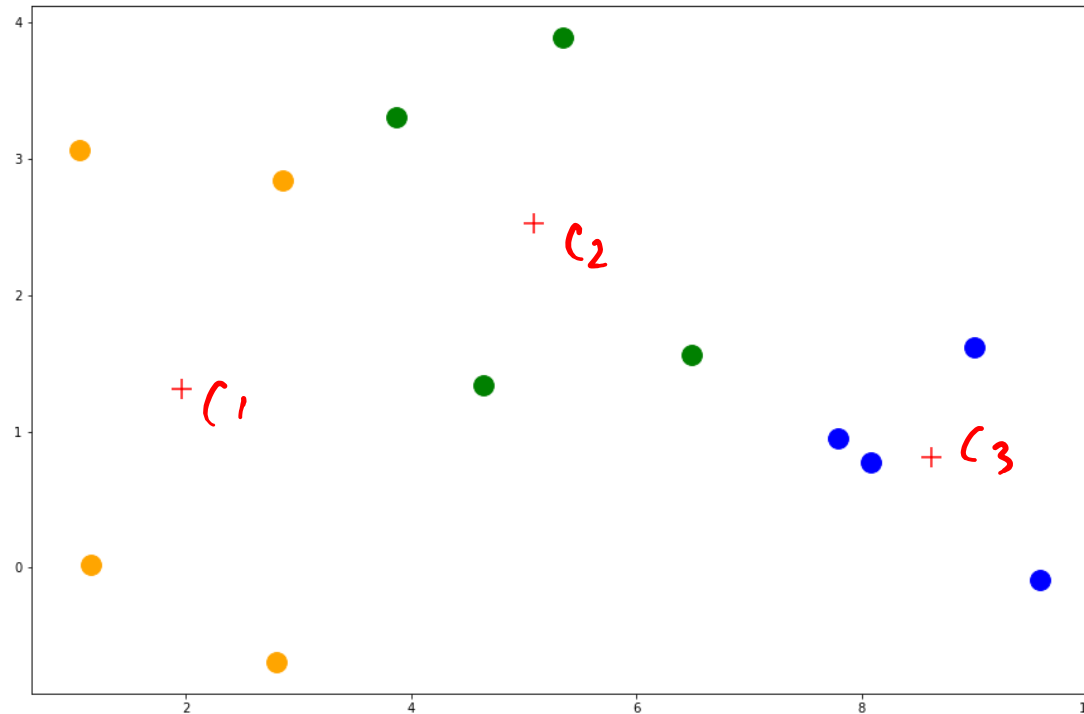
K-Means Clustering (Example) – Iteration 3

- We will re-calculate the new cluster representatives. The total **inertia** for these clusters are 18.07.



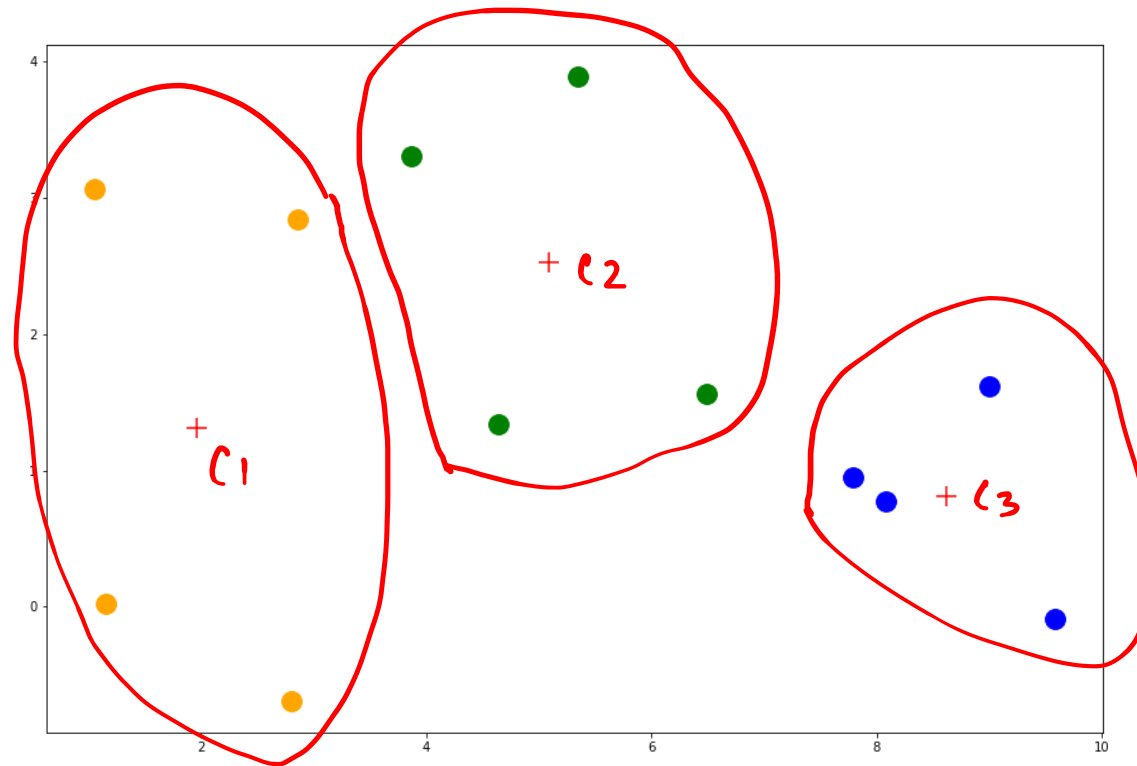
K-Means Clustering (Example) – Iteration 4

- We will repeat steps 2 and 3 with the new cluster representatives. The new clusters are:



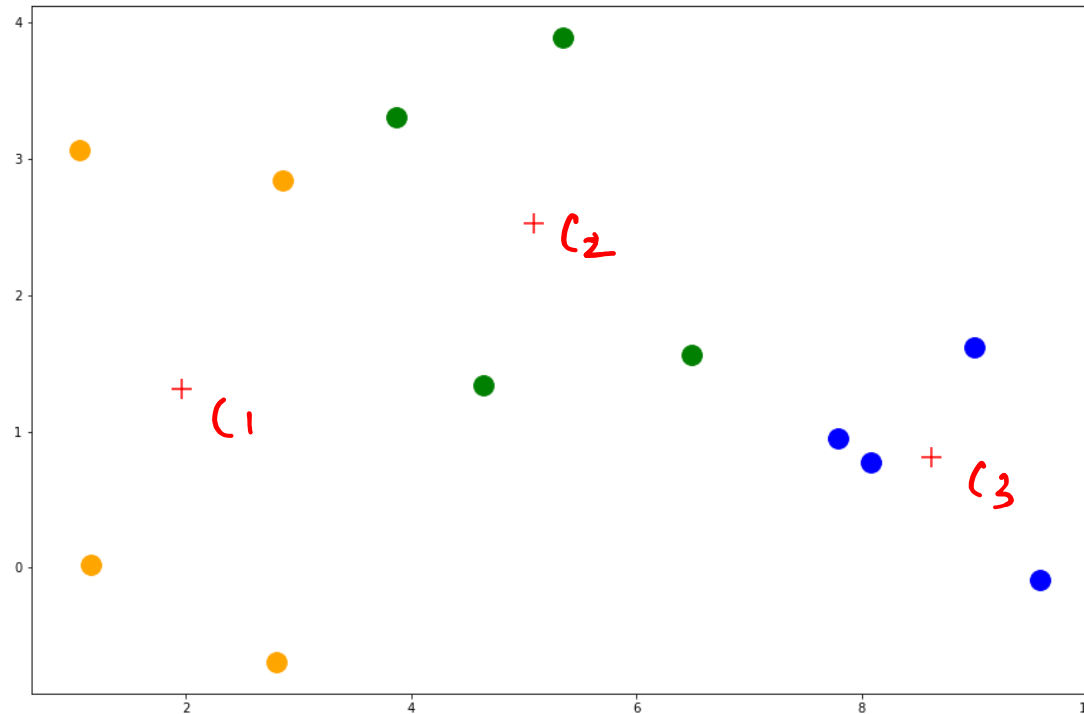
K-Means Clustering (Example) – Iteration 4

- We will re-calculate the new cluster representatives. The total **inertia** for these clusters are 17.25.



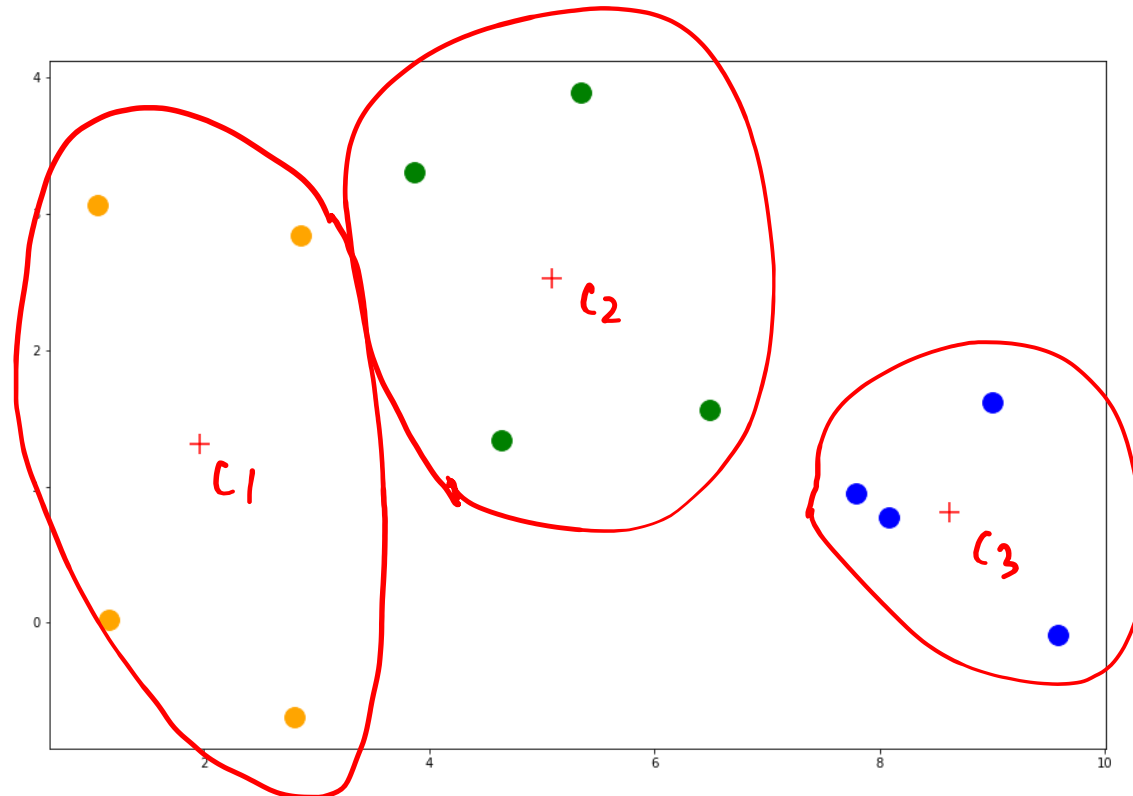
K-Means Clustering (Example) – Iteration 5

- We will repeat steps 2 and 3 with the new cluster representatives. The new clusters are:



K-Means Clustering (Example) – Iteration 5

- We will re-calculate the new cluster representatives. The total **inertia** for these clusters are 16.83.
- We have reached our estimated number of iterations and the cluster **representatives** remain the **same**.
- We will stop the iterations here. These are our **final clusters**:



Agenda

1. What is Clustering?

2. K-Means Clustering

3. When to use K-Means Clustering?

4. What is K?

5. Euclidean Distance

6. Inertia

7. Stopping Criteria

8. K-Means Clustering Example

9. Elbow Method

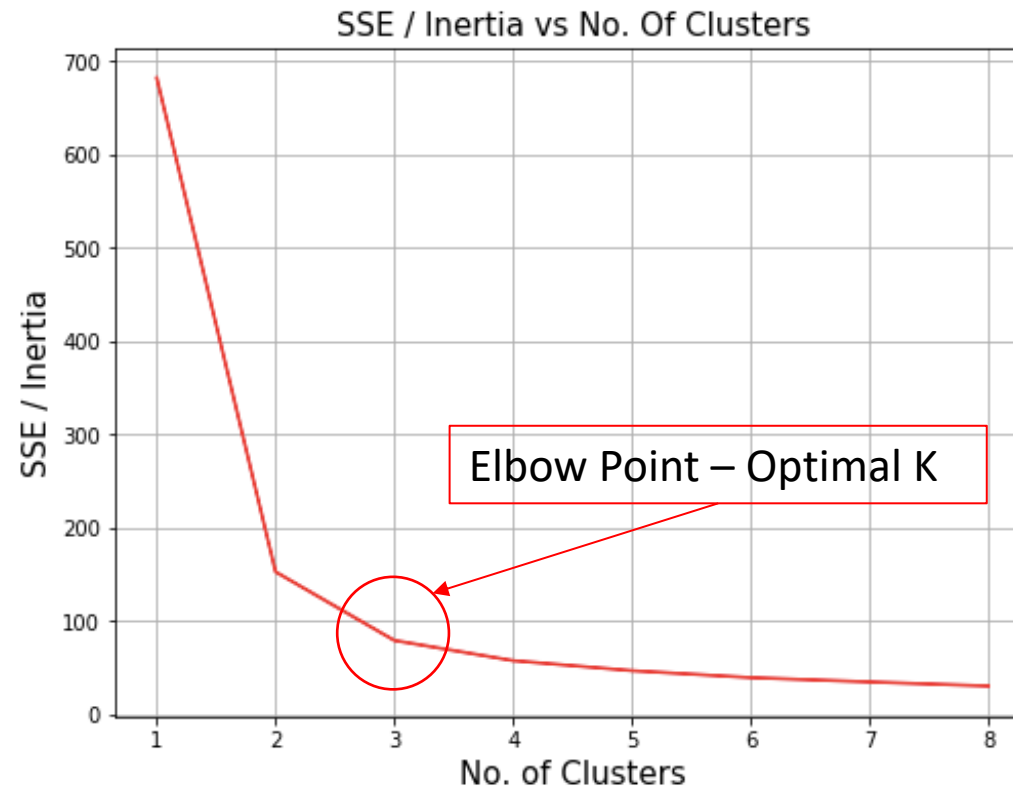
10. Advantages

11. Limitations

12. Applications

Elbow Method – Optimal Value

- It is one of the most popular methods to determine the optimal value of K.
- We use it to choose a K when we observe negligible change in the inertial values between different values of K.



Agenda

1. What is Clustering?

2. K-Means Clustering

3. When to use K-Means Clustering?

4. What is K?

5. Euclidean Distance

6. Inertia

7. Stopping Criteria

8. K-Means Clustering Example

9. Elbow Method

10. Advantages

11. Limitations

12. Applications

Advantages



- Relatively simple to implement.
- Scales to large data sets.
- Guarantees convergence.
- Can warm-start the positions of centroids.
- Easily adapts to new examples.
- Generalizes to clusters of different shapes and sizes, such as elliptical clusters.

Agenda

1. What is Clustering?

2. K-Means Clustering

3. When to use K-Means Clustering?

4. What is K?

5. Euclidean Distance

6. Inertia

7. Stopping Criteria

8. K-Means Clustering Example

9. Elbow Method

10. Advantages

11. Limitations

12. Applications

Limitations



- Choosing k manually. (Unsupervised ML)
- Dependent on initial values.
- Prone to varying sizes and density of clusters.
- Prone to outliers.
- Prone to Curse of Dimensionality. (huge no of cols.)
- Convergence to a local minimum may produce wrong results.

Agenda

1. What is Clustering?

2. K-Means Clustering

3. When to use K-Means Clustering?

4. What is K?

5. Euclidean Distance

6. Inertia

7. Stopping Criteria

8. K-Means Clustering Example

9. Elbow Method

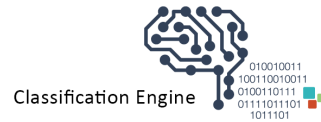
10. Advantages

11. Limitations

12. Applications

Applications – Document Classification

- By classifying text, we are aiming to assign **multiple categories** to a **document**, making it easier to manage and sort.
- This is especially useful for publishers, news sites, blogs or anyone who deals with a lot of content.



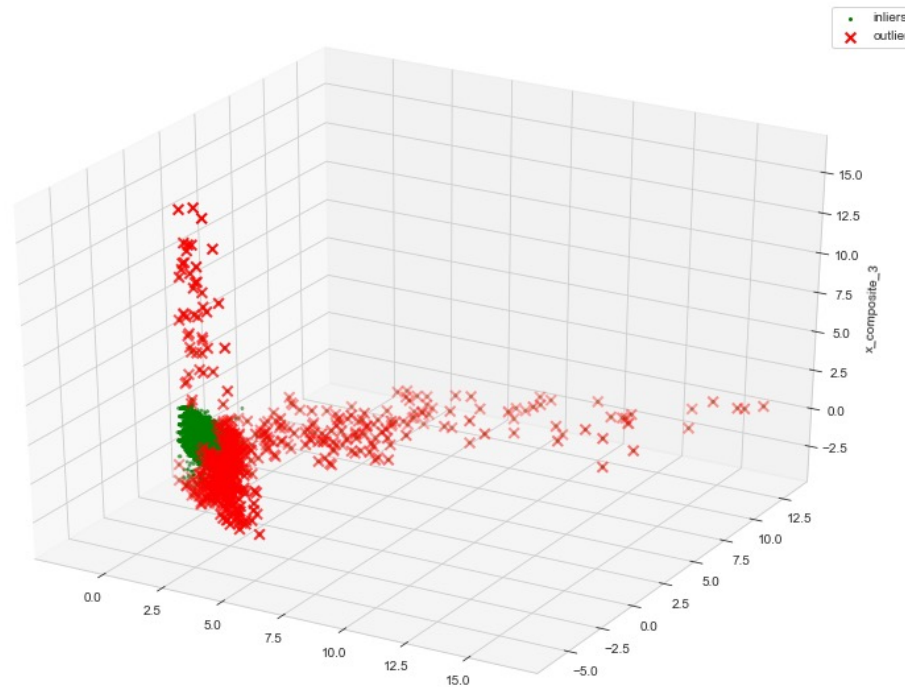
Applications – Image Compression

- Image compression is reducing the size that an image takes while storing or transmitting.
- While compressing, the colors are clustered towards some major colors by grouping them towards the major colors.



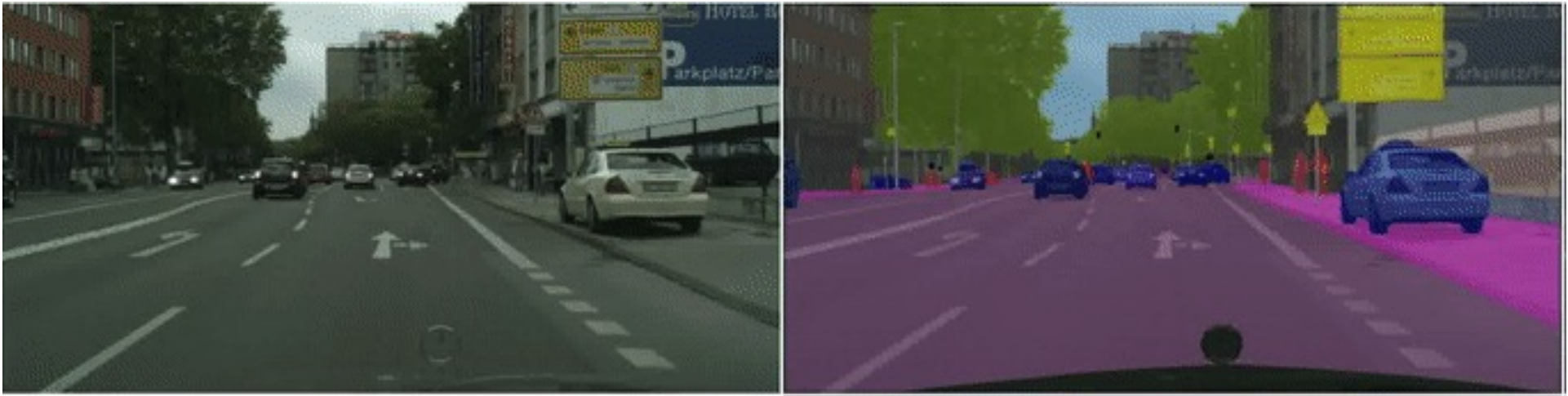
Applications – Outlier Detection

- In the k-means based outlier detection technique the **data points** are partitioned into k groups by assigning them to the **closest cluster centers**.



Applications – Image Segmentation

- Image Segmentation is the process of assigning a **label** to **every pixel** in an image such that **pixels** with the **same label** **share** certain **characteristics**.



Thank
you