

The IBM logo is displayed in the top left corner of the slide. The background of the slide features a blue-toned image of a computer keyboard with binary code (0s and 1s) overlaid on it.

IBM

Artificial Intelligence Practicals

Training Module

1.Introduction to Clustering

- **Definition:** Clustering is the process of grouping similar data points into clusters or groups based on certain features. The goal is to have data points in the same group more similar to each other than to data points in other groups.
 - **Types of Clustering:**
 - **Hard Clustering:** Each data point is assigned to exactly one cluster (e.g., K-means).
 - **Soft Clustering:** Data points can belong to multiple clusters with a degree of membership (e.g., Fuzzy C-Means).
 - **Applications of Clustering:**
 - **Customer Segmentation:** Grouping customers based on purchasing behavior.
 - **Image Compression:** Reducing the size of images by grouping similar pixels.
 - **Anomaly Detection:** Identifying outliers or rare events in data (e.g., fraud detection).
 - **Document Clustering:** Grouping similar documents based on content.
-

2. K-means Clustering: Theoretical Understanding

2.1 Overview of K-means Algorithm

- K-means is an unsupervised machine learning algorithm used to partition a dataset into K distinct clusters. Each cluster is represented by its centroid, which is the average of all the points within the cluster.

2.2 How K-means Works

1. **Choose K:** Select the number of clusters, K, beforehand.
2. **Initialize Centroids:** Randomly choose K points in the dataset as the initial centroids.
3. **Assign Data Points to Clusters:** Assign each data point to the nearest centroid using a distance metric (usually Euclidean distance).
4. **Recompute Centroids:** After assigning all points to clusters, recalculate the centroids of each cluster (i.e., the mean of all the points in the cluster).

- Repeat: Repeat steps 3 and 4 until the centroids no longer change (i.e., convergence is reached).

2.3 Key Concepts

- **Centroid:** The mean of all data points in a cluster, representing the "center" of the cluster.
 - **Inertia (SSE):** The sum of squared errors (SSE), or the sum of squared distances between each point and its assigned centroid.
 - **Convergence:** The algorithm stops when the centroids stabilize and don't change between iterations.
-

3. Mathematical Foundations of K-means

3.1 Objective Function: Minimize Inertia

The goal of K-means is to minimize the sum of squared distances between the points and their corresponding cluster centroids.

$$J = \sum_{i=1}^K \sum_{x_j \in C_i} \|x_j - \mu_i\|^2$$

Where:

- J is the cost function (sum of squared distances).
- K is the number of clusters.
- C_i is the set of points in the i th cluster.
- μ_i is the centroid of the i th cluster.
- x_j represents a data point in cluster C_i .

3.2 Distance Metrics

The K-means algorithm relies on a distance metric to measure how close data points are to the centroids. The most common distance metric is Euclidean distance:

$$d(x, \mu) = \sqrt{\sum_{i=1}^n (x_i - \mu_i)^2}$$

Where:

- x_i is the feature of the data point.
- μ_i is the corresponding feature of the centroid.

Other distance metrics can be used in certain contexts, such as Manhattan distance

or cosine similarity.

3.3 Choosing K (The Number of Clusters)

- The value of K (number of clusters) must be specified in advance, but there are several ways to determine the optimal K:
 - Elbow Method: Plot the SSE against different values of K and choose the K where the curve starts to flatten (the "elbow").
 - Silhouette Score: Measures how similar a point is to its own cluster compared to other clusters, with higher values indicating better clustering.
 - Gap Statistic: Compares the within-cluster dispersion with that expected under a null reference distribution of the data.
-

4. K-means Algorithm: Step-by-Step Implementation

4.1 Pseudocode

1. Choose K and initialize centroids randomly.
2. Assign each data point to the nearest centroid.
3. Recalculate the centroids based on the points assigned to them.
4. Repeat steps 2 and 3 until convergence.

4.2 Python Code Example

Here is a Python implementation using the popular scikit-learn library:

python

Copy code

```
import numpy as np
import matplotlib.pyplot as plt
from sklearn.cluster import KMeans
from sklearn.datasets import make_blobs

# Generate synthetic data
X, _ = make_blobs(n_samples=300, centers=4, random_state=42)

# Apply K-means
kmeans = KMeans(n_clusters=4)
```

```
kmeans.fit(X)
```

```
# Get cluster centers and labels
```

```
centroids = kmeans.cluster_centers_
```

```
labels = kmeans.labels_
```

```
# Plotting
```

```
plt.scatter(X[:, 0], X[:, 1], c=labels, cmap='viridis')
```

```
plt.scatter(centroids[:, 0], centroids[:, 1], c='red', marker='x', s=100)
```

```
plt.title('K-means Clustering')
```

```
plt.show()
```

This example uses the KMeans class from scikit-learn to generate clusters and visualize the results.

5. Evaluating K-means Clustering

6.1 Metrics for Evaluation

- Silhouette Score: Measures how well-separated the clusters are.

```
python
```

Copy code

```
from sklearn.metrics import silhouette_score
```

```
score = silhouette_score(X, labels)
```

```
print(f'Silhouette Score: {score}')
```

- Inertia: Sum of squared distances from each point to its assigned cluster's centroid.

```
print(f'Inertia: {kmeans.inertia_}')
```

- Adjusting Hyperparameters:

- K: The number of clusters (use Elbow or Silhouette methods).
- Initialization: K-means can be sensitive to the initial choice of centroids.

Use the k-means++ initialization method to improve convergence.

6. Advanced Topics

7.1 K-means Variants

- Mini-batch K-means: A variant of K-means that uses small random batches of data for each iteration to speed up the process.
- K-means++: An improved version of K-means that helps to spread out the initial centroids more effectively, improving convergence and reducing the chances of poor local minima.

7.2 K-means for Large Datasets

- Scalability: K-means can struggle with very large datasets. In such cases, you can:
 - Use Mini-batch K-means for faster processing.
 - Use distributed versions of K-means like Apache Spark's MLlib.
-