

Jaiwin Shah  
2019130058  
Te Comps  
Batch C

## Experiment 5

**Aim:** To train and test a machine learning model using K-Means algorithm

**Theory:**

K-Means Clustering is an unsupervised learning algorithm that is used to solve the clustering problems in machine learning or data science. It groups the unlabeled dataset into different clusters. Here K defines the number of pre-defined clusters that need to be created in the process, as if  $K=2$ , there will be two clusters, and for  $K=3$ , there will be three clusters, and so on.

How does K-Means algorithm work?

The working of the K-Means algorithm is explained in the below steps:

*Step-1:* Select the number K to decide the number of clusters.

*Step-2:* Select random K points or centroids. (It can be other from the input dataset).

*Step-3:* Assign each data point to their closest centroid, which will form the predefined K clusters.

*Step-4:* Calculate the variance and place a new centroid of each cluster.

*Step-5:* Repeat the third steps, which means reassign each datapoint to the new closest centroid of each cluster.

*Step-6:* If any reassignment occurs, then go to step-4 else go to FINISH.

*Step-7:* The model is ready.

**Code:**

```

import sys
import sklearn
import matplotlib
import numpy as np
import matplotlib.pyplot as plt
get_ipython().run_line_magic('matplotlib', 'inline')
from keras.datasets import mnist

(x_train, y_train), (x_test, y_test) = mnist.load_data()

# create figure with 3x3 subplots using matplotlib.pyplot
fig, axs = plt.subplots(3, 3, figsize = (12, 12))
plt.gray()

# loop through subplots and add mnist images
for i, ax in enumerate(axs.flat):
    ax.matshow(x_train[i])
    ax.axis('off')
    ax.set_title('Number {}'.format(y_train[i]))

# display the figure
fig.show()

# preprocessing the images

# convert each image to 1 dimensional array
X = x_train.reshape(len(x_train),-1)
Y = y_train

# normalize the data to 0 - 1
X = X.astype(float) / 255.

print(X.shape)
print(X[0].shape)

from sklearn.cluster import MiniBatchKMeans

n_digits = len(np.unique(y_test))
print(n_digits)

# Initialize KMeans model
kmeans = MiniBatchKMeans(n_clusters = n_digits)

# Fit the model to the training data
kmeans.fit(X)

kmeans.labels_

```

```

def infer_cluster_labels(kmeans, actual_labels):

    inferred_labels = {}

    for i in range(kmeans.n_clusters):

        # find index of points in cluster
        labels = []
        index = np.where(kmeans.labels_ == i)

        # append actual labels for each point in cluster
        labels.append(actual_labels[index])

        # determine most common label
        if len(labels[0]) == 1:
            counts = np.bincount(labels[0])
        else:
            counts = np.bincount(np.squeeze(labels))

        # assign the cluster to a value in the inferred_labels dictionary
        if np.argmax(counts) in inferred_labels:
            # append the new number to the existing array at this slot
            inferred_labels[np.argmax(counts)].append(i)
        else:
            # create a new array in this slot
            inferred_labels[np.argmax(counts)] = [i]

    return inferred_labels

def infer_data_labels(X_labels, cluster_labels):

    # empty array of len(X)
    predicted_labels = np.zeros(len(X_labels)).astype(np.uint8)

    for i, cluster in enumerate(X_labels):
        for key, value in cluster_labels.items():
            if cluster in value:
                predicted_labels[i] = key

    return predicted_labels

# test the infer_cluster_labels() and infer_data_labels() functions
cluster_labels = infer_cluster_labels(kmeans, Y)
X_clusters = kmeans.predict(X)
predicted_labels = infer_data_labels(X_clusters, cluster_labels)
print(predicted_labels[:20])

print(Y[:20])

from sklearn import metrics

def calculate_metrics(estimator, data, labels):

```

```

    # Calculate and print metrics
    print('Number of Clusters: {}'.format(estimator.n_clusters))
    print('Inertia: {}'.format(estimator.inertia_))
    print('Homogeneity: {}'.format(metrics.homogeneity_score(labels,
estimator.labels_)))

clusters = [10, 16, 36, 64, 144, 256]

# test different numbers of clusters
for n_clusters in clusters:
    estimator = MiniBatchKMeans(n_clusters = n_clusters)
    estimator.fit(X)

    # print cluster metrics
    calculate_metrics(estimator, X, Y)

    # determine predicted labels
    cluster_labels = infer_cluster_labels(estimator, Y)
    predicted_Y = infer_data_labels(estimator.labels_, cluster_labels)

    # calculate and print accuracy
    print('Accuracy: {}'.format(metrics.accuracy_score(Y, predicted_Y)))

X_test = x_test.reshape(len(x_test),-1)

# normalize the data to 0 - 1
X_test = X_test.astype(float) / 255.

# initialize and fit KMeans algorithm on training data
kmeans = MiniBatchKMeans(n_clusters = 256)
kmeans.fit(X)
cluster_labels = infer_cluster_labels(kmeans, Y)

# predict labels for testing data
test_clusters = kmeans.predict(X_test)
predicted_labels = infer_data_labels(kmeans.predict(X_test), cluster_labels)

# calculate and print accuracy
print('Accuracy: {}'.format(metrics.accuracy_score(y_test, predicted_labels)))

# Initialize and fit KMeans algorithm
kmeans = MiniBatchKMeans(n_clusters = 36)
kmeans.fit(X)

# record centroid values
centroids = kmeans.cluster_centers_

# reshape centroids into images
images = centroids.reshape(36, 28, 28)
images *= 255
images = images.astype(np.uint8)

# determine cluster labels

```

```
cluster_labels = infer_cluster_labels(kmeans, Y)

# create figure with subplots using matplotlib.pyplot
fig, axs = plt.subplots(6, 6, figsize = (20, 20))
plt.gray()

# loop through subplots and add centroid images
for i, ax in enumerate(axs.flat):

    # determine inferred label using cluster_labels dictionary
    for key, value in cluster_labels.items():
        if i in value:
            ax.set_title('Inferred Label: {}'.format(key))

    # add image to subplot
    ax.matshow(images[i])
    ax.axis('off')

# display the figure
fig.show()
```

## Output:

```
Number of Clusters: 10  
Inertia: 2364902.5552591607  
Homogeneity: 0.47099046206872164  
Accuracy: 0.5895333333333334
```

```
Number of Clusters: 16  
Inertia: 2201472.0373784294  
Homogeneity: 0.5668314517523975  
Accuracy: 0.6611666666666667
```

```
Number of Clusters: 36  
Inertia: 1982790.0415831676  
Homogeneity: 0.6603098306853722  
Accuracy: 0.7408333333333333
```

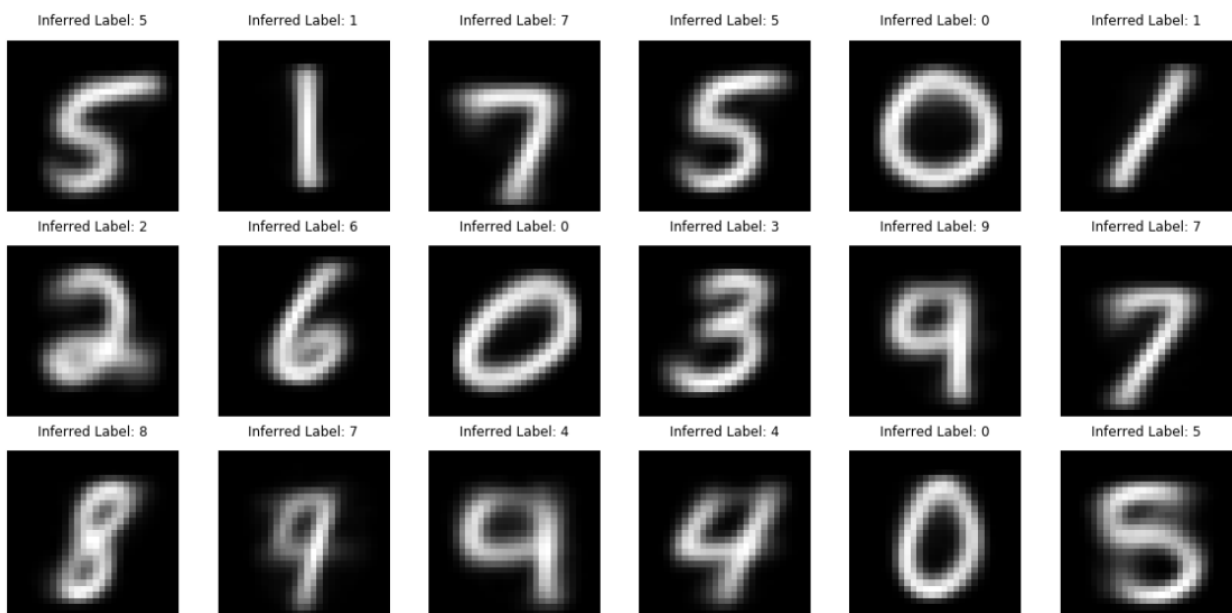
```
Number of Clusters: 64  
Inertia: 1816582.813291647  
Homogeneity: 0.7337832032469938  
Accuracy: 0.8078333333333333
```

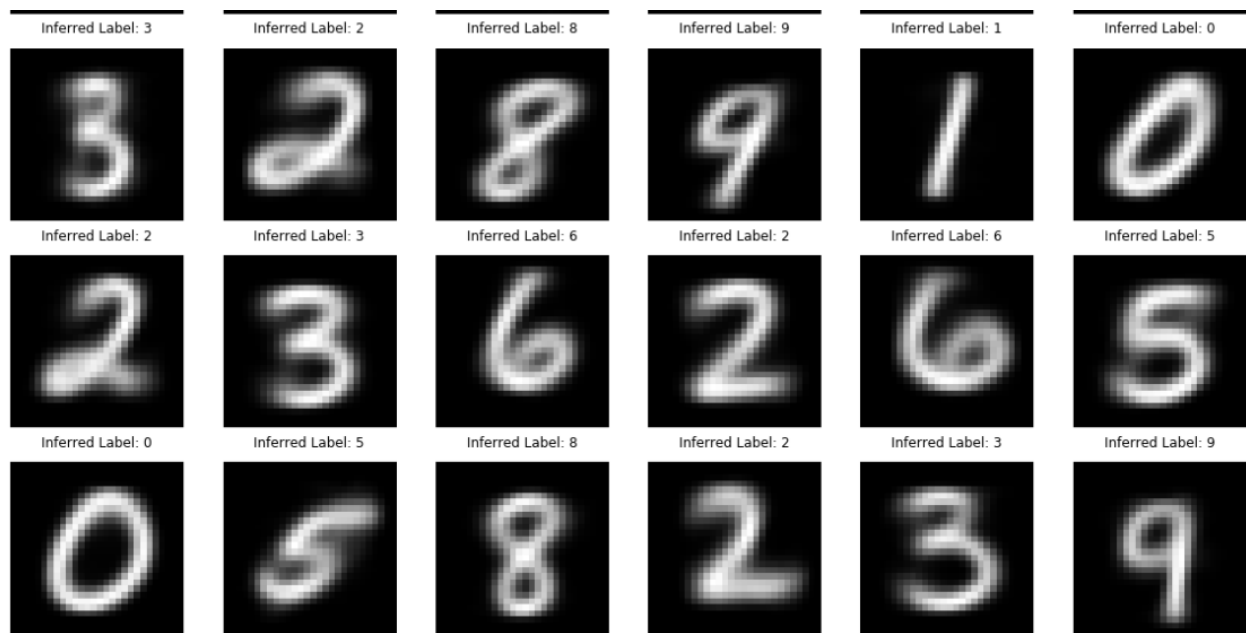
```
Number of Clusters: 144  
Inertia: 1635990.0070627772  
Homogeneity: 0.8043141694406153  
Accuracy: 0.8693
```

```
Number of Clusters: 256  
Inertia: 1514333.1809242023  
Homogeneity: 0.8419171450187173  
Accuracy: 0.8956666666666667
```

Accuracy: 0.903

D:\Users\jaiwi\anaconda3\lib\site-packages\ipykernel\_launcher.py:33: UserWarning: Matplotlib is currently using module://ipykernel.pylab.backend\_inline, which is a non-GUI backend, so cannot show the figure.





## Conclusion:

From the above experiment, I learned about the basics of the K-Means algorithm. It is a centroid-based algorithm, where each cluster is associated with a centroid.

The main aim of this algorithm is to minimize the sum of distances between the data point and their corresponding clusters.

It determines the best value for K center points or centroids by an iterative process and assigns each data point to its closest k-center. Those data points which are near to the particular k-center, create a cluster.

The algorithm takes the unlabeled dataset as input, divides the dataset into k-number of clusters, and repeats the process until it finds the best cluster. The value of k is predetermined in this algorithm.

The accuracy of the algorithm varies with the number of clusters selected.