

MIT WORLD PEACE UNIVERSITY

Data Science for Cybersecurity and Forensics  
Third Year B. Tech, Semester 6

---

---

IMPLEMENTATION OF K-MEANS CLUSTERING IN  
PYTHON

---

---

ASSIGNMENT 7

Prepared By

Krishnaraj Thadesar  
Cyber Security and Forensics  
Batch A1, PA 10

April 16, 2024

# Contents

<b>1</b>	<b>Aim</b>	<b>1</b>
<b>2</b>	<b>Objectives</b>	<b>1</b>
<b>3</b>	<b>Theory</b>	<b>1</b>
3.1	K-Means Clustering . . . . .	1
3.2	K-Means Clustering Algorithm . . . . .	1
<b>4</b>	<b>Procedure</b>	<b>2</b>
<b>5</b>	<b>Platform</b>	<b>2</b>
<b>6</b>	<b>Requirements</b>	<b>2</b>
<b>7</b>	<b>Code</b>	<b>2</b>
<b>8</b>	<b>FAQs</b>	<b>3</b>
<b>9</b>	<b>Conclusion</b>	<b>3</b>

## 1 Aim

Implement K-Means Clustering in Python using IOT Based Attacks dataset.

## 2 Objectives

1. To understand the concept of K-Means Clustering.
2. To implement K-Means Clustering in Python.

## 3 Theory

### 3.1 K-Means Clustering

K-Means clustering is a popular unsupervised machine learning algorithm used for partitioning a dataset into K distinct clusters. The algorithm aims to minimize the variance within each cluster and maximize the variance between clusters.

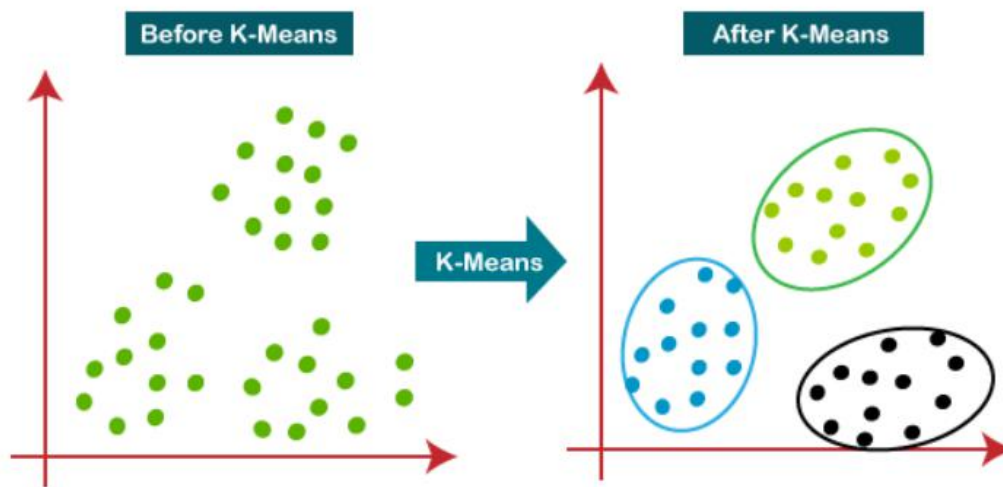


Figure 1: K Means Clustering Example

### 3.2 K-Means Clustering Algorithm

The K-Means clustering algorithm can be summarized in the following steps:

1. **Initialization:** Randomly initialize K cluster centroids.
2. **Assign Data Points to Nearest Centroid:** Assign each data point to the cluster with the nearest centroid.
3. **Update Centroids:** Recalculate the centroids of each cluster based on the mean of the data points assigned to that cluster.
4. **Repeat Steps 2 and 3:** Iterate the process of assigning data points to clusters and updating centroids until convergence.

5. **Convergence Criteria:** The algorithm converges when the centroids no longer change significantly between iterations or when a specified number of iterations is reached.

The objective function of K-Means clustering can be defined as minimizing the within-cluster sum of squared distances:

$$\operatorname{argmin}_C \sum_{i=1}^K \sum_{x \in C_i} \|x - \mu_i\|^2$$

Where: -  $C$  represents the set of clusters. -  $C_i$  represents the data points assigned to cluster  $i$ . -  $\mu_i$  represents the centroid of cluster  $i$ . -  $\|\cdot\|$  denotes the Euclidean distance.

The algorithm converges to a locally optimal solution, and the quality of the clustering depends on the initial placement of centroids and the choice of  $K$ . K-Means is efficient and scalable for large datasets but may converge to suboptimal solutions depending on the initial centroids and data distribution.

## 4 Procedure

1. Import the required python packages.
2. Load the dataset.
3. Data analysis.
4. Split the dataset into dependent/independent variables.
5. Split data into Train/Test sets.
6. Train the regression model.
7. Predict the result.

## 5 Platform

**Operating System:** Windows 11

**IDEs or Text Editors Used:** Visual Studio Code

**Compilers or Interpreters:** Python 3.10.1

## 6 Requirements

```
1 python==3.10.1
2 matplotlib==3.8.3
3 numpy==1.26.4
4 pandas==2.2.2
5 seaborn==0.13.2
```

## 7 Code

## 8 FAQs

### 1. What do you understand by K-Means method?

- K-Means is a popular clustering algorithm used in machine learning and data mining.
- It aims to partition a dataset into K distinct, non-overlapping clusters, where each data point belongs to the cluster with the nearest mean.
- The algorithm iteratively assigns each data point to the nearest cluster centroid and recalculates the centroids until convergence.
- K-Means is commonly used for data exploration, pattern recognition, and image segmentation tasks.

### 2. Discuss on how can we get data from IoT devices for Cyber Security?

- IoT (Internet of Things) devices generate vast amounts of data that can be leveraged for cybersecurity purposes.
- Data from IoT devices can be obtained through various means, including direct data collection from sensors, network traffic monitoring, and device logs.
- Security protocols and standards such as MQTT (Message Queuing Telemetry Transport) and HTTPS (Hypertext Transfer Protocol Secure) can be used to securely transmit data from IoT devices to centralized servers or cloud platforms.
- Data preprocessing techniques, such as data cleaning, normalization, and feature extraction, may be applied to IoT data to prepare it for cybersecurity analysis.

### 3. Can K-Means method be used for Anomaly Detection? Explain how?

- While K-Means is primarily a clustering algorithm, it can be adapted for anomaly detection in certain scenarios.
- One approach is to use K-Means to cluster normal data points and identify clusters with fewer data points, which may indicate anomalies.
- Another approach is to calculate the distance of each data point to its nearest cluster centroid and flag data points with distances above a certain threshold as anomalies.
- However, K-Means may not be suitable for detecting complex or nonlinear anomalies, and other anomaly detection techniques such as Isolation Forest or One-Class SVM may be more appropriate in such cases.

## 9 Conclusion

In this Assignment, we implemented the K-Means Clustering algorithm in Python using the IOT Based Attacks dataset. We loaded the dataset, performed data analysis, split the data into dependent and independent variables, and trained the K-Means model. We visualized the clusters and analyzed the results. K-Means clustering is a powerful technique for partitioning data into distinct clusters based on similarity and can be applied to various domains, including cybersecurity and forensics.