

CSCI 6516 - Machine Learning for Big Data

(Fall 2023)

Assignment 3

Submitted by:

Anhant Dugar
800917961
ar968345@dal.ca

```
In [116]:
import pandas as pd
import seaborn as sns
import numpy as np
import matplotlib.pyplot as plt
import matplotlib.image as mpimg
from sklearn.preprocessing import MinMaxScaler
from sklearn.cluster import KMeans
from sklearn.model_selection import train_test_split
from sklearn.metrics import confusion_matrix, accuracy_score, precision_score, recall_score, f1_score, classification_r
from sklearn.decomposition import PCA
import tensorflow as tf
import tensorflow.keras as tfk
import tensorflow.keras.layers import Conv2D, MaxPooling2D, Flatten, Dense
from tensorflow.keras.optimizers import Adam
```

[1] Prepare the dataset by applying data transformation. Which method of transformation did you choose? Explain your rationality behind it.

```
In [2]:
training_data = pd.read_csv('data/sign_mnist_train.csv')
testing_data = pd.read_csv('data/sign_mnist_test.csv')
training_data.head()

Out[2]:
   label  pixel1  pixel2  pixel3  pixel4  pixel5  pixel6  pixel7  pixel8  pixel9  ...  pixel775  pixel776  pixel777  pixel778  pixel779
0      3         6    155   157   158   127   134   139   143   146   ...      207      207      207      207      207
1      2         6    187   188   188   187   186   187   188   187   ...      232      201      201      200      199
2      3         2    211   211   212   212   211   210   211   210   ...      235      234      233      231      236
4     13        13    164   167   170   172   176   179   180   184   ...      92      105      105      108      133

5 rows x 785 columns

In [3]:
print('\n\b[1:31m'+ "Check for NaN values in data : '+'\b[0m'")
training_data.isna().sum()

Check for NaN values in data :
Out[3]:
label      0
pixel1     0
pixel2     0
pixel3     0
pixel4     0
...
pixel1780  0
pixel1781  0
pixel1782  0
pixel1783  0
pixel1784  0
Length: 785, dtype: int64

In [4]:
print('\n\b[1:31m'+ "Total Records : '+'\b[0m", len(training_data))
print('\n\b[1:31m'+ "Duplicated Records : '+'\b[0m", len(training_data[training_data.duplicated()]))

Total Records : 27455
Duplicated Records : 0

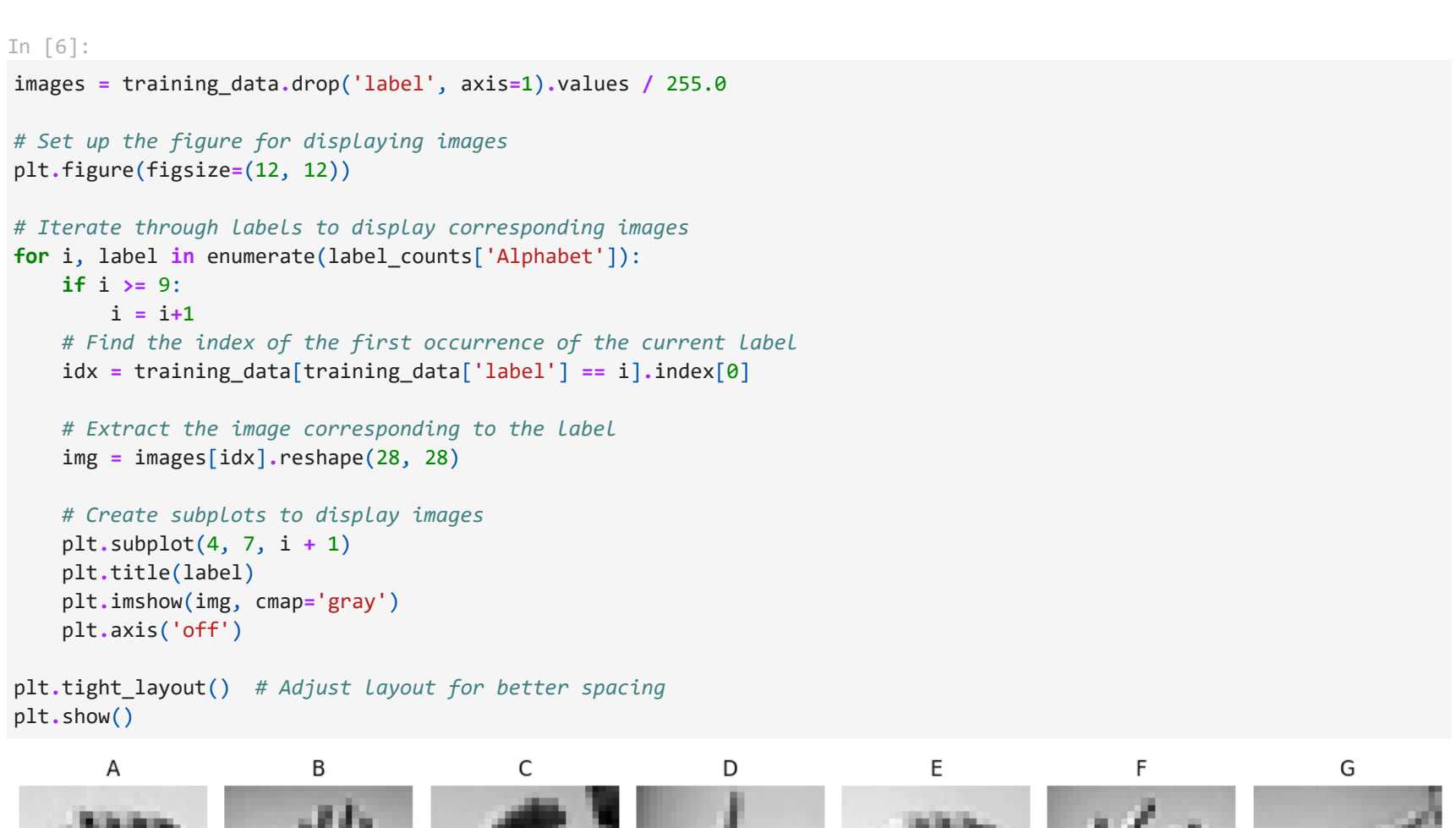
In [5]:
label_counts = training_data['label'].value_counts().reset_index().sort_values(by='label')
label_counts.reset_index(drop=True, inplace=True)
label_counts['Alphabet'] = label_counts['label'].apply(lambda x: chr(x + ord('A')))
label_counts.columns = ['label', 'count', 'Alphabet']
print(label_counts)

print('\n\n\b[1:31m'+ "Visual Analysis of data training data : '+'\b[0m\n\n")

plt.figure(figsize=(10,8))
sns.barplot(data=label_counts, x="Alphabet", y="count")
plt.title('Frequency of Alphabets/Labels in Training data')
plt.show()

Labels      Count  Alphabet
0           0    1126      A
1           1    1010      B
2           2    1144      C
3           3    1196      D
4           4     957      E
5           5    1284      F
6           6    1090      G
7           7    1013      H
8           8    1162      I
9          10    1114      K
10          11    1241      L
11          12    1055      M
12          13    1151      N
13          14    1196      O
14          15    1088      P
15          16    1279      Q
16          17    1294      R
17          18    1199      S
18          19    1186      T
19          20    1161      U
20          21    1082      V
21          22    1225      W
22          23    1164      X
23          24    1118      Y
```

Visual Analysis of data training data :



```
In [6]:
images = training_data.drop('label', axis=1).values / 255.0

# Set up the figure for displaying images
plt.figure(figsize=(12, 12))

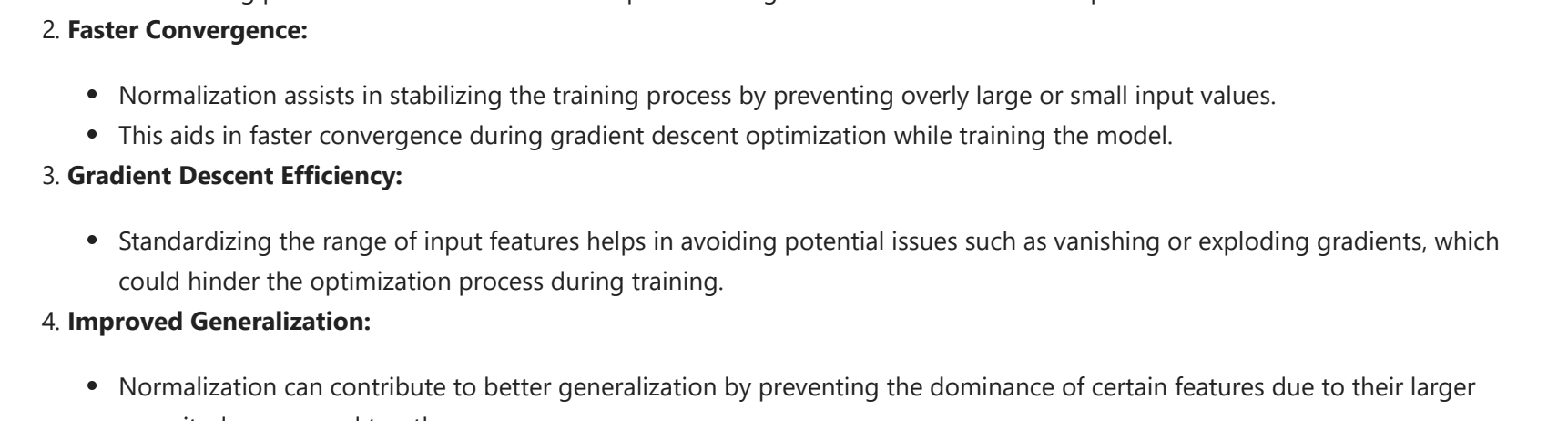
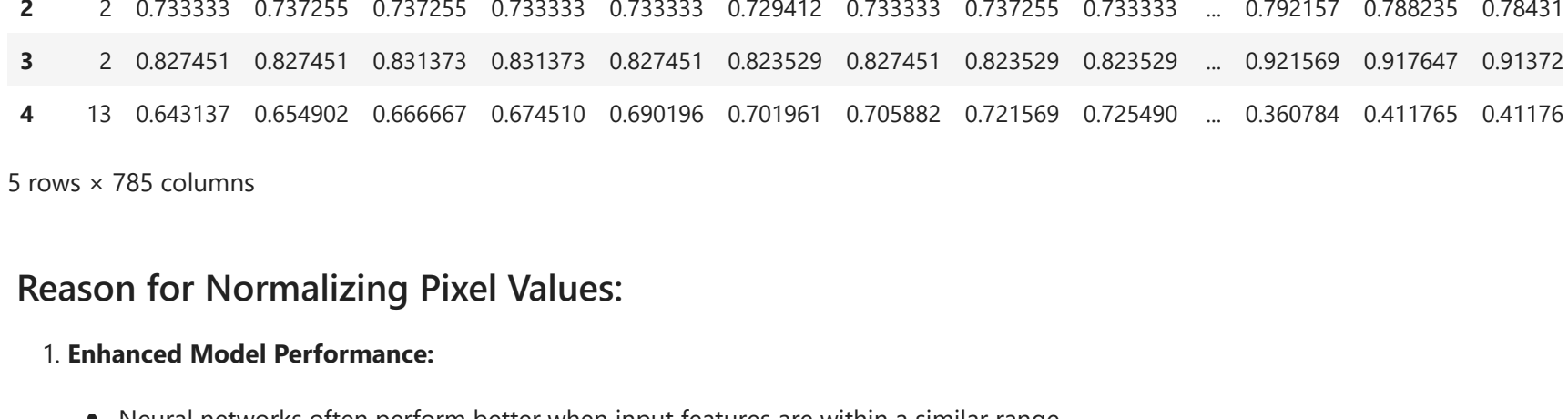
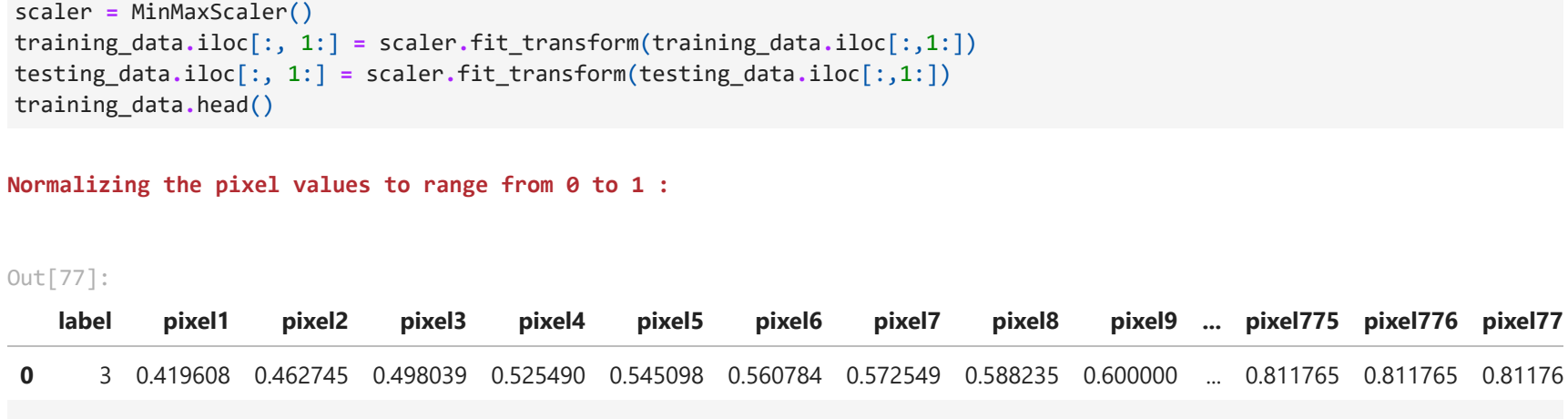
# Iterate through labels to display corresponding images
for i, label in enumerate(label_counts['Alphabet']):
    if i >= 9:
        # Find the index of the first occurrence of the current label
        idx = training_data[training_data['label'] == i].index[0]

        # Extract the image corresponding to the label
        img = images[idx].reshape(28, 28)

        # Create subplots to display images
        plt.subplot(4, 7, i + 1)
        plt.title(label)
        plt.imshow(img, cmap='gray')
        plt.axis('off')

plt.tight_layout() # Adjust layout for better spacing
plt.show()

A B C D E F G
```



```
In [7]:
# Normalizing the pixel values to range from 0 to 1, aiding in faster convergence during model training
scaler = MinMaxScaler()
training_data.iloc[:, 1:] = scaler.fit_transform(training_data.iloc[:, 1:])
testing_data.iloc[:, 1:] = scaler.fit_transform(testing_data.iloc[:, 1:])
training_data.head()

Normalizing the pixel values to range from 0 to 1 :

Out[7]:
   label  pixel1  pixel2  pixel3  pixel4  pixel5  pixel6  pixel7  pixel8  pixel9  ...  pixel775  pixel776  pixel777
0      3  0.019608  0.462745  0.498039  0.525490  0.545098  0.560784  0.617565  0.588235  0.600000  ...  0.811765  0.811765  0.811765
1      2  0.007843  0.615688  0.617656  0.611765  0.611765  0.615686  0.611765  0.619608  0.619608  ...  0.270588  0.584314  0.584314
2      2  0.273333  0.732755  0.737255  0.733333  0.733333  0.729412  0.733333  0.737255  0.733333  ...  0.792157  0.788235  0.788235
3      2  0.0827451  0.827451  0.831373  0.831373  0.827451  0.823529  0.827451  0.823529  0.823529  ...  0.921569  0.917647  0.917647
4     13  0.643137  0.654902  0.666667  0.674510  0.690196  0.701961  0.705882  0.721569  0.725490  ...  0.360784  0.411765  0.411765

5 rows x 785 columns
```

Reasons for Normalizing Pixel Values:

- Enhanced Model Performance:**
 - Neural networks often perform better when input features are within a similar range.
 - Normalizing pixel values between 0 and 1 helps in creating a consistent scale for the input features.
- Faster Convergence:**
 - Normalization assists in stabilizing the training process by preventing overly large or small input values.
 - This aids in faster convergence during gradient descent optimization while training the model.
- Gradient Descent Efficiency:**
 - Standardizing the range of input features helps in avoiding potential issues such as vanishing or exploding gradients, which could hinder the optimization process during training.
- Improved Generalization:**
 - Normalization can contribute to better generalization by preventing the dominance of certain features due to their larger magnitude compared to others.

```
In [8]:
print('\n\n\b[1:31m'+ "Visual Analysis of data training data after Normalization : '+'\b[0m\n\n")

# Set up the figure for displaying images
plt.figure(figsize=(12, 12))

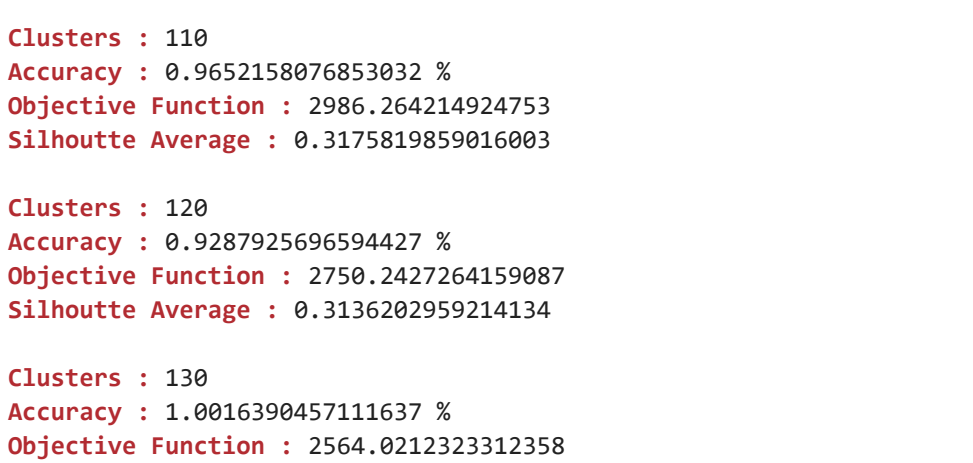
# Iterate through labels to display corresponding images
for i, label in enumerate(label_counts['Alphabet']):
    if i >= 9:
        # Find the index of the first occurrence of the current label
        idx = training_data[training_data['label'] == i].index[0]

        # Extract the image corresponding to the label
        img = images[idx].reshape(28, 28)

        # Create subplots to display images
        plt.subplot(4, 7, i + 1)
        plt.title(label)
        plt.imshow(img, cmap='gray')
        plt.axis('off')

plt.tight_layout()
plt.show()

A B C D E F G
```



```
In [165]:
X_train = training_data.drop('label', axis=1).values
y_train = training_data['label'].values

X_test = testing_data.drop('label', axis=1).values
y_test = testing_data['label'].values

# Split the data into training and validation sets (80% training, 20% testing)
X_train, X_val, y_train, y_val = train_test_split(X_train, y_train, test_size=0.2)

pca = PCA(n_components=2)
X_train, X_test, X_val = pca.fit_transform(X_train), pca.fit_transform(X_test), pca.fit_transform(X_val)

I have used PCA to reduce the number of features (or dimensions) while retaining most of the important information. In this case, by setting n_components to 2, PCA is reducing the high-dimensional data to a two-dimensional space. This reduction is beneficial for visualization purposes, especially when dealing with datasets with many features.
```

[2] Apply the k-means algorithm to Sign Language MNIST dataset.

a. Change the number of clusters from 10 to 200 with the step size of 10. Show the performance of the algorithm based on accuracy and the objective function value for each cluster number

```
In [26]:
# Initialize lists to store accuracy, inertia and silhouette Avg
accuracy_list_kmeans = []
inertia_list_kmeans = []
silhouette_avg_kmeans = []

# Range of clusters from 10 to 200 with step size 10
for i in range(10, 201, 10):
    kmeans = KMeans(n_clusters=i, random_state=42, n_init='auto')
    kmeans.fit(X_train)

    y_pred = kmeans.predict(X_val)

    # Calculate accuracy and inertia
    accuracy = accuracy_score(y_val, y_pred)
    inertia = kmeans.inertia_

    # Append accuracy and inertia to the lists
    accuracy_list_kmeans.append(accuracy)
    inertia_list_kmeans.append(inertia)
    silhouette_avg_kmeans.append(silhouette_score(X_val, y_pred))

In [27]:
j = 0
for i in range(10, 201, 10):
    print('\n\b[1:31m'+ "Clusters : '+'\b[0m", i)
    print('\n\b[1:31m'+ "Accuracy : '+'\b[0m", accuracy_list_kmeans[j]*100, "%")
    print('\n\b[1:31m'+ "Objective Function : '+'\b[0m", inertia_list_kmeans[j])
    print('\n\b[1:31m'+ "Silhouette Average : '+'\b[0m", silhouette_avg_kmeans[j])
    j+=1

Clusters : 10
Accuracy : 6.064469131385773 %
Objective Function : 30287.278941159886
Silhouette Average : 0.3240383634253947

Clusters : 20
Accuracy : 5.026468847568749 %
Objective Function : 15852.608961665832
Silhouette Average : 0.3252865386064282

Clusters : 30
Accuracy : 3.64232388258685 %
Objective Function : 10766.456262885801
Silhouette Average : 0.31561292562687624

Clusters : 40
Accuracy : 2.62247317861956 %
Objective Function : 8217.989480130326
Silhouette Average : 0.3185133829115644

Clusters : 50
Accuracy : 2.1671826625378 %
Objective Function : 6528.9788099390555
Silhouette Average : 0.31807936517529123

Clusters : 60
Accuracy : 1.2938249499180477 %
Objective Function : 5254.264048379399
Silhouette Average : 0.3219407257113174

Clusters : 70
Accuracy : 1.238390892879257 %
Objective Function : 4776.53934340957
Silhouette Average : 0.31498434049626365

Clusters : 80
Accuracy : 1.2019668548533966 %
Objective Function : 4314.9554860180572
Silhouette Average : 0.3209349862576019

Clusters : 90
Accuracy : 1.0246971407395148 %
Objective Function : 3724.3950145859084
Silhouette Average : 0.3182178535113167

Clusters : 100
Accuracy : 1.0388622837370241 %
Objective Function : 3295.2788113445415
Silhouette Average : 0.3133536793382463

Clusters : 110
Accuracy : 0.9652156876857832 %
Objective Function : 2986.246214927497
Silhouette Average : 0.3176819859016083

Clusters : 120
Accuracy : 0.9287925696594427 %
Objective Function : 2750.2427264159887
Silhouette Average : 0.313627295514134

Clusters : 130
Accuracy : 1.0013094857111637 %
Objective Function : 2564.0212323323358
Silhouette Average : 0.31135576937382463

Clusters : 140
Accuracy : 0.874157712628652 %
Objective Function : 2356.2113417410624
Silhouette Average : 0.3179985732798658

Clusters : 150
Accuracy : 0.47358209433618645 %
Objective Function : 2187.38908359469
Silhouette Average : 0.31586464048913891

Clusters : 160
Accuracy : 0.546348578379874 %
Objective Function : 2048.4871403782354
Silhouette Average : 0.31956330240823756

Clusters : 170
Accuracy : 0.589925332362047 %
Objective Function : 1936.7904113543414
Silhouette Average : 0.3193162593266838

Clusters : 180
Accuracy : 0.4917137134911673 %
Objective Function : 1816.904849087888
Silhouette Average : 0.31505814543116906

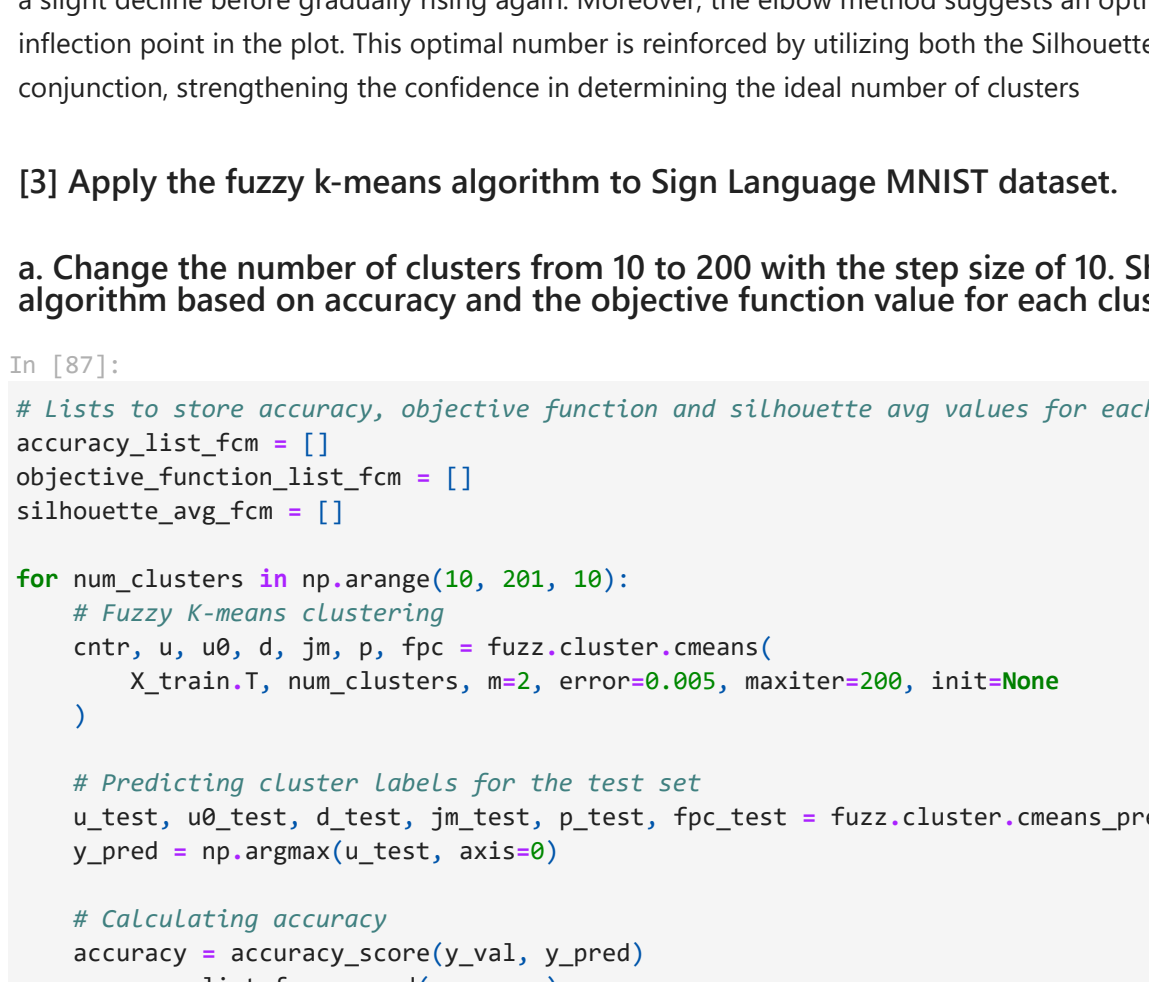
Clusters : 190
Accuracy : 0.3824439992715352 %
Objective Function : 1709.2160338308057
Silhouette Average : 0.3193565583317822

Clusters : 200
Accuracy : 0.291859042606984 %
Objective Function : 1627.5466091272528
Silhouette Average : 0.3176014854813783
```

b. What is the optimal number of clusters? Justify your answer.

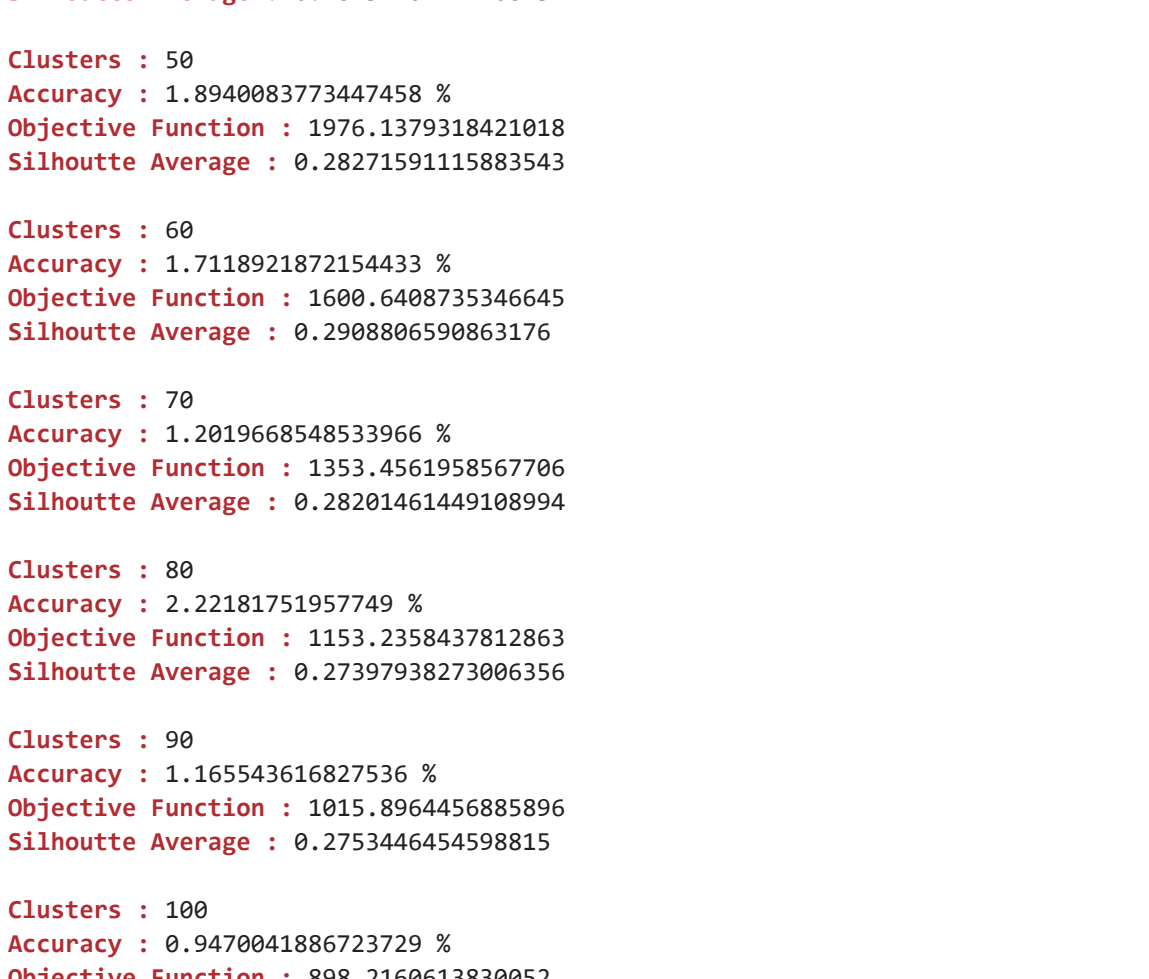
```
In [28]:
print('\n\n\b[1:31m'+ "Silhouette Analysis For Optimal k : '+'\b[0m\n\n")
plt.plot(range(10, 201, 10), silhouette_avg_kmeans, marker='o')
plt.xlabel('Values of k')
plt.ylabel('Silhouette Score')
plt.title('Silhouette Analysis For Optimal k')
plt.show()

Silhouette Analysis For Optimal k :
```



```
In [29]:
print('\n\n\b[1:31m'+ "Elbow Method For Optimal k : '+'\b[0m\n\n")
plt.plot(range(10, 201, 10), inertia_list_kmeans, marker='o')
plt.xlabel('Values of k')
plt.ylabel('Sum of squared distances/Inertia')
plt.title('Elbow Method For Optimal k')
plt.show()

Elbow Method For Optimal k :
```



```
In [91]:
num_clusters = 20
kmeans = KMeans(n_clusters=20, random_state=42, n_init='auto')
kmeans.fit(X_train)

y_pred_kmeans = kmeans.predict(X_test)

# Calculate accuracy and inertia
accuracy = accuracy_score(y_test, y_pred_kmeans)
inertia = kmeans.inertia_

# Append accuracy and inertia to the lists
print('\n\b[1:31m'+ "Clusters : '+'\b[0m", num_clusters)
print('\n\b[1:31m'+ "Accuracy : '+'\b[0m", accuracy*100, "%")
print('\n\b[1:31m'+ "Inertia : '+'\b[0m", inertia)
print('\n\b[1:31m'+ "Silhouette Average : '+'\b[0m", silhouette_score(X_val, y_pred_kmeans))

Clusters : 20
Accuracy : 3.8761851645287226 %
Inertia : 15860.71562139212
Silhouette Average : 0.3194365824699914

In [93]:
# Plot before clustering
plt.figure(figsize=(10,4))

plt.subplot(1, 2, 1)
colors = plt.cm.rainbow(np.linspace(0, 1, 20))
plt.scatter(X_test['x_test==1', 0], X_test['y_test==1', 1], marker='o', color = c, label=1, alpha=0.7)
plt.legend(fontsize=5, title='X-axis')
plt.title('Before K-means Clustering')
plt.xlabel('X-axis')
plt.ylabel('Y-axis')

# Plot after clustering
plt.subplot(1, 2, 2)
colors = plt.cm.rainbow(np.linspace(0, 1, 20))
plt.scatter(X_test['x_test==1', 0], X_test['y_pred_kmeans==1', 1], marker='o', color = c, label=1, alpha=0.7)
plt.legend(fontsize=5, title='X-axis')
plt.title('After K-means Clustering')
plt.xlabel('X-axis')
plt.ylabel('Y-axis')

plt.tight_layout()
plt.show()

Before K-means Clustering After K-means Clustering
```

The most suitable cluster count appears to be 20, as observed from the point where the Silhouette score reaches its peak, followed by a slight decline before gradually rising again. Moreover, the elbow method suggests an optimal value around 20, marking a notable inflection point in the plot. This optimal number is reinforced by utilizing both the Silhouette Method and the Elbow Method in conjunction, strengthening the confidence in determining the ideal number of clusters

[3] Apply the fuzzy k-means algorithm to Sign Language MNIST dataset.

a. Change the number of clusters from 10 to 200 with the step size of 10. Show the performance of the algorithm based on accuracy and the objective function value for each cluster number.

```
In [97]:
# Lists to store accuracy, objective function and silhouette avg values for each cluster count
accuracy_list_fcm = []
objective_function_list_fcm = []
silhouette_avg_fcm = []

for num_clusters in np.arange(10, 201, 10):
    # Fuzzy k-means clustering
    u, m, d, Jm, p, fcm = fuzz.cluster.cmeans(
        X_train, num_clusters, m=2, error=0.005, maxiter=200, init=None)

    # Predicting cluster labels for the test set
    u_test, m_test, d_test, Jm_test, p_test, fcm_test = fuzz.cluster.cmeans_predict(X_val, u, cntr, 2, error=0.005, maxiter=200, init=None)

    # Calculating accuracy
    accuracy = accuracy_score(y_val, y_pred)
    accuracy_list_fcm.append(accuracy)

    objective_function_list_fcm.append(Jm-1)
    silhouette_avg_fcm.append(silhouette_score(X_val, y_pred))

In [98]:
j = 0
for i in range(10, 201, 10):
    print('\n\b[1:31m'+ "Clusters : '+'\b[0m", i)
    print('\n\b[1:31m'+ "Accuracy : '+'\b[0m", accuracy_list_fcm[j]*100, "%")
    print('\n\b[1:31m'+ "Objective Function : '+'\b[0m", objective_function_list_fcm[j])
    print('\n\b[1:31m'+ "Silhouette Average : '+'\b[0m", silhouette_avg_fcm[j])
    j+=1

Clusters : 10
Accuracy : 3.07776163185212 %
Objective Function : 13841.2160338308057
Silhouette Average : 0.3678972280781195

Clusters : 20
Accuracy : 3.9781329448187948 %
Objective Function : 5712.478111992527
Silhouette Average : 0.2955477331794089

Clusters : 30
Accuracy : 3.7515935166636316 %
Objective Function : 3564.561226643499
Silhouette Average : 0.2899862570077815

Clusters : 40
Accuracy : 2.2404691385094267 %
Objective Function : 2552.7054999784608
Silhouette Average : 0.2893462121390873

Clusters : 50
Accuracy : 1.894008373447458 %
Objective Function : 1976.179318421818
Silhouette Average : 0.282751911583543

Clusters : 60
Accuracy : 0.6191959464396285 %
Objective Function : 1600.640873163645
Silhouette Average : 0.290886598663176

Clusters : 70
Accuracy : 1.2019668548533966 %
Objective Function : 1353.4561958567780
Silhouette Average : 0.2828614492810894

Clusters : 80
Accuracy : 2.22181751957749 %
Objective Function : 1153.2354317812863
Silhouette Average : 0.27397938273006356

Clusters : 90
Accuracy : 1.165543616827536 %
Objective Function : 1015.896456885896
Silhouette Average : 0.2753446454598815

Clusters : 100
Accuracy : 0.9470841886723729 %
Objective Function : 898.2160338308057
Silhouette Average : 0.28149367637521844

Clusters : 110
Accuracy : 0.801311236568931 %
Objective Function : 806.9501282340429
Silhouette Average : 0.28748049254587045

Clusters : 120
Accuracy : 0.656161828464849 %
Objective Function : 733.2342362732920
Silhouette Average : 0.285316667593188

Clusters : 130
Accuracy : 0.72846476051721 %
Objective Function : 673.798626584333
Silhouette Average : 0.27934807916172684

Clusters : 140
Accuracy : 0.874157712628652 %
Objective Function : 621.5397793480993
Silhouette Average : 0.27994387916172684

Clusters : 150
Accuracy : 0.6191959464396285 %
Objective Function : 575.8386137813851
Silhouette Average : 0.27588790286649987

Clusters : 160
Accuracy : 0.6738299834784192 %
Objective Function : 534.46238480415
Silhouette Average : 0.2794094742559547

Clusters : 170
Accuracy : 0.801311236568931 %
Objective Function : 504.76881663519043
Silhouette Average : 0.275171188265187

Clusters : 180
Accuracy : 0.47358209433618645 %
Objective Function : 479.358153280265
Silhouette Average : 0.2770317334931073

Clusters : 190
Accuracy : 0.4917137134911673 %
Objective Function : 430.7244621158965
Silhouette Average : 0.2669810921574156

Clusters : 200
Accuracy : 0.3824439992715352 %
Objective Function : 430.7244621158965
Silhouette Average : 0.2669810921574156
```

```
In [99]:
print('\n\n\b[1:31m'+ "Silhouette Analysis For Optimal k : '+'\b[0m\n\n")
plt.plot(range(10, 201, 10), silhouette_avg_fcm, marker='o')
plt.xlabel('Values of k')
plt.ylabel('Silhouette Score')
plt.title('Silhouette Analysis For Optimal k')
plt.show()

Silhouette Analysis For Optimal k :
```