

# Supervised Machine Learning: Classification

## Introduction

We will be using the Human Activity Recognition with Smartphones database, which was built from the recordings of study participants performing activities of daily living (ADL) while carrying a smartphone with an embedded inertial sensors. The objective is to classify activities into one of the six activities (walking, walking upstairs, walking downstairs, sitting, standing, and lying) to properly activities.

For each record in the dataset it is provided:

Triaxial acceleration from the accelerometer (total acceleration) and the estimated body acceleration. Triaxial Angular velocity from the gyroscope. A 561-feature vector with time and frequency domain variables. Its activity label.

```
In [1]: # Import the necessary libraries
import pandas as pd, numpy as np, matplotlib.pyplot as plt, seaborn as sns, os, sys

C:\Users\raman\AppData\Local\Programs\Python\Python37\lib\site-packages\numpy\distributor_in
it.py:32: UserWarning: loaded more than 2 Dbl from .libs:
C:\Users\raman\AppData\Local\Programs\Python\Python37\lib\site-packages\numpy\.libs\libopenbl
as.TXa6YQSD3QCQC22GEQ54J2UDCXXDhM, gfortran-win_and64.dll
C:\Users\raman\AppData\Local\Programs\Python\Python37\lib\site-packages\numpy\.libs\libopenbl
as.WCDJW7YWPW2QZME2Z2HJ3R33JXK0B7, gfortran-win_and64.dll
stacklevel=1)

In [2]: # Read the File
filepath = "Final Project\Human_Activity_Recognition_Using_Smartphones_Data.csv"
data = pd.read_csv(filepath, sep=',')

In [3]: data.head()
```

	tBodyAcc- mean(X)	tBodyAcc- mean(Y)	tBodyAcc- mean(Z)	tBodyAcc- std(X)	tBodyAcc- std(Y)	tBodyAcc- std(Z)	tBodyAcc- mag(X	tBodyAcc- mag(Y	tBodyAcc- mag(Z)	tBodyAcc- mag(X
0	0.289589	-0.020294	-0.132905	-0.995278	-0.983111	-0.913526	-0.995112	-0.983185	-0.923927	-0.934724
1	0.278439	-0.016411	-0.132905	-0.996545	-0.975300	-0.960322	-0.998807	-0.974814	-0.957696	-0.943068
2	0.278651	-0.019467	-0.113462	-0.995380	-0.967187	-0.978944	-0.996520	-0.963668	-0.977469	-0.939692
3	0.273174	-0.010501	-0.132193	-0.996991	-0.963403	-0.990765	-0.991709	-0.982750	-0.989202	-0.939692
4	0.276629	-0.016570	-0.115362	-0.998139	-0.980817	-0.990482	-0.998321	-0.979672	-0.990441	-0.942469

5 rows × 562 columns

## Data Exploration

```
In [4]: data.shape
Out[4]: (10299, 562)

The data columns are all floats except for the activity label.

In [5]: data.dtypes.value_counts()
Out[5]: float64    561
         object    1
         dtype: int64

In [6]: data.dtypes.tail()
Out[6]: angle(tBodyGyroJerkMean,gravityMean)    float64
         angle(X,gravityMean)                  float64
         angle(Y,gravityMean)                  float64
         angle(Z,gravityMean)                  float64
         Activity                             object
         dtype: object

The data are all scaled from -1 (minimum) to 1.0 (maximum).

In [7]: data.iloc[:, :-1].min().value_counts()
Out[7]: -1.0    561
         dtype: int64

In [8]: data.iloc[:, :-1].max().value_counts()
Out[8]: 1.0    561
         dtype: int64

Examine the breakdown of activities-they are relatively balanced.

In [9]: data.Activity.value_counts()
Out[9]: LAYING           1944
         STANDING       1986
         SITTING        1777
         WALKING        1722
         WALKING_UPSTAIRS  1544
         WALKING_DOWNSTAIRS 1486
         Name: Activity, dtype: int64

In [10]: # Calculate the correlation values
feature_cols = data.columns[:-1]
corr_values = data[feature_cols].corr()

# Simplify by emptying all the data below the diagonal
tril_index = np.tril_indices_from(corr_values)

# Make the unused values NaNs
corr_array = np.array(corr_values)
corr_array[np.tril_indices_from(corr_values)] = np.nan

# recreate correlation pandas dataframe
corr_values = pd.DataFrame(corr_array,columns = corr_values.columns, index = corr_values.ind
ex)

# Stack the data and convert to a dataframe
corr_values = (corr_values
               .stack()
               .to_frame()
               .reset_index()
               .rename(columns={'level_0':'feature1',
                               'level_1':'feature2',
                               0: 'correlation'}))

# Get the absolute values for sorting
corr_values['abs_correlation'] = corr_values.correlation.abs()
```

```
In [11]: from sklearn.model_selection import StratifiedShuffleSplit

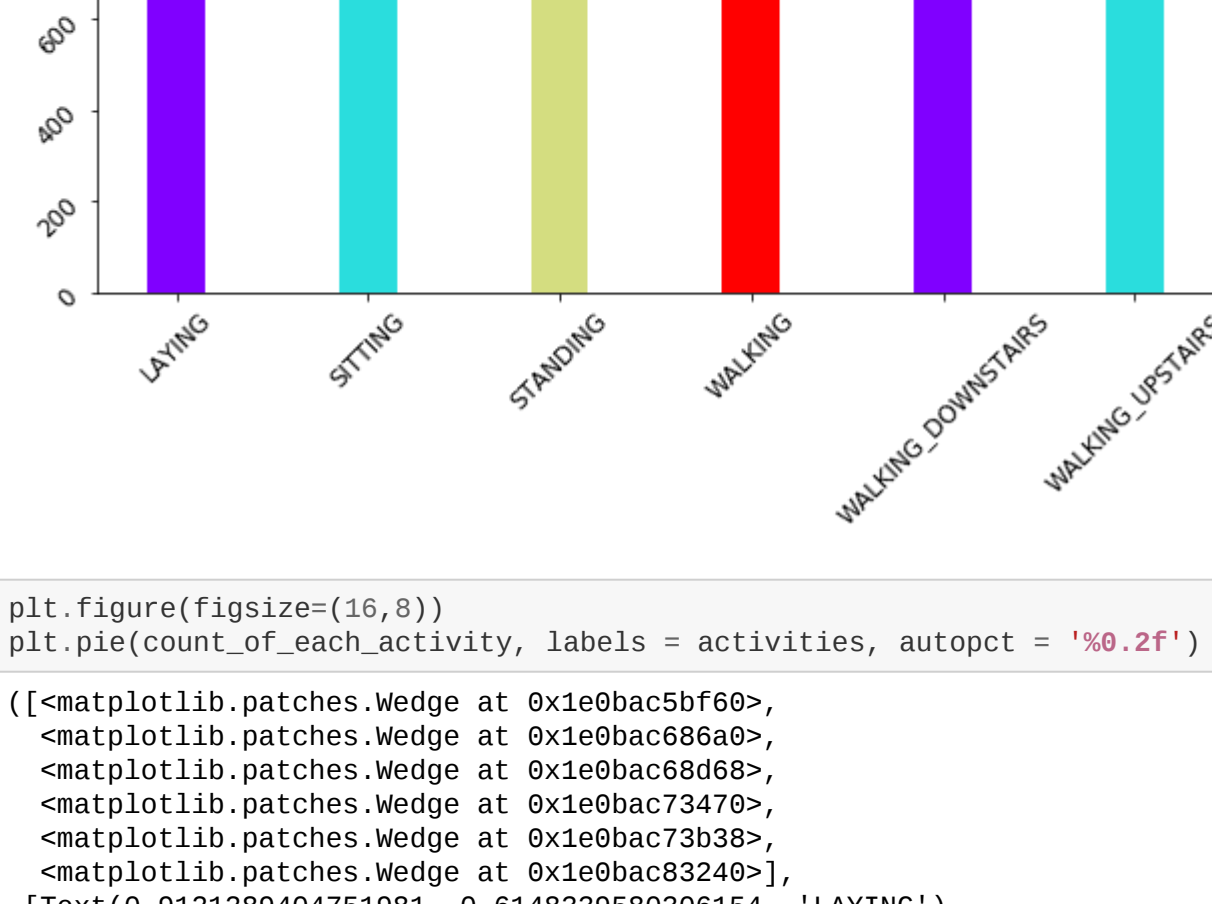
# Get the split indexes
strat_shuf_split = StratifiedShuffleSplit(n_splits=1,
                                         test_size=0.3,
                                         random_state=42)

train_idx, test_idx = next(strat_shuf_split.split(data[feature_cols], data.Activity))

# Create the dataframes
X_train = data.loc[train_idx, feature_cols]
y_train = data.loc[train_idx, 'Activity']
X_test = data.loc[test_idx, feature_cols]
y_test = data.loc[test_idx, 'Activity']

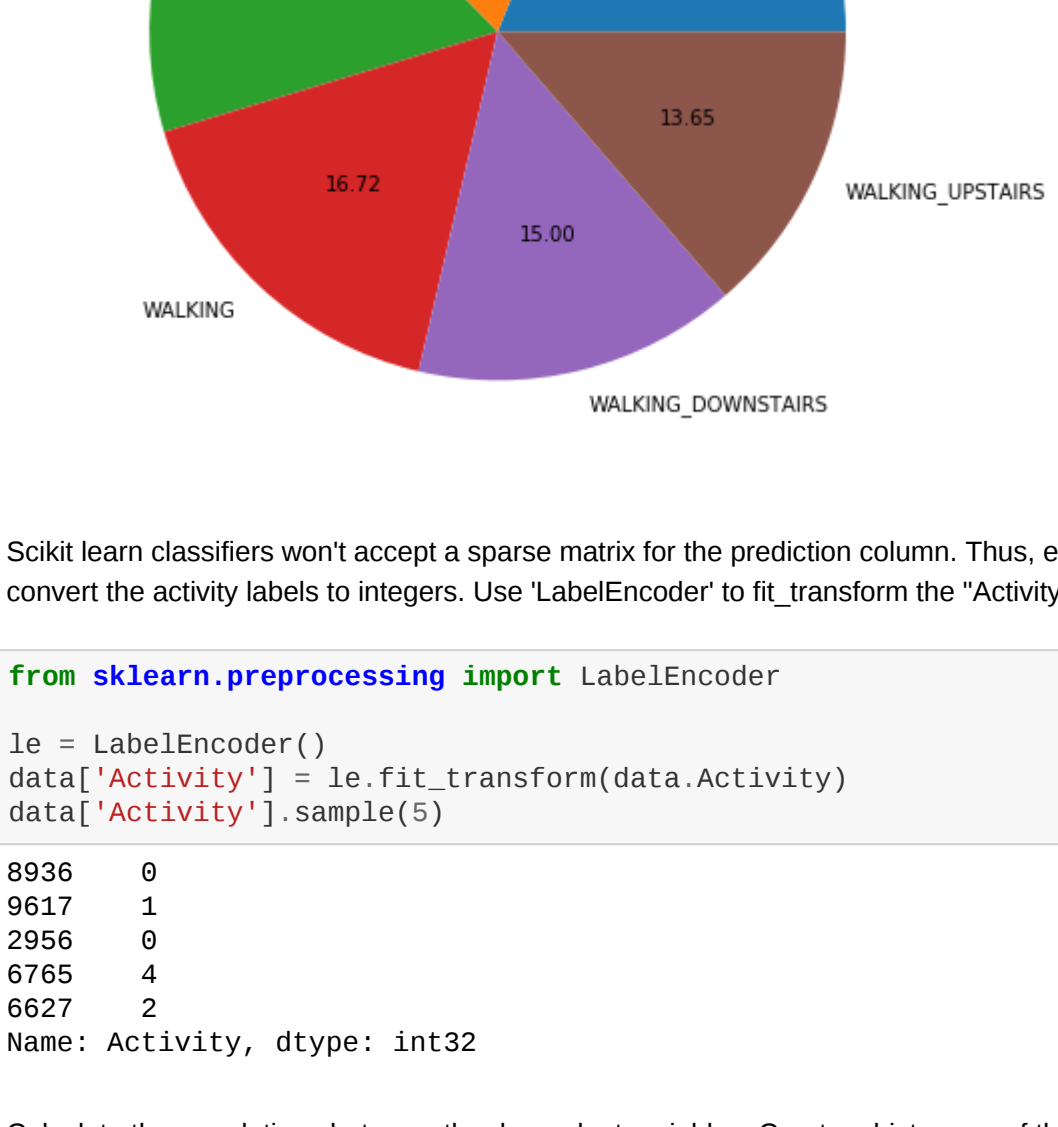
In [12]: # Data Visualization
import matplotlib.cm as cm
count_of_each_activity = np.array(y_train.value_counts())
activities = sorted(y_train.unique())
colors = cm.rainbow(np.linspace(0, 1, 4))
plt.figure(figsize=(10,6))
plt.bar(activities,count_of_each_activity,width=0.3,color=colors)
plt.xticks(rotation=45, fontsize=12)
plt.yticks(rotation=45, fontsize=12)
```

Out[12]: (array([ 0., 200., 400., 600., 800., 1000., 1200., 1400., 1600.]),  
<a list of 9 Text yticklabel objects>)



```
In [13]: plt.figure(figsize=(10,9))
plt.pie(count_of_each_activity, labels = activities, autopct = '%0.2f')

Out[13]: ([<matplotlib.patches.Wedge at 0x1e0bac5bf60>,
<matplotlib.patches.Wedge at 0x1e0bac686a0>,
<matplotlib.patches.Wedge at 0x1e0bac686a0>,
<matplotlib.patches.Wedge at 0x1e0bac73470>,
<matplotlib.patches.Wedge at 0x1e0bac73b38>,
<matplotlib.patches.Wedge at 0x1e0bac83240>],
[Text(0.912289404751385, 0.6148395898386154, 'LAYING'),
Text(-0.21504804519819878, 1.8787766559824995, 'SITTING'),
Text(-1.0656462070182775, 0.27276826372541656, 'STANDING'),
Text(-0.75212007716833865, 0.8017530555229055, 'WALKING'),
Text(0.26369462766475395, -1.0679256263152164, 'WALKING_DOWNSTAIRS'),
Text(1.00040494732078233, -0.457372812463889, 'WALKING_UPSTAIRS')],
[Text(0.4075524076622353, 0.335338977180769835, '18.48'),
Text(-0.11729457610810841, 0.588423383283815, '18.50'),
Text(-0.5812615674645149, 0.148778325668409, '17.26'),
Text(-0.410793969518167, -0.4373262848386735, '16.72'),
Text(0.34383434327168346, 0.5825646871810271, '15.40'),
Text(0.545674399315399, -0.249476079525744, '13.65')])
```



Scikit learn classifiers won't accept a sparse matrix for the prediction column. Thus, either 'LabelEncoder' needs to be used to convert the activity labels to integers. Use 'LabelEncoder' to fit 'transform' the "Activity" column, and look at 5 random values.

```
In [13]: from sklearn.preprocessing import LabelEncoder

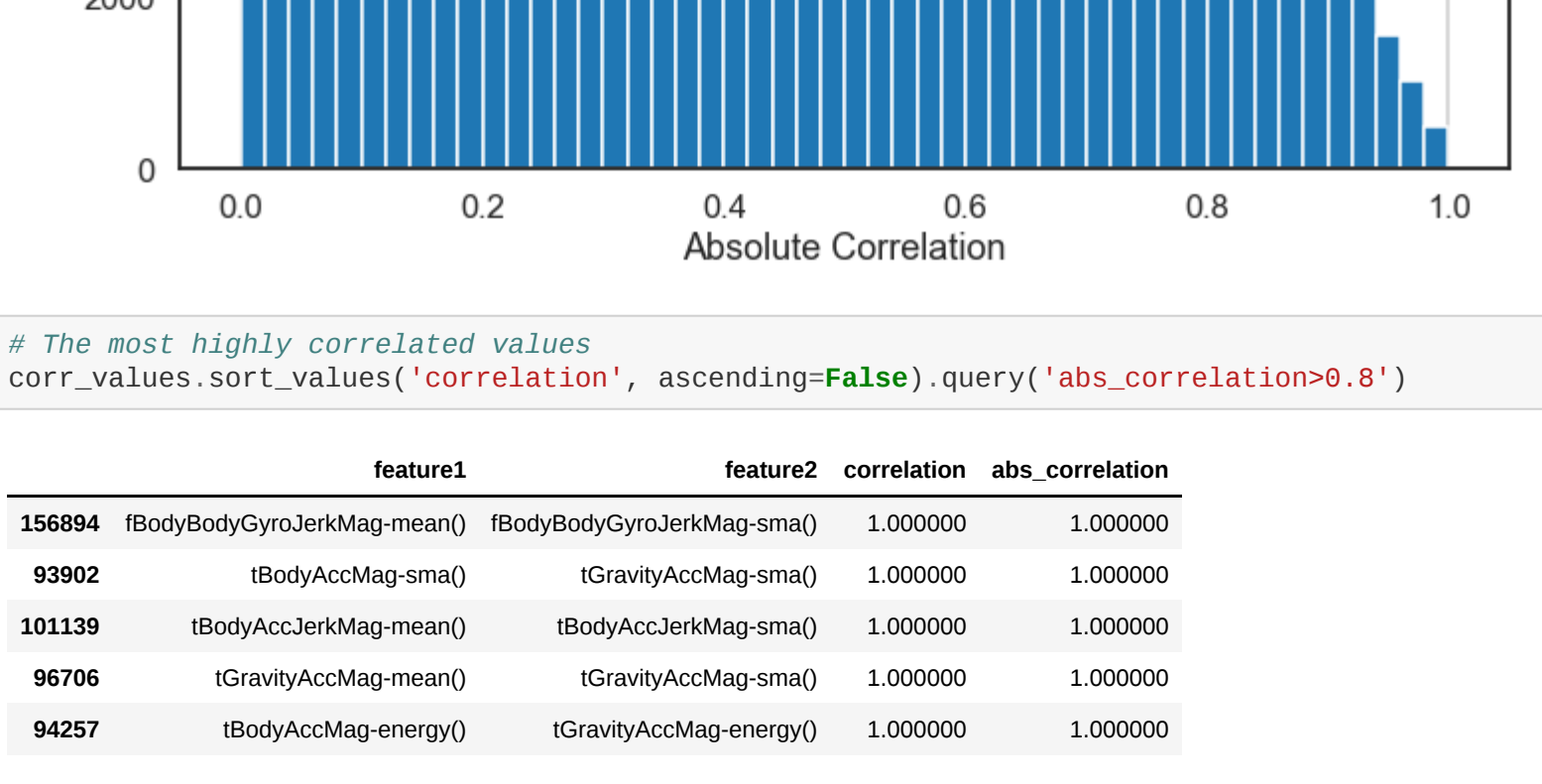
le = LabelEncoder()
data['Activity'] = le.fit_transform(data.Activity)
data['Activity'].sample(5)

Out[14]: 8936    0
         9617    1
         2956    0
         6765    0
         6627    2
         Name: Activity, dtype: int32

Calculate the correlations between the dependent variables. Create a histogram of the correlation values. Identify those that are most correlated (either positively or negatively).
```

```
In [15]: sns.set_context('talk')
sns.set_style('white')

ax = corr_values. abs_correlation.hist(bins=50, figsize=(12, 8))
ax.set(xlabel='Absolute Correlation', ylabel='Frequency')
```



```
In [16]: # The most highly correlated values
corr_values.sort_values('correlation', ascending=False).query('abs_correlation>0.8')

Out[16]:
```

	feature1	feature2	correlation	abs_correlation
156894	tBodyBodyGyroJerkMag-mean()	tBodyBodyGyroJerkMag-sma()	1.000000	1.000000
83902	tBodyAccMag-sma()	tGravityAccMag-sma()	1.000000	1.000000
101139	tBodyAccJerkMag-sma()	tBodyAccJerkMag-sma()	1.000000	1.000000
96706	tGravityAccMag-sma()	tGravityAccMag-sma()	1.000000	1.000000
94257	tBodyAccMag-energy()	tGravityAccMag-energy()	1.000000	1.000000
...	...	...	...	...
22657	tGravityAcc-mean()-Y	angle(tGravityMean)	-0.993425	0.993425
39225	tGravityAcc-acCoeff()-Z.3	tGravityAcc-acCoeff()-Z.4	-0.994267	0.994267
38739	tGravityAcc-acCoeff()-Z.2	tGravityAcc-acCoeff()-Z.3	-0.994628	0.994628
23176	tGravityAcc-mean()-Z	angle(tGravityMean)	-0.994764	0.994764
38252	tGravityAcc-acCoeff()-Z.1	tGravityAcc-acCoeff()-Z.2	-0.995195	0.995195

22815 rows × 4 columns

```
In [17]: y_train.value_counts(normalize=True)

Out[17]: LAYING           0.188792
         STANDING       0.185946
         SITTING        0.172562
         WALKING        0.167314
         WALKING_UPSTAIRS 0.149851
         WALKING_DOWNSTAIRS 0.136496
         Name: Activity, dtype: float64

In [18]: y_test.value_counts(normalize=True)

Out[18]: LAYING           0.188673
         STANDING       0.185113
         SITTING        0.172492
         WALKING        0.167314
         WALKING_UPSTAIRS 0.149838
         WALKING_DOWNSTAIRS 0.136576
         Name: Activity, dtype: float64
```

## Summary of training at least three different classifier models

```
In [19]: from sklearn.linear_model import LogisticRegression

# Standard logistic regression
lr = LogisticRegression(solver='liblinear').fit(X_train, y_train)

In [20]: y_pred = lr.predict(X_test)
```

```
In [21]: # Making the Confusion Matrix
from sklearn.metrics import confusion_matrix, accuracy_score, recall_score, f1_score
cm = confusion_matrix(y_test, y_pred)
accuracy_score=accuracy_score(y_test, y_pred)
recall_score=recall_score(y_test, y_pred, average='weighted')
f1_score=f1_score(y_test, y_pred, average='weighted')
print(y_pred)
print(cm)
print(accuracy_score)
print(recall_score)
print(f1_score)

[ 'WALKING' 'WALKING_UPSTAIRS' 'WALKING' ... 'SITTING' 'SITTING'
[ [583  0  0  0  0]
[  0 512 21  0  0]
[  0 22 550  0  0]
[  0  0  0 515  1  1]
[  0  0  0  1 420  1]
[  0  0  0  1  1 461]]
0.9841423948220665
0.9841423948220665
0.984142828415666
```

```
In [22]: accuracy_scores = np.zeros(4)

In [23]: from sklearn.metrics import accuracy_score
accuracy_scores[0] = accuracy_score(y_test, y_pred)*100
print('Logistic Regression accuracy: {0}%'.format(accuracy_scores[0]))

Logistic Regression accuracy: 98.41423948220666%
```

```
In [24]: # K Nearest Neighbors
from sklearn.neighbors import KNeighborsClassifier
clf = KNeighborsClassifier().fit(X_train, y_train)

In [25]: prediction = clf.predict(X_test)
```

```
In [26]: accuracy_scores[1] = accuracy_score(y_test, prediction)*100
print('K Nearest Neighbors Classifier accuracy: {0}%'.format(accuracy_scores[1]))

K Nearest Neighbors Classifier accuracy: 96.44012944983818%
```

```
In [27]: # Making the Confusion Matrix
from sklearn.metrics import confusion_matrix, accuracy_score, recall_score, f1_score
cm = confusion_matrix(y_test, prediction)
accuracy_score=accuracy_score(y_test, prediction)
recall_score=recall_score(y_test, prediction, average='weighted')
f1_score=f1_score(y_test, prediction, average='weighted')
print(prediction)
print(cm)
print(accuracy_score)
print(recall_score)
print(f1_score)

[ 'WALKING' 'WALKING_UPSTAIRS' 'WALKING' ... 'SITTING' 'SITTING'
[ [582  1  0  0  0]
[  1 479 52  0  0]
[  0 44 528  0  0]
[  0  0  0 517  0]
[  0  0  0  1 418  3]
[  0  0  0  1  2 451]]
0.9644012944983819
0.9644012944983819
0.9643655770411913
```

```
In [28]: # Support Vector Classifier
from sklearn.svm import SVC
svc = SVC().fit(X_train, y_train)

In [29]: prediction = svc.predict(X_test)
```

```
from sklearn.metrics import accuracy_score
accuracy_scores[2] = accuracy_score(y_test, prediction)*100
print('Support Vector Classifier accuracy: {0}%'.format(accuracy_scores[2]))

Support Vector Classifier accuracy: 97.47572815533981%
```

```
In [30]: # Making the Confusion Matrix
from sklearn.metrics import confusion_matrix, accuracy_score, recall_score, f1_score
cm = confusion_matrix(y_test, prediction)
accuracy_score=accuracy_score(y_test, prediction)
recall_score=recall_score(y_test, prediction, average='weighted')
f1_score=f1_score(y_test, prediction, average='weighted')
print(prediction)
print(cm)
print(accuracy_score)
print(recall_score)
print(f1_score)

[ 'WALKING' 'WALKING_UPSTAIRS' 'WALKING' ... 'SITTING' 'SITTING'
[ [582  0  0  0  0]
[  1 501 30  0  0]
[  0 37 535  0  0]
[  0  0  0 516  0  1]
[  0  0  0  1 418  3]
[  0  0  0  1  3 459]]
0.9747572815533981
0.9747572815533981
0.974743495290193
```

```
In [32]: # Random Forest
from sklearn.ensemble import RandomForestClassifier
rfc = RandomForestClassifier().fit(X_train, y_train)

In [33]: y_pred = rfc.predict(X_test)
```

```
In [34]: from sklearn.metrics import accuracy_score
accuracy_scores[3] = accuracy_score(y_test, y_pred)*100
print('Random Forest Classifier accuracy: {0}%'.format(accuracy_scores[3]))

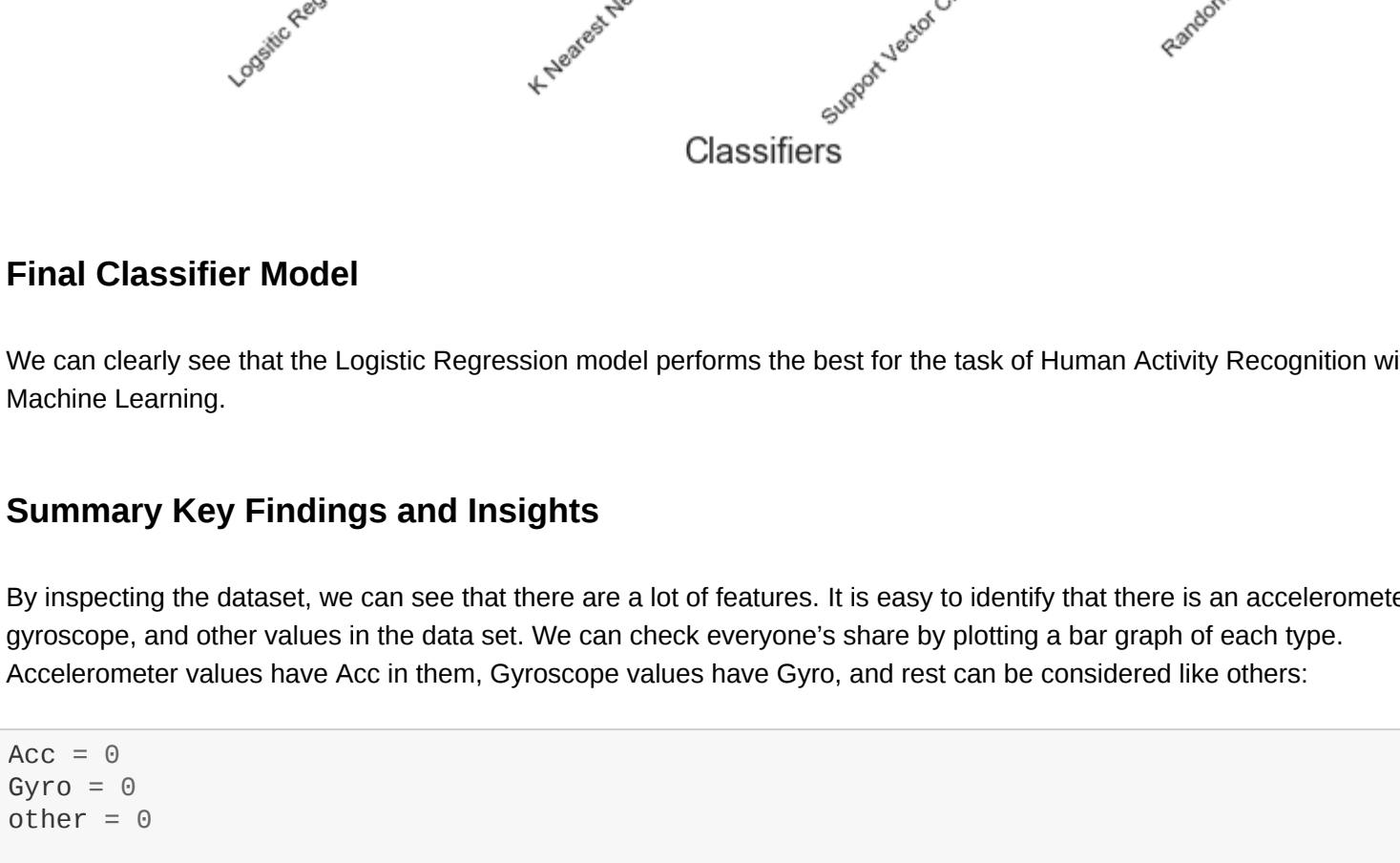
Random Forest Classifier accuracy: 97.57281553398859%
```

```
In [35]: # Making the Confusion Matrix
from sklearn.metrics import confusion_matrix, accuracy_score, recall_score, f1_score
cm = confusion_matrix(y_test, y_pred)
accuracy_score=accuracy_score(y_test, y_pred)
recall_score=recall_score(y_test, y_pred, average='weighted')
f1_score=f1_score(y_test, y_pred, average='weighted')
print(y_pred)
print(cm)
print(accuracy_score)
print(recall_score)
print(f1_score)

[ 'WALKING' 'WALKING_UPSTAIRS' 'WALKING' ... 'SITTING' 'SITTING'
[ [582  0  0  0  0]
[  1 511 21  0  0]
[  0 20 552  0  0]
[  0  0  0 506  4  7]
[  0  0  0  2 489 11]
[  0  0  0  0 8 455]]
0.9757281553398858
0.9757281553398858
0.97575729676549
```

```
In [36]: # Performance of different classifiers
import matplotlib.cm as cm
plt.figure(figsize=(12,8))
colors = cm.rainbow(np.linspace(0, 1, 4))
labels = ['Logistic Regression', 'K Nearest Neighbors', 'Support Vector Classifier', 'Random Forest']
plt.bar(labels,
        accuracy_scores,
        color = colors)
plt.xlabel('Classifiers', fontsize=18)
plt.ylabel('Accuracy', fontsize=18)
plt.title('Accuracy of various algorithms', fontsize=20)
plt.xticks(rotation=45, fontsize=12)
plt.yticks(rotation=45, fontsize=12)
```

Out[36]: (array([ 0., 20., 40., 60., 80., 100., 120.]),  
<a list of 7 Text yticklabel objects>)



The accelerometer provides the maximum functionality, followed by the gyroscope. The other features are much less so.

## Suggestions for next steps in analyzing this data

- More data can be helpful.
- More data exploration and feature engineering can also be good.
- Using different models like LSTM, neural networks, to see the performance of the data on other models.