# ML Second Project Documentation

Professor: Dr. Kiani

Ilya Shabanpour Fouladi

4003613041

GitHub Rep. link:

https://github.com/ilya-shabanpour/ML_Classification-Clustering

1402 – 1403 second semester

# Table of Contents

# Foreword

In this Documentation I want to elaborate on the implementation of my second Machine Learning Project.

Our task is to find a good classification and clustering model based on some features that were given to us in the form of a csv dataset.

Also, we had to extract features from about 400 leaf pictures and add those features to the first csv file.

In my implementation I've used these Python Libraries:

- Numpy
- Pandas
- Sklearn
- MatPlotLib
- Seaborn
- Skimage

# Feature Extraction

Our original dataset has 14 features. In order to add other features based on the pictures of leaves given to us in leaves.rar file, we need to use specific algorithms.

The First one is Hog (Histogram of Oriented Gradients). This method "captures edge orientations and intensity gradients within localized portions of an image to represent its structural features".

The second one is LBP (Local Binary Pattern). This method "encodes local texture information by thresholding neighborhood pixel values and generating a binary pattern for each pixel".

I used these two methods to extract features from the pictures.

Here is the code:

```python
import os
import pandas as pd
from skimage.feature import hog, local_binary_pattern
from skimage.io import imread
from skimage import img_as_ubyte, transform
import numpy as np
import re

# Function to extract HOG features
def extract_hog_features(image):
    # Increase pixels_per_cell and cells_per_block to reduce the number of features
    features, hog_image = hog(image, pixels_per_cell=(32, 32),
cells_per_block=(2, 2), visualize=True, multichannel=False)
    return features

# Function to extract LBP features
def extract_lbp_features(image):
    # Parameters for LBP
    radius = 3
    n_points = 8 * radius
    lbp = local_binary_pattern(image, n_points, radius, method='uniform')
    # Calculate the histogram of the LBP
    lbp_hist, _ = np.histogram(lbp.ravel(), bins=np.arange(0, n_points + 3),
range=(0, n_points + 2))
    # Normalize the histogram
    lbp_hist = lbp_hist.astype("float")
    lbp_hist /= (lbp_hist.sum() + 1e-6)
    return lbp_hist

# Function to extract both HOG and LBP features
def extract_features(image_path):
    # Read image using skimage
```

```python
    img = imread(image_path, as_gray=True)

    # Reduce image resolution to reduce the number of HOG features
    img_resized = transform.resize(img, (128, 128))

    # Convert image to unsigned byte (0-255)
    img_uint8 = img_as_ubyte(img_resized)

    # Extract HOG features
    hog_features = extract_hog_features(img_uint8)

    # Extract LBP features
    lbp_features = extract_lbp_features(img_uint8)

    # Combine HOG and LBP features
    combined_features = np.hstack((hog_features, lbp_features))

    return combined_features

# Directory containing images
image_dir = '/content/drive/MyDrive/All_Leaves'

# Initialize lists to store features and corresponding labels
all_features = []
all_labels = []
all_ex_labels = []
pattern = r'C(\d+)'
pattern2 = r'(\d+)'

# Iterate over each image in the directory
for filename in os.listdir(image_dir):
    if filename.endswith('.JPG'):
        # Extract label from filename (format is 'C01.jpg', 'C02.png', etc.)
        split = filename.split('_')
        label = split[1]
        ex_label = split[2]

        label = re.findall(pattern, label)[0]
        label = int(label)

        match = re.search(pattern2, ex_label)
        if match:
            ex_label = match.group(0)
            ex_label = int(ex_label)

        # Construct image path
        image_path = os.path.join(image_dir, filename)

        # Extract features
        features = extract_features(image_path)

        # Append features and label to lists
        all_features.append(features)
        all_labels.append(label)
        all_ex_labels.append(ex_label)

# Determine the number of features dynamically
```

```
num_features = len(all_features[0])
feature_columns = [f'feature_{i}' for i in range(num_features)]

# Create a pandas DataFrame for features and labels
leaf_data = pd.DataFrame(all_features, columns=feature_columns)
leaf_data['label'] = all_labels
leaf_data['ex_label'] = all_ex_labels

# Save DataFrame to CSV
csv_file = '/content/drive/MyDrive/leaf_hog_lbp_features.csv'
leaf_data.to_csv(csv_file, index=False)

print(f'HOG and LBP features extracted and saved to {csv_file}.')
```

I used Google Colab and uploaded the images in my drive. Then for each picture I extracted Hog and LBP features. In each image's name there is the number of the class it belongs to and the number of the sample it is.

So, I added two columns for 'Label' and 'sample' two each row so that later it would be possible to merge the first dataset and this new one.

In the end I put the extracted features and labels in a .csv file.

# Data Preprocessing

We should preprocess our data so that it will be ready to be given to our models.

At first let's look at our dataset:





As you can see the first two columns are related to labels and 364 columns are features, 14 of which are given in the original dataset and 350 columns are the extracted features. There are no missing values or NaNs in this dataset.

One Thing about this dataset is that it only has 30 unique labels. But the numbers are from 1 to 36 so we have to map each label to make it in 0 to 29 order. This code does exactly that:

```python
y = df.pop("label")
y = np.array(y)
class_samples = df.pop("ex_label")

unique_labels = np.unique(y)

label_mapping = {old_label: new_label for new_label, old_label in
enumerate(unique_labels)}
y = np.vectorize(label_mapping.get)(y)
```

# Classification

First, we split the dataset into train and test (with test_size = 0.1):

```python
x_train, x_test, y_train, y_test = train_test_split(x, y, shuffle=True,
test_size=0.1)
```

then we scale our features using Standard Scaling:

```python
scaler = StandardScaler()
x_train = scaler.fit_transform(x_train)
x_test = scaler.transform(x_test)
```

we could make Polynomial Features and add them to the features:

```python
pf = PolynomialFeatures(degree=2)
x_train_poly = pf.fit_transform(x_train)
x_test_poly = pf.transform(x_test)

x_train_poly = scaler.fit_transform(x_train_poly)
x_test_poly = scaler.transform(x_test_poly)
```

*** but after checking the test accuracies I found out that the Polynomial features are reducing my test accuracy, so I deleted them.

I used PCA and LDA for Dimension Reduction. PCA with n_components = 0.95

Keeps 95% of the most important created features. And LDA reduces the number of features to the number of classes minus 1. So in the end we have 29 refined features.

## PCA

```python
from sklearn.decomposition import PCA
pca = PCA(n_components=0.95)
x_train = pca.fit_transform(x_train)
x_test = pca.transform(x_test)
```
Executed at 2024.06.28 20:19:55 in 96ms

## LDA

```python
from sklearn.discriminant_analysis import LinearDiscriminantAnalysis as LDA
lda = LDA()
x_train = lda.fit_transform(x_train, y_train)
x_test = lda.transform(x_test)
```
Executed at 2024.06.28 20:19:55 in 35ms

## Models

I Used KNN, GNB, Adaboost, Random Forest and SVM.

### KNN:

```
KNN

from sklearn.neighbors import KNeighborsClassifier
clf = KNeighborsClassifier(n_neighbors=1, metric='euclidean')
clf.fit(x_train, y_train)
y_pred = clf.predict(x_test)
accuracy = accuracy_score(y_test, y_pred)
print("Accuracy:", accuracy)
Executed at 2024.06.28 20:19:55 in 15ms

  Accuracy: 0.9117647058823529
```

KNN finds the K Nearest neighbors of a given test data and based on the label of the neighbors judges the label of the test data. I've set the number of neighbors to 1 and distance metric to 'Euclidean' because after testing these parameters get the best result in my dataset.

### GNB:

```
NB

gnb = GaussianNB()
gnb.fit(x_train, y_train)
y_pred = gnb.predict(x_test)
accuracy = accuracy_score(y_test, y_pred)
print("Accuracy GNB:", accuracy)
Executed at 2024.06.28 20:19:55 in 39ms

  Accuracy GNB: 0.9411764705882353
```

## Random Forest:

```
Random Forest

rf = RandomForestClassifier(n_estimators=200, max_depth=20, random_state=0)
rf.fit(x_train, y_train)
y_pred = rf.predict(x_test)
accuracy = accuracy_score(y_test, y_pred)
print("Accuracy:", accuracy)
Executed at 2024.06.28 20:19:56 in 700ms

   Accuracy: 0.9411764705882353
```

Random Forest is a bagging method for decision trees. It checks many different decision trees and returns the one with the highest accuracy.

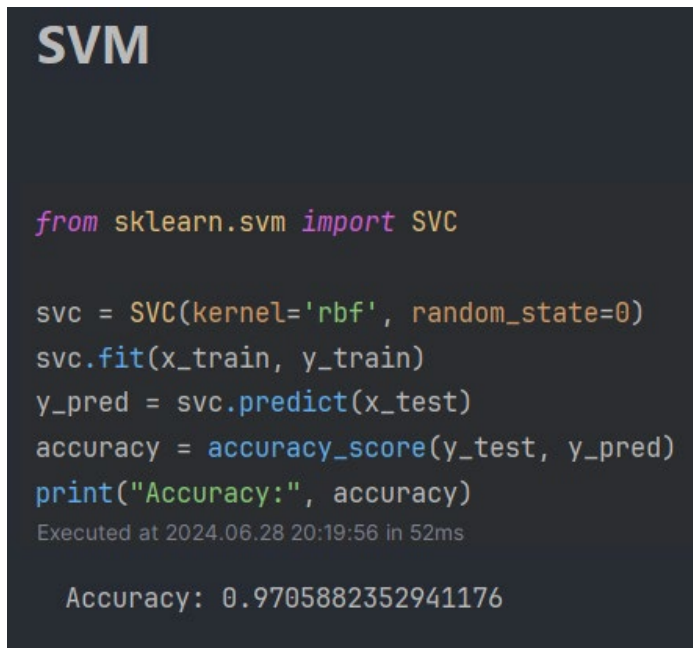Max_depth shows how deep the tree search goes.

## AdaBoost:

```
AdaBoost

from sklearn.tree import DecisionTreeClassifier
base = rf

adaBoost = AdaBoostClassifier(estimator=base, n_estimators=200, learning_rate=0.01, random_state=0)
adaBoost.fit(x_train, y_train)
y_pred = adaBoost.predict(x_test)
accuracy = accuracy_score(y_test, y_pred)
print("Accuracy:", accuracy)
Executed at 2024.06.28 21:40:24 in 546ms

   Accuracy: 0.8529411764705882
```

AdaBoost takes a base Model like decision trees or Random Forest and tries to improve the learning by focusing on the data with error in the last epoch. I gave it Random Forest as the base model. There is a chance that this method leads to Overfit.

SVM:

```
SVM

from sklearn.svm import SVC

svc = SVC(kernel='rbf', random_state=0)
svc.fit(x_train, y_train)
y_pred = svc.predict(x_test)
accuracy = accuracy_score(y_test, y_pred)
print("Accuracy:", accuracy)
Executed at 2024.06.28 20:19:56 in 52ms

  Accuracy: 0.9705882352941176
```

Support Vector Machine Model uses one with one method for more than 2 classes. The kernel is set to 'rbf' to detect non-linear relationships.

This method and Random Forest had the best test accuracies in my dataset.

## Assembling Models Using Bagging and Voting

After defining the different Models, we can use Bagging and Voting for getting a better final result:

```python
Ensemble Bagging

from sklearn.ensemble import BaggingClassifier, VotingClassifier

clf1 = KNeighborsClassifier(n_neighbors=1, metric='euclidean')
clf2 = GaussianNB()
clf3 = RandomForestClassifier(n_estimators=200, max_depth=20, random_state=0)
clf4 = SVC(kernel='rbf', random_state=0)

bagged_clf1 = BaggingClassifier(estimator=clf1, n_estimators=10, random_state=0)
bagged_clf2 = BaggingClassifier(estimator=clf2, n_estimators=10, random_state=0)
bagged_clf3 = BaggingClassifier(estimator=clf3, n_estimators=10, random_state=0)
bagged_clf4 = BaggingClassifier(estimator=clf4, n_estimators=10, random_state=0)

ensemble = VotingClassifier(estimators=[
    ('bagged_clf1', bagged_clf1),
    ('bagged_clf2', bagged_clf2),
    ('bagged_clf3', bagged_clf3),
    ('bagged_clf4', bagged_clf4)
], voting='soft')

ensemble.fit(x_train, y_train)
y_pred = ensemble.predict(x_test)
accuracy = accuracy_score(y_test, y_pred)
print(f'Ensemble Model Accuracy: {accuracy}')
Executed at 2024.06.28 21:40:28 in 3s 885ms
```

***Ensemble Model Accuracy: 0.8823529411764706***

First, we redefine the Previous models and bag them. Then we create a Voting Classifier with those bagged models and I used "soft" voting method. Soft voting means that a probability based on each models accuracy is calculated and then the weighted vote leads to the decision.

# Clustering

For Clustering I used the extracted and preprocessed features as I mentioned earlier.

In clustering there is no need for train and test splitting because we need all the data and their features.

However, we should scale the data using Standard Scaling.

## Defining Algorithms

First, we define the algorithms we want to use. I used K-Means, Agglomerative Clustering, DBSCAN and Gaussian Mixture Model:

```python
clustering_algorithms = [
    KMeans(n_clusters=30, random_state=42, n_init=20),
    AgglomerativeClustering(n_clusters=30),
    DBSCAN(eps=2.5, min_samples=5),
    GaussianMixture(n_components=30, random_state=42)
]
```

the n_components shows the K clusters we want in each algorithm. I've used the number of true labels for the number of clusters.

Esp in DBSCAN show the distance range that two points can have to form a cluster and Min_Sample shows the number of points each cluster needs to become one.

## Predicting Labels and Silhouette Score

For each algorithm we compute the clusters and their members, then we calculate their Silhouette Score:

```python
for algorithm in clustering_algorithms:
    name = type(algorithm).__name__

    if name == 'GaussianMixture':
        algorithm.fit(X_scaled)
        y_pred = np.argmax(algorithm.predict_proba(X_scaled), axis=1)
    else:
        algorithm.fit(X_scaled)
        if hasattr(algorithm, 'labels_'):
            y_pred = algorithm.labels_
        else:
```

```
        y_pred = algorithm.predict(X_scaled)

    if len(np.unique(y_pred)) == 1:
        print(f'{name} did not find any clusters.')
        continue

    if len(np.unique(y_pred)) > 1:
        silhouette = silhouette_score(X_scaled, y_pred)
    else:
        silhouette = -1
```

## Mapping Clusters to Labels

Then we map the clusters to the true labels of our dataset. The label of each cluster is the label of its most repeated member label. For example, if there are 8 samples of class 4 in a cluster and it is the highest rank, then the label of this cluster is 4:

```
def map_clusters_to_labels(y_true, y_pred):
    labels = np.unique(y_true)
    clusters = np.unique(y_pred)
    matrix = np.zeros((len(labels), len(clusters)))

    for i, label in enumerate(labels):
        for j, cluster in enumerate(clusters):
            matrix[i, j] = np.sum((y_true == label) & (y_pred == cluster))

    return np.argmax(matrix, axis=0)
```

Then we create the confusion matrix of that algorithm and store the values in a dictionary:

```
conf_matrix = confusion_matrix(y, y_pred_mapped)
accuracy = accuracy_score(y, y_pred_mapped)

results[name] = {
    'confusion_matrix': conf_matrix,
    'accuracy': accuracy,
    'silhouette_score': silhouette
}
```

Finally, we display the accuracy, Silhouette Score and Confusion Matrix for each Clustering Algorithm using MatPlotLib and Seaborn:

```
for name, result in results.items():
    print(f'{name} - Accuracy: {result["accuracy"]}, Silhouette Score:
{result["silhouette_score"]}')
    plt.figure(figsize=(10, 8))
```

```
    sns.heatmap(result['confusion_matrix'], annot=True, fmt='d',
cmap='Blues')
    plt.title(f'{name} Confusion Matrix')
    plt.xlabel('Predicted Labels')
    plt.ylabel('True Labels')
    plt.show()
```

These are the final results:

*DBSCAN did not find any clusters.*

*K-Means - Accuracy: 0.6352941176470588, Silhouette Score: 0.08869227933873192*

*Agglomerative Clustering - Accuracy: 0.6852941176470588, Silhouette Score: 0.09290624713834027*

*Gaussian Mixture - Accuracy: 0.6058823529411764, Silhouette Score: 0.0786121524800283*

DBSCAN doesn't find any clusters. DBSCAN is a Density based clustering method and it can't create clusters based on my scaled data.

The other algorithms have a low Silhouette Score but good accuracies. This is because this high dimensional multifeatured and densely close dataset of mine, makes the distance metric less meaningful and therefore gives a low Silhouette Score.

The accuracy of each algorithm is based on their matching with the true labels. This will be shown in the confusion matrix.
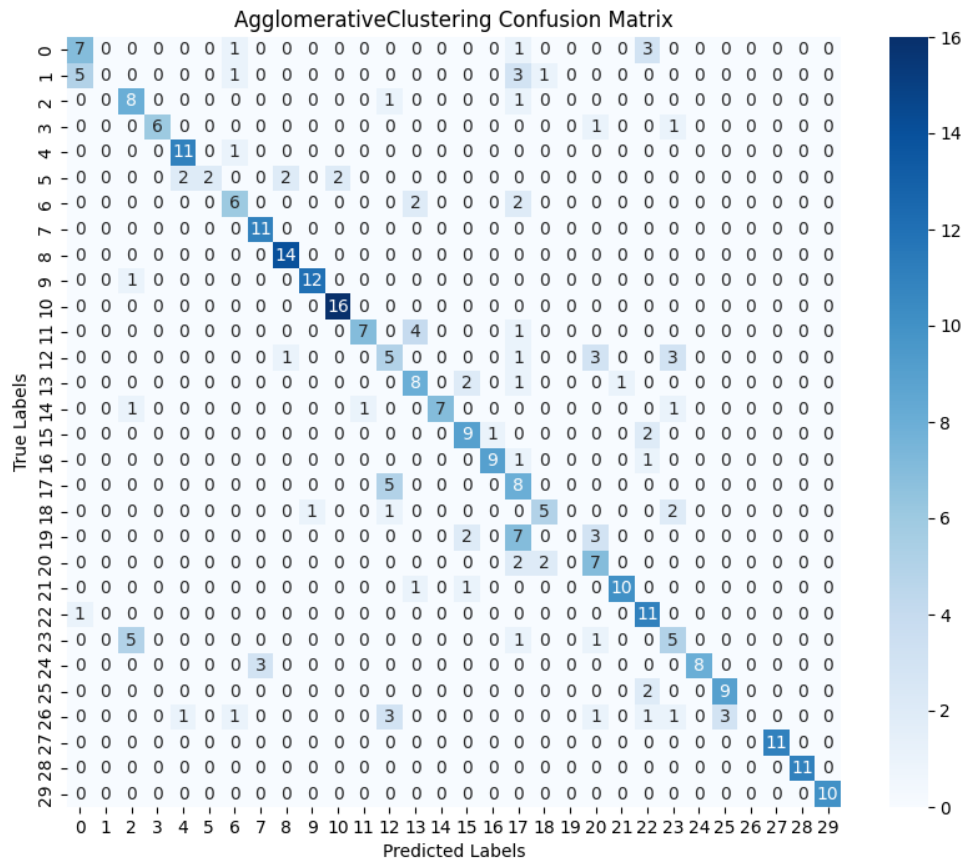
The best algorithm based on accuracy was:

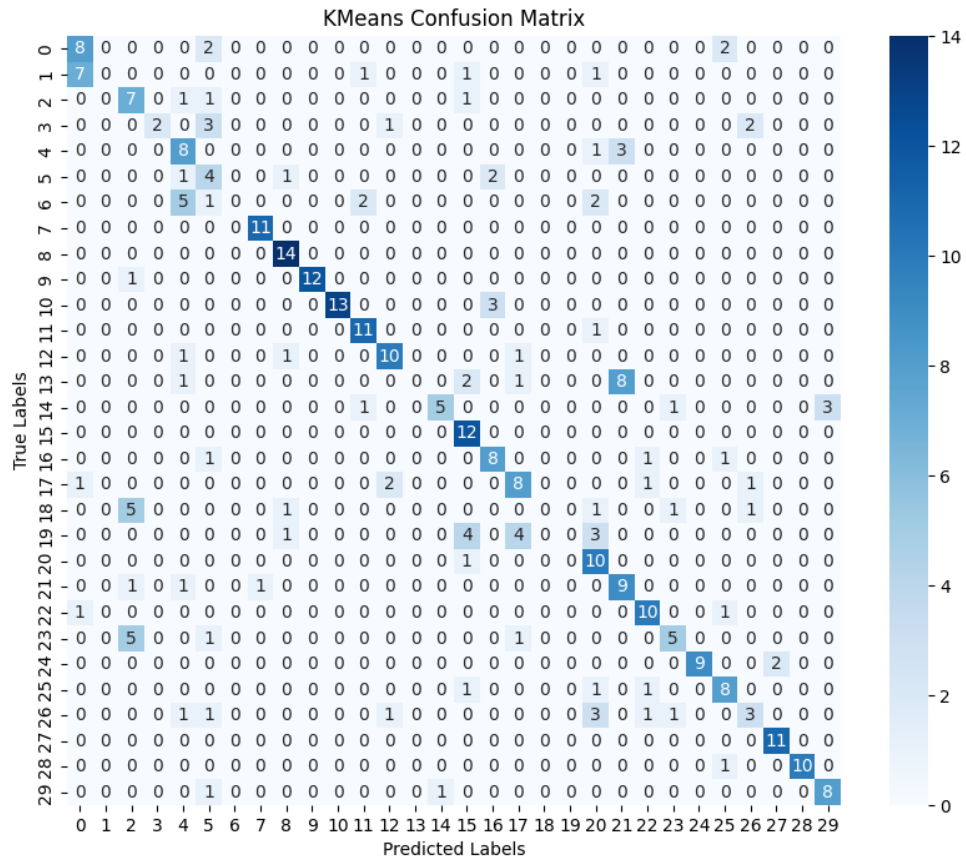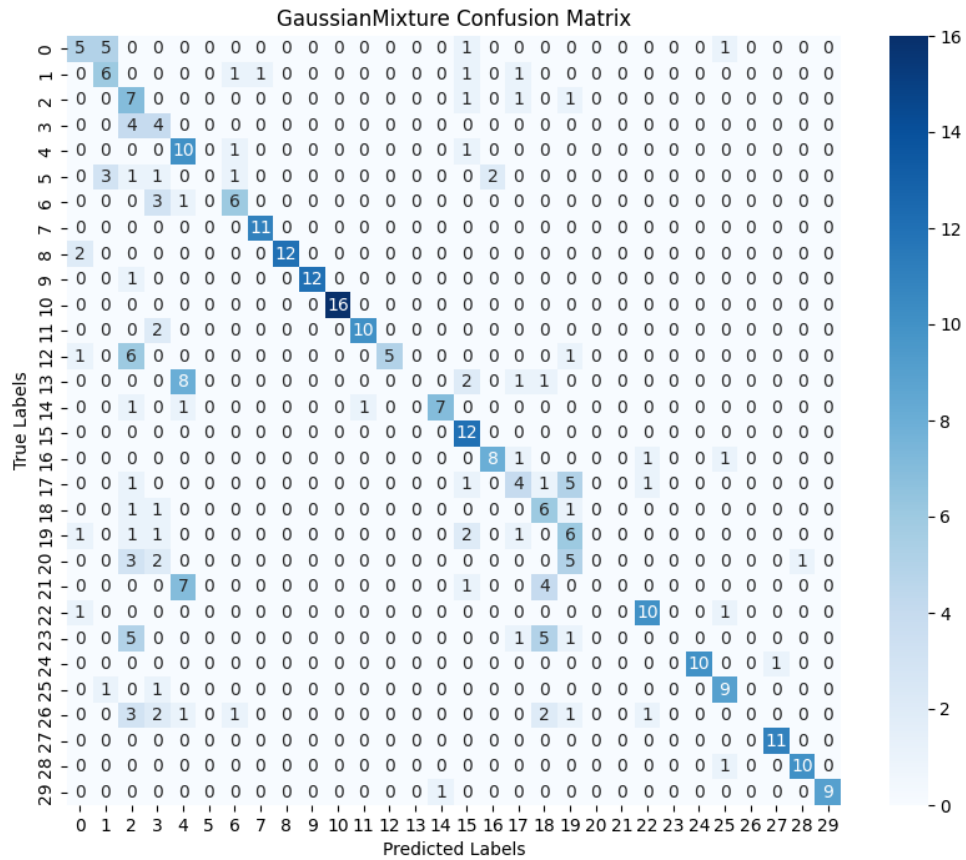1 – Agglomerative Clustering

2 – K-Means

3 – GMM

# confusion matrix



*Agglomerative Clustering - Accuracy: 0.685*

The main diameter of confusion matrix shows the truly predicted clusters. For instance, all 16 samples of the 11th class of leaves (the 10th label here), are in a single cluster and this raises the total accuracy. But the samples of 13th class of leaves (the 12th class here) have been scattered across multiple clusters and that lowers the total accuracy.

But generally, for the majority of classes, most of their samples are in one cluster and this shows a good clustering.

*K-Means - Accuracy: 0.635*

*Gaussian Mixture - Accuracy: 0.605*