

# Statistical Analysis of Human Activity Recognition

2023-02-27

Report By:

1. MAVILLAPALLI VENKATA TARUN KUMAR (P37000126)
2. GATTEM PRIYA MADHURI (P37000162)
3. KARRI RAHUL REDDY (P37000157)

Table of Contents:

## 1. Case Study Challenge:

1.1 Problem Statement and Importance of the Problem 1.2 Scope of Work 1.3 Data Overview 1.4 Data Understanding about features

## 2. Exploratory Data Analysis(EDA)

2.1 Data Pre-Processing and Cleaning 2.1.1 Outlier Detection 2.2 Descriptive Statistics 2.3 Multidimensional Data Analysis 2.3.1 Clustering(Hierarchical Clustering) 2.3.2 Correlation 2.3.3 Removing Correlated data 2.3.4 Lasso on uncorrelated data

## 3. Confirmatory Data Analysis(CDA)

3.1 Modeling 3.2 Checking Work Hypotheses 3.3 Fit Model(s) to Data 3.4 Interpretation 3.5 Model Assessment 3.6 Prediction: Extract and Exploit Knowledge 3.7 Model Selection and Generalization. 3.7.1 SVM 3.7.2 Logistic Regression 3.7.3 Decision Tree 3.7.4 Random Forest

## 4. Statistical Data Analysis(SDA)

4.1 Outcome Value 4.2 Future Actions 4.3 Innovation 4.4 How the beneficiary/stakeholders utilize the outcome of statistical data analysis.

-> Case Study:

Activity Recognition (AR) is monitoring the liveliness of a person by using smart phone. Smart phones are equipped in detecting sensors like gyroscope and accelerometer.

Human Activity Recognition (HAR) framework collects the raw data from sensors and observes the human movement using different approaches. These approaches are proposed to identify motions of humans with plausible high accuracy by using sensed data.

This dataset is collected from 30 persons (referred as subjects in this dataset), performing different activities with a smartphone to their waists. The data is recorded with the help of sensors (accelerometer and Gyroscope) in that smartphone. This experiment was video recorded to label the data manually.

This dataset can be used for developing a model that accurately classifies the six human activities (walking, walking upstairs, walking downstairs, sitting, standing, and laying) based on sensor data.

The goal of this study is to improve the accuracy and precision of the model by using advanced techniques such as Dimensionality Reduction, regularization, Classification models and Cross Validation.

-> Challenges:

This is a complex and high-dimensional dataset, which can make it challenging to preprocess and analyze the data effectively.

The data collected by the sensors may have a lot of noise, missing values, and outliers, which can affect the accuracy of the model.

The dataset may suffer from class imbalance, where some activities are over-represented in the data, while others are under-represented. This can lead to a biased model that performs poorly on the under-represented activities.

Another challenge is the need for effective feature selection techniques to reduce the dimensionality of the data and improve the performance of the model.

It is important to explore multiple models and compare their results with one other get best performing model.

-> Problem Statement and Importance of the Problem:

The problem statement for Human Activity Recognition using low cost sensor inputs is to develop precise and accurate prediction systems that can leverage the continuous stream of data received from accelerometer and gyroscopic sensors available from mobile phones. The challenge lies in developing algorithms that can accurately classify human activities such as WALKING, WALKING\_UPSTAIRS, WALKING\_DOWNSTAIRS, SITTING, STANDING, LAYING. The research aims to bridge the gap in existing solutions and improve real-time patient health monitoring and assistive feedback for treatment, especially for patients with heart diseases, obesity, and diabetes. The goal is to build an app for Healthcare Monitoring which uses sensors(accelerometer, gyroscope) to collect data from our mobile phone to predict Physical Activity accurately. we used the following measures(Accuracy, kappa statistics and precision ) to find the final model which best suits and predicts activity correctly. Which can positively impact the lives of millions of people across different industries, including healthcare, sports, fitness.

-> Scope of Work:

The data has been acquired in raw form. The scope of the project will be all stages including data cleaning, dimension reduction, model building, and model selection. Using validation methods, we aim to produce an accurate and precise predictive model.

-> Data Understanding:

- How data was recorded ?

1. 30 participants within an age bracket of 19-48 years(referred as subjects in this dataset) performed activities of daily living while carrying a waist-mounted smartphone. The phone was configured to record two implemented sensors (accelerometer and gyroscope). For these time series the directors of the underlying study performed feature generation and generated the dataset by moving a fixed-width window of 2.56s over the series. Since the windows had 50% overlap the resulting points are equally spaced (1.28s).This experiment was video recorded to label the data manually.
2. By using the sensors(Gyroscope and accelerometer) in a smartphone, they have captured ‘3-axial linear acceleration’(tAcc-XYZ) from accelero-meter and ‘3-axial angular velocity’ (tGyro-XYZ) from Gyroscope with several variations.

- prefix ‘t’ in those metrics denotes time.
- suffix ‘XYZ’ represents 3-axial signals in X , Y, and Z directions.

3.Each person performed six activities (WALKING, WALKING\_UPSTAIRS, WALKING\_DOWNSTAIRS, SITTING, STANDING, LAYING) wearing a smartphone (Samsung Galaxy S II) on the waist. Using its embedded accelerometer and gyroscope, we captured 3-axial linear acceleration and 3-axial angular velocity at a constant rate of 50Hz.

→ Train and Test ratio:

Data is randomly split to 70% of the volunteers were taken as training data and remaining 30% subjects' recordings were taken for test data

→ Attribute Information:

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.
  - An identifier of the subject who carried out the experiment.

—> Data Understanding about features:

As can be seen above, this dataset has 563 features with 10299 observations. That's a lot! Let's see the meaning of these features step by step.

1. Smartphone sensors captured triaxial linear acceleration (tAcc-XYZ) from accelero-meter and triaxial angular velocity (tGyro-XYZ) from gyroscope with several variations. Prefix t in those variables denotes time. Suffix XYZ represents 3-axial signals in X, Y, and Z directions.
2. These sensor signals are pre processed by applying noise filters and then sampled in fixed-width windows (sliding windows) of 2.56 seconds each with 50% overlap, that is, each window has 128 readings.
3. From each window, a feature vector was obtained by calculating variables from the time and frequency domain. In the dataset, each datapoint represents a window with different readings.
4. The acceleration signals were separated into Body and Gravity acceleration signals (tBodyAcc-XYZ and tGravityAcc-XYZ) using some low pass filter with corner frequency of 0.3 Hz.
5. After that, the body linear acceleration and angular velocity were derived in time to obtain jerk signals (tBodyAccJerk-XYZ and tBodyGyroJerk-XYZ).
6. The magnitude of these 3-dimensional signals were calculated using Euclidean norm. These magnitudes are represented as features with names like tBodyAccMag, tGravityAccMag, tBodyAccJerkMag, tBodyGyroMag, and tBodyGyroJerkMag.
7. Finally, frequency domain signals were calculated from some of the available signals by applying an FFT (Fast Fourier Transform). These obtained signals were labeled with prefix ‘f’ just like original signals with prefix ‘t’. These signals are labeled as fBodyAcc-XYZ, fBodyGyroMag, etc.
8. These are the signals that we got so far.

- tBodyAcc-XYZ
  - tGravityAcc-XYZ
  - tBodyAccJerk-XYZ
  - tBodyGyro-XYZ
  - tBodyGyroJerk-XYZ
  - tBodyAccMag
  - tGravityAccMag
  - tBodyAccJerkMag
  - tBodyGyroMag
  - tBodyGyroJerkMag
  - fBodyAcc-XYZ
  - fBodyAccJerk-XYZ
  - fBodyGyro-XYZ
  - fBodyAccMag
  - fBodyAccJerkMag
  - fBodyGyroMag
  - fBodyGyroJerkMag
9. Some set of variables are then estimated from the above signals, that is, the following properties are estimated on each and every signal that was recorded so far.
- mean(): mean value
  - std(): standard deviation
  - mad(): median absolute deviation
  - max(): largest value in array
  - min(): smallest value in array
  - sma(): signal magnitude area
  - energy(): energy measure (sum of the squares divided by the number of values)
  - iqr(): interquartile range
  - entropy(): signal entropy
  - arCoeff(): autoregresion coefficient with Burg order equal to 4
  - correlation(): correlation coefficient between two signals
  - maxInds(): index of the frequency component with the largest magnitude
  - meanFreq(): weighted average of the frequency components to obtain a mean frequency
  - skewness(): skewness of the frequency domain signal

- kurtosis(): kurtosis of the frequency domain signal
  - bandsEnergy(): energy of a frequency interval within 64 bins of the FFT of each window
  - angle(): angle between two vectors
10. Some other vectors are obtained by taking the average of signals in a single window sample. These are used on the angle() variable.
- gravityMean
  - tBodyAccMean
  - tBodyAccJerkMean
  - tBodyGyroMean
  - tBodyGyroJerkMean
11. subject feature denotes the person id. There are 30 unique ids, each for one of 30 people.
12. Activity feature represents the activity a subject was doing, consists of:

- WALKING
- WALKING\_UPSTAIRS
- WALKING\_DOWNSTAIRS
- SITTING
- STANDING
- LAYING

—> Exploratory Data Analysis(EDA):

-> Data Pre-Processing and Cleaning:

Lets start our journey into statistics with data loading and cleaning.

```
#Listing out necessary libraries.

library(randomForest)

## randomForest 4.7-1.1

## Type rfNews() to see new features/changes/bug fixes.

library(caret)

## Loading required package: ggplot2

##
## Attaching package: 'ggplot2'
```

```

## The following object is masked from 'package:randomForest':
##
##      margin

## Loading required package: lattice

library(kernlab)

##
## Attaching package: 'kernlab'

## The following object is masked from 'package:ggplot2':
##
##      alpha

library(naivebayes)

## naivebayes 0.9.7 loaded

library(glmnet)

## Loading required package: Matrix

## Loaded glmnet 4.1-6

library(ggplot2)
library(dplyr)

##
## Attaching package: 'dplyr'

## The following object is masked from 'package:randomForest':
##
##      combine

## The following objects are masked from 'package:stats':
##
##      filter, lag

## The following objects are masked from 'package:base':
##
##      intersect, setdiff, setequal, union

library(Rtsne)
library(e1071)
library(rpart)
library(rattle)

## Loading required package: tibble

```

```

## Loading required package: bitops

##
## Attaching package: 'bitops'

## The following object is masked from 'package:Matrix':
##
##      %&%

## Rattle: A free graphical interface for data science with R.
## Version 5.5.1 Copyright (c) 2006-2021 Togaware Pty Ltd.
## Type 'rattle()' to shake, rattle, and roll your data.

##
## Attaching package: 'rattle'

## The following object is masked from 'package:randomForest':
##
##      importance

library(class)
library(forecast)

## Registered S3 method overwritten by 'quantmod':
##   method           from
##   as.zoo.data.frame zoo

library(tidyverse)

## -- Attaching packages ----- tidyverse 1.3.2 --

## v tidyverse 1.3.2     v stringr 1.4.1
## v readr    2.1.3      vforcats 0.5.2
## v purrr    0.3.5
## -- Conflicts ----- tidyverse_conflicts() --
## x kernlab::alpha() masks ggplot2::alpha()
## x dplyr::combine() masks randomForest::combine()
## x purrr::cross()   masks kernlab::cross()
## x tidyverse::expand() masks Matrix::expand()
## x dplyr::filter()  masks stats::filter()
## x dplyr::lag()     masks stats::lag()
## x purrr::lift()    masks caret::lift()
## x ggplot2::margin() masks randomForest::margin()
## x tidyverse::pack() masks Matrix::pack()
## x tidyverse::unpack() masks Matrix::unpack()

library(corrplot)

## corrplot 0.92 loaded

```

```

library(Hmisc)

## Loading required package: survival
##
## Attaching package: 'survival'
##
## The following object is masked from 'package:caret':
##
##      cluster
##
## Loading required package: Formula
##
## Attaching package: 'Hmisc'
##
## The following object is masked from 'package:e1071':
##
##      impute
##
## The following objects are masked from 'package:dplyr':
##
##      src, summarize
##
## The following objects are masked from 'package:base':
##
##      format.pval, units

library(rpart)
library(MLmetrics)

##
## Attaching package: 'MLmetrics'
##
## The following objects are masked from 'package:caret':
##
##      MAE, RMSE
##
## The following object is masked from 'package:base':
##
##      Recall

library(mlbench)
library(Boruta)
library(dendextend)

##
## -----
## Welcome to dendextend version 1.16.0
## Type citation('dendextend') for how to cite the package.
##
## Type browseVignettes(package = 'dendextend') for the package vignette.
## The github page is: https://github.com/talgalili/dendextend/

```

```

## 
## Suggestions and bug-reports can be submitted at: https://github.com/talgalili/dendextend/issues
## You may ask questions at stackoverflow, use the r and dendextend tags:
##   https://stackoverflow.com/questions/tagged/dendextend
##
## To suppress this message use: suppressPackageStartupMessages(library(dendextend))
## -----
##
##
## Attaching package: 'dendextend'
##
## The following object is masked from 'package:rpart':
## 
##     prune
##
## The following object is masked from 'package:stats':
## 
##     cutree

library(plyr)

## -----
## You have loaded plyr after dplyr - this is likely to cause problems.
## If you need functions from both plyr and dplyr, please load plyr first, then dplyr:
## library(plyr); library(dplyr)
## -----
##
## Attaching package: 'plyr'
##
## The following objects are masked from 'package:Hmisc':
## 
##     is.discrete, summarize
##
## The following object is masked from 'package:purrr':
## 
##     compact
##
## The following objects are masked from 'package:dplyr':
## 
##     arrange, count, desc, failwith, id, mutate, rename, summarise,
##     summarize

library(reshape2)

## 
## Attaching package: 'reshape2'
## 
## The following object is masked from 'package:tidyR':
## 
##     smiths

```

```

library(lattice)
library(MASS)

## 
## Attaching package: 'MASS'
##
## The following object is masked from 'package:dplyr':
## 
##     select

library(psych)

## 
## Attaching package: 'psych'
##
## The following object is masked from 'package:MLmetrics':
## 
##     AUC
## 
## The following object is masked from 'package:Hmisc':
## 
##     describe
## 
## The following object is masked from 'package:kernlab':
## 
##     alpha
## 
## The following objects are masked from 'package:ggplot2':
## 
##     %+%, alpha
## 
## The following object is masked from 'package:randomForest':
## 
##     outlier

#Setting the working directory path to access the file
setwd("C:/Users/PRIYA/Downloads/UCI HAR Dataset")

#Reading the file "X_train.txt" from the "train" sub-directory of the "UCI HAR Dataset" directory, and st
train_X<-read.table("UCI HAR Dataset/train/X_train.txt")

#Reading the file "y_train.txt" from the "train" sub-directory of the "UCI HAR Dataset" directory, and st
train_Y<-read.table("UCI HAR Dataset/train/y_train.txt")

#Reading the file "X_test.txt" from the "test" sub-directory of the "UCI HAR Dataset" directory, and st
test_X<-read.table("UCI HAR Dataset/test/X_test.txt")

#Reading the file "y_test.txt" from the "test" sub-directory of the "UCI HAR Dataset" directory, and st
test_Y<-read.table("UCI HAR Dataset/test/y_test.txt")

```

```

#Reading the file "features.txt" from the "UCI HAR Dataset" directory and storing its contents in a data frame
col<- read.table("UCI HAR Dataset/features.txt")

#Reading the file "features.txt" from the "UCI HAR Dataset" directory and stores its contents as a character vector
col_names <- readLines("UCI HAR Dataset/features.txt")

#Here we are setting the column names of the data frame "train_X" to a cleaned-up version of the character vector
colnames(train_X)<-make.names(col_names)

#Here we are setting the column names of the data frame "test_X" to a cleaned-up version of the character vector
colnames(test_X)<-make.names(col_names)

#Setting the name of the single column in the "train_Y" data frame to "activity".
colnames(train_Y)<-"activity"

#Setting the name of the single column in the "test_Y" data frame to "activity".
colnames(test_Y)<-"activity"

#Reading the file "subject_test.txt" from the "test" sub-directory of the "UCI HAR Dataset" directory and storing its contents in a data frame
temp5 <- read.table("UCI HAR Dataset/test/subject_test.txt")

#Reading the file "subject_train.txt" from the "train" sub-directory of the "UCI HAR Dataset" directory and storing its contents in a data frame
temp6 <- read.table("UCI HAR Dataset/train/subject_train.txt")

#Setting the name of the single column in the "temp5" data frame to "subject".
colnames(temp5)<-"subject"

#Setting the name of the single column in the "temp6" data frame to "subject".
colnames(temp6)<-"subject"

#Combining the "test_Y" and "test_X" data frames by column-binding them together, and storing the result in a data frame
X <- cbind(test_Y, test_X)

#Combining the "train_Y" and "train_X" data frames by column-binding them together, and storing the result in a data frame
Y <- cbind(train_Y, train_X)

#Combining the "temp6" and "Y" data frames by column-binding them together, and storing the result in a data frame
train <- cbind( temp6, Y)

#Combining the "temp5" and "X" data frames by column-binding them together, and storing the result in a data frame
test<- cbind(temp5,X)

#Adding three new columns to the "train" and "test" data frames: "subjectID", "activityID", and "partition"
#The "transform()" function is used to create the "subjectID" and "activityID" columns, which convert the categorical variables into factors
train <- transform(train, subjectID = factor(subject), activityID = factor(activity))
test <- transform(test, subjectID = factor(subject), activityID = factor(activity))

```

```

#Setting the value of the "partition" column in the "train" data frame to the character string "train".
train$partition = "train"
#Setting the value of the "partition" column in the "test" data frame to the character string "test".
test$partition = "test"

#Combining the "train" and "test" data frames by row-binding them together, and storing the result in a
final_data <- rbind(train,test)

#converting the "activity" column in the "final_data" data frame to a factor.
final_data$activity<-factor(final_data$activity)

#Setting the levels of the "activity" factor in the "final_data" data frame to a new set of labels, spe
levels(final_data$activity) <- c("WALKING", "WALKING_UPSTAIRS", "WALKING_DOWNSTAIRS", "SITTING", "STAND

#Returning a vector containing the dimensions of the "final_data" data frame in the form (number of row
dim(final_data)

## [1] 10299 566

```

As we can see from above the final data has 10299 rows and 566 columns, in which 4 columns are added for our understanding purpose and visualization which are :

1. Subject : This column helps us to identity of data subjects(i.e 30 subjects).
2. Subject ID: Is the factor of subject
3. Activity ID : Is the factor of activity
4. Partition : Identity the testing and training the data

—> Checking duplicates

Let's check if there are any duplicated observations or missing values.

```
cat("Number of duplicated rows:", sum(duplicated(final_data)))
```

```
## Number of duplicated rows: 0
```

```
cat("Number of missing values:", sum(is.na(final_data)))
```

```
## Number of missing values: 0
```

Great! None of them exists. so we are proceeding forward with descriptive statistics.

—> Descriptive Statistics:

```
# Create contingency table of activity counts
activity_table <- table(final_data$activity)
activity_table
```

```
##
##          WALKING   WALKING_UPSTAIRS   WALKING_DOWNSTAIRS           SITTING
##          1722                 1544                  1406                1777
##          STANDING               LAYING
##          1906                 1944
```

Contingency tables are used to summarize and display the distribution of categorical variables.

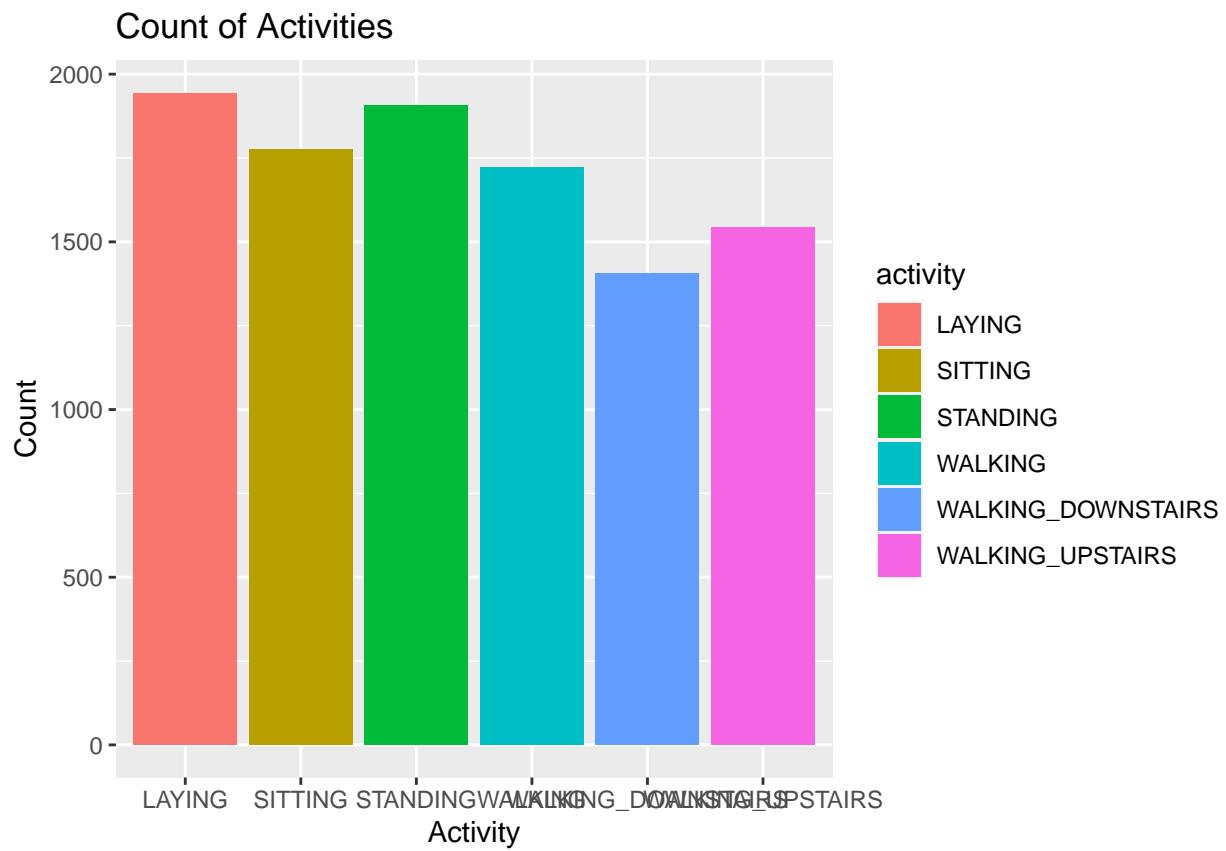
Creating a contingency table of activity counts allows us to see the frequency of each activity in the dataset.

```
# Convert to data frame
activity_counts <- data.frame(activity = names(activity_table),
                                count = as.numeric(activity_table))
```

Creating a bar plot

This code creates a bar plot using ggplot that displays the count of activities on the y-axis and the name of each activity on the x-axis. Each bar is filled with the corresponding activity color. The plot title is “Count of Activities”.

```
ggplot(activity_counts, aes(x = activity, y = count, fill = activity)) +
  geom_bar(stat = "identity") +
  xlab("Activity") +
  ylab("Count") +
  ggtitle("Count of Activities")
```



We have observed that laying and standing having the highest count.

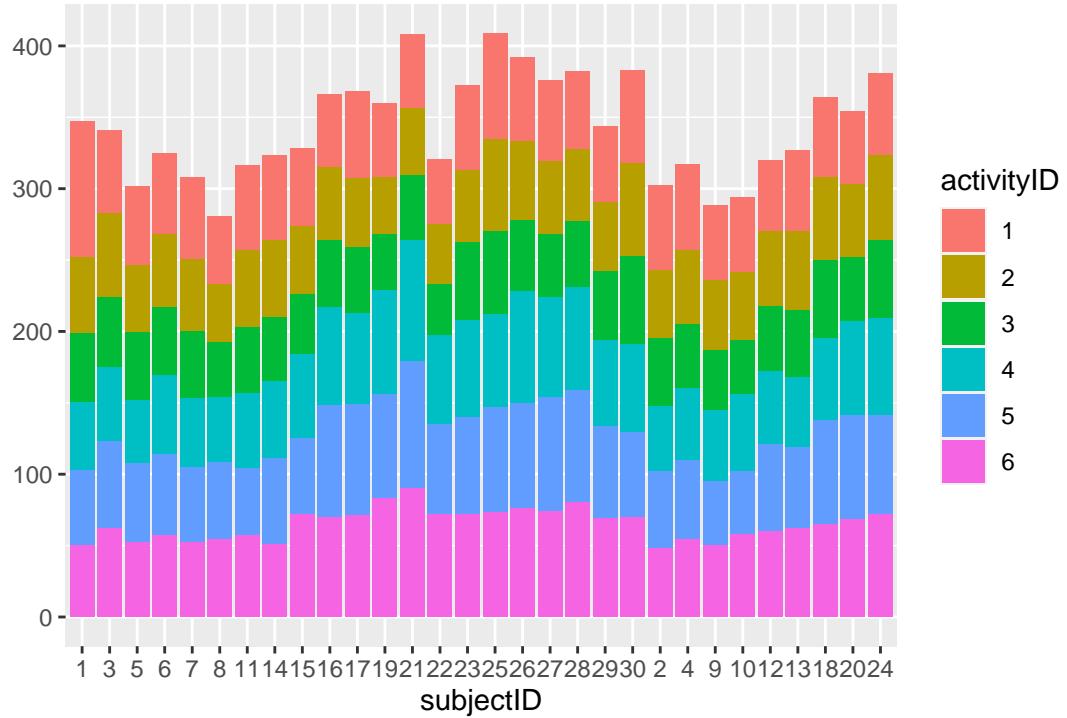
-> checking for class imbalance

Class imbalance refers to a situation where the number of observations in different classes or categories of a categorical variable is not evenly distributed. Class imbalance can affect the accuracy of statistical models because the model tends to be biased towards the majority class, resulting in poor performance for the minority class.

SO lets check for class imbalance in categorical variable(activity) in our data set.

```
par(mfrow=c(1,2))
#shows the distribution of subjects across different activities
qplot(data= final_data,x= subjectID,fill= activityID) + theme(plot.margin= unit(c(1,1,1,1),"cm"))

## Warning: 'qplot()' was deprecated in ggplot2 3.4.0.
```

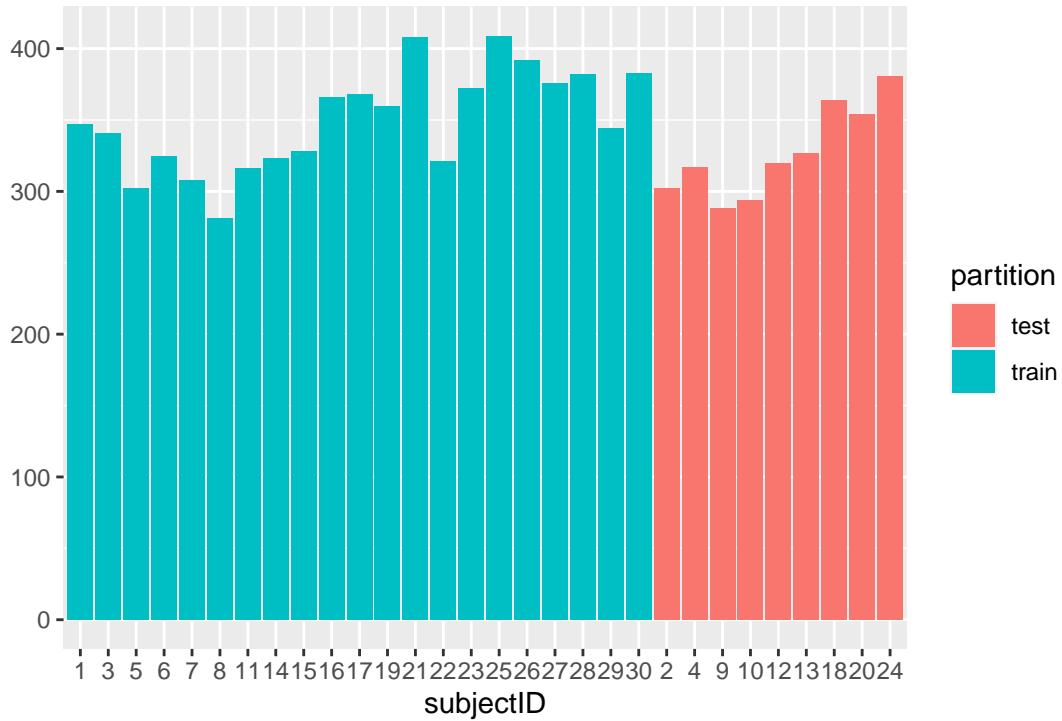


From the above graph we understood activities of different subjects.

There is no significant difference in terms of activity count for each subject. Hence, all target classes are balanced.

Now lets move forward and understand the test and train split of our data.

```
#shows the distribution of subjects across train and test partitions
qplot(data = final_data, x = subjectID, fill = partition) + theme(plot.margin= unit(c(1,1,1,1),"cm"))
```



```
par(mfrow=c(1,1))
```

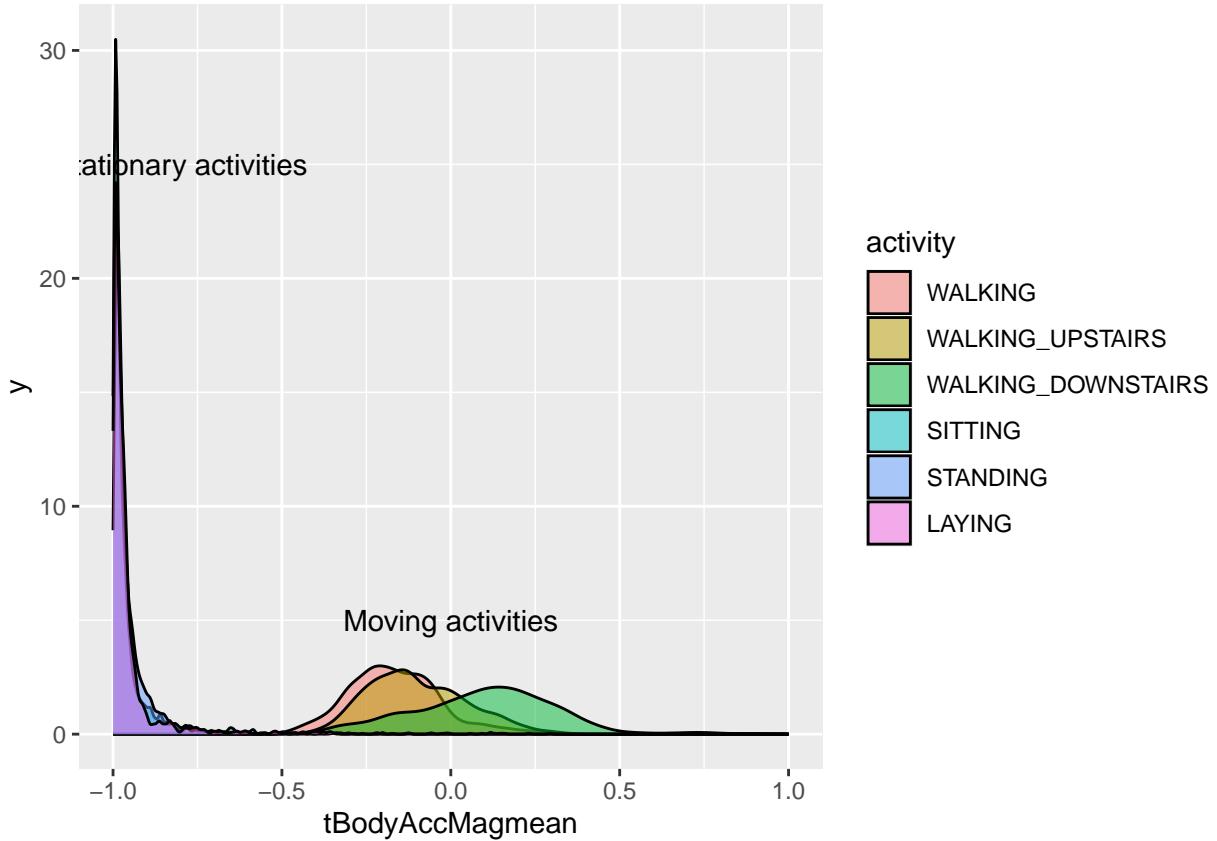
This graph tells about the split of train and test data among the subjects. which is in the ratio of 70 and 30 percent.

lets dig deep into different activity types using tBodyAccMagmean column.

tBodyAccMagmean -> Describes the mean of the magnitude of the body acceleration signal measured by the accelerometer sensor in the time domain. It is a summary statistic that provides an estimate of the average amount of body acceleration across all three axes during a specific time period. This feature can be useful for predicting the type of activity being performed based on the patterns of body acceleration observed by the smartphone's accelerometer sensor.

```
colnames(final_data)[203] = "tBodyAccMagmean"

ggplot(final_data,
       aes(x =tBodyAccMagmean, group = activity, fill = activity)) +
  geom_density(alpha = .5) +
  annotate('text', x = -.8, y = 25, label = "Stationary activities") +
  annotate('text', x = -.0, y = 5, label = "Moving activities")
```



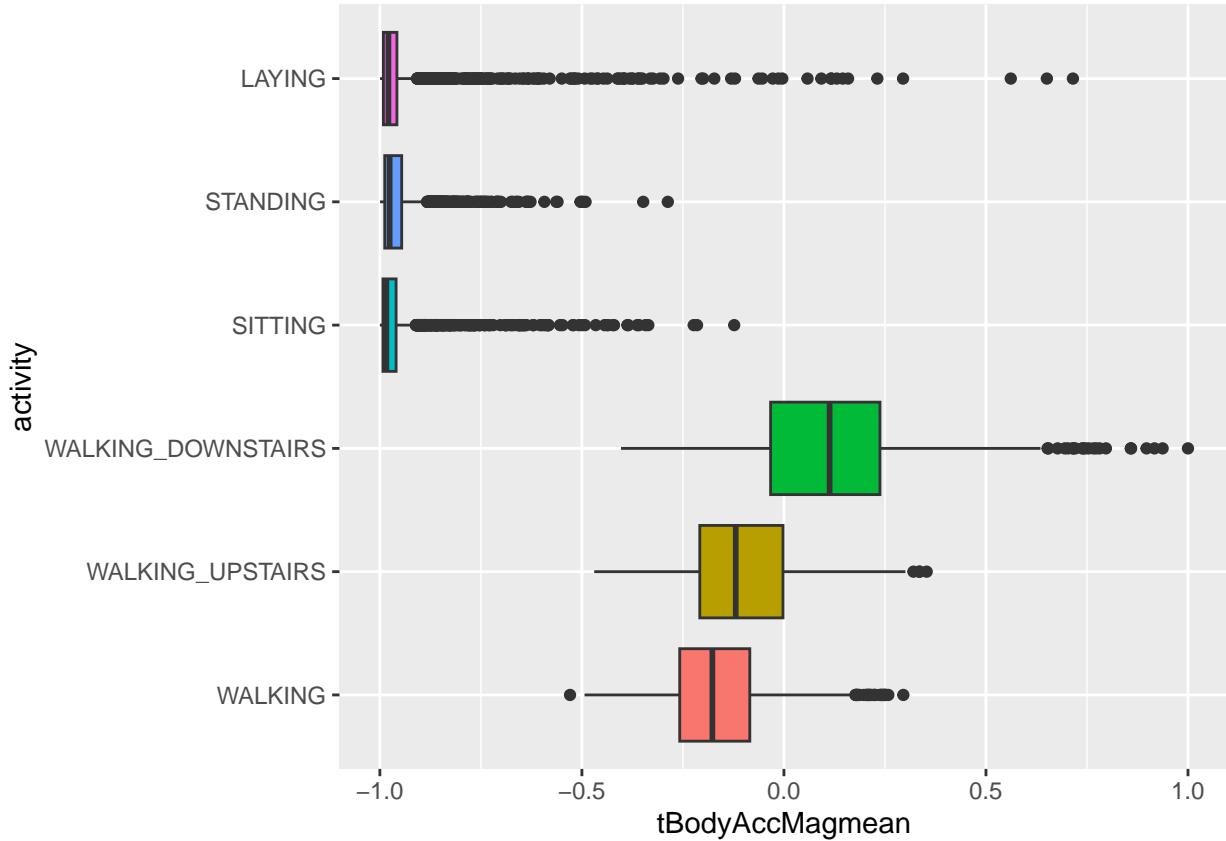
using the mantra from the above graph we understood:

The Activities are categorized into:

- stationary activities (such as laying, sitting, and standing) and
- moving activities (such as walking, walking downstairs, and walking upstairs).

Let's see the mean of triaxial magnitude of body acceleration signals (tBodyAccMagmean).

```
ggplot(final_data,
       aes(y = activity, x = tBodyAccMagmean, group = activity, fill = activity)) +
  geom_boxplot(show.legend = FALSE)
```



We can see a clear distinction between the two groups:

stationary activities have very small body movements compared to those of moving activities.

-if  $t\text{BodyAccMagmean} > -0.5$ , then the activity will probably be either walking, walking upstairs, or walking downstairs.

-if  $t\text{BodyAccMagmean} < -0.5$ , then the activity is most probably either laying, standing, or sitting.

Now, there should be also a distinction between laying and other activities in terms of phone orientation.

While laying, unlike other activities, people tend to have the phone fairly horizontal on their waist.

So, let's see whether this hypothesis is true by comparing the angle between X, Y, and Z axis to the mean of gravity acceleration signals on each axis ( $\text{angleXgravityMean}$ ,  $\text{angleYgravityMean}$ , and  $\text{angleZgravityMean}$ ).

“ $\text{angleXgravitymean}$ ”, “ $\text{angleYgravitymean}$ ”, and “ $\text{angleZgravitymean}$ ” are features that describe the mean angles between the body and gravity vectors, measured in radians. These angles can provide information on the orientation of the body during different activities, which are useful for identifying and classifying different types of movements.

—> Outlier detection using IQR

The below code creates a loop that plots boxplots for each of the columns in the final\_data. The loop here calculates the interquartile range (IQR) and the upper and lower whiskers for each column, identifies any outliers based on these values, and prints the number of outliers for each column.

```
colnames(final_data)[561] = "angleXgravityMean"
colnames(final_data)[562] = "angleZgravityMean"
colnames(final_data)[563] = "angleYgravityMean"
```

```

# Create a list of the columns to plot
coor <- c('angleXgravityMean', 'angleZgravityMean', 'angleYgravityMean')

# Loop over the columns and plot the boxplot
for (i in coor) {
  # Create the boxplot
  p <- ggplot(final_data, aes_string(y = 'activity', x = i, group = 'activity', fill = 'activity')) +
    geom_boxplot(show.legend = FALSE)

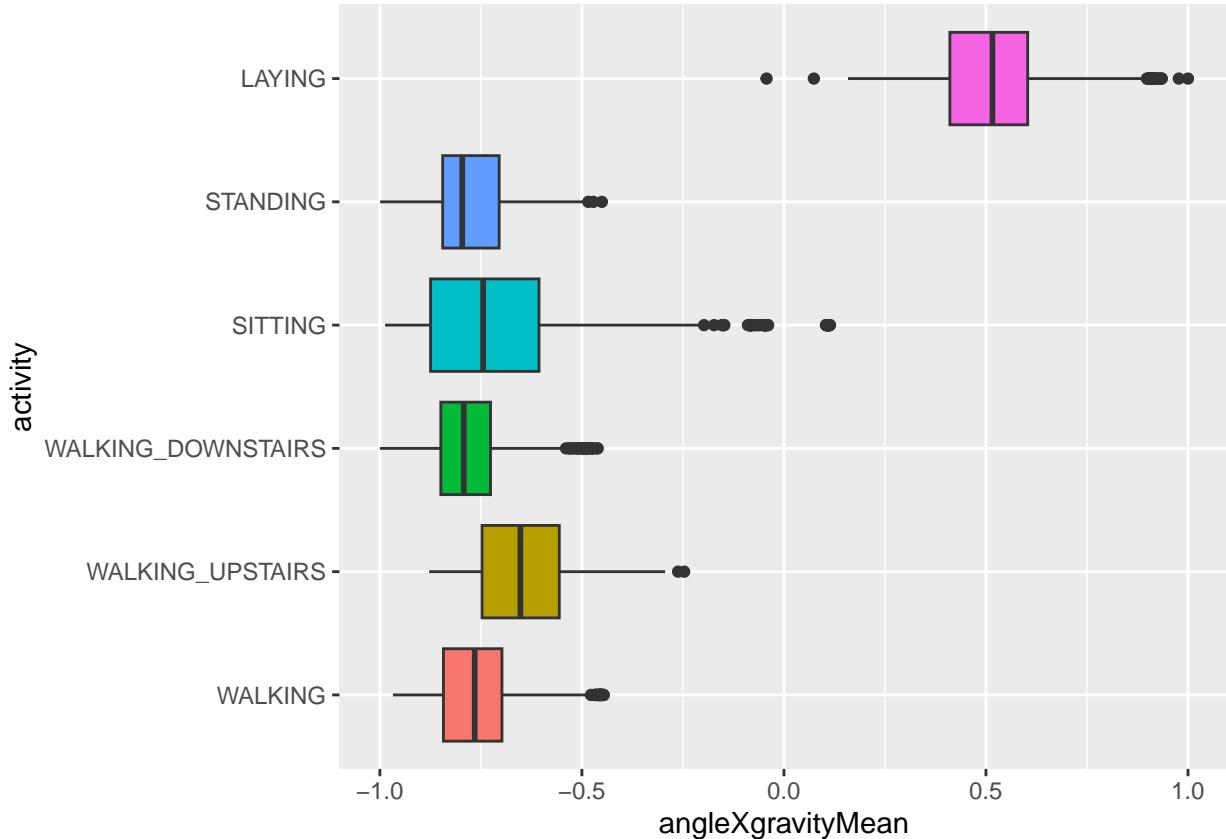
  # Calculate the IQR and whiskers
  q <- quantile(final_data[[i]], probs = c(0.25, 0.75), na.rm = TRUE)
  iqr <- q[2] - q[1]
  upper_whisker <- min(q[2] + 1.5*iqr, max(final_data[[i]], na.rm = TRUE))
  lower_whisker <- max(q[1] - 1.5*iqr, min(final_data[[i]], na.rm = TRUE))

  # Identify the outliers
  outliers <- final_data[final_data[[i]] > upper_whisker | final_data[[i]] < lower_whisker, ]

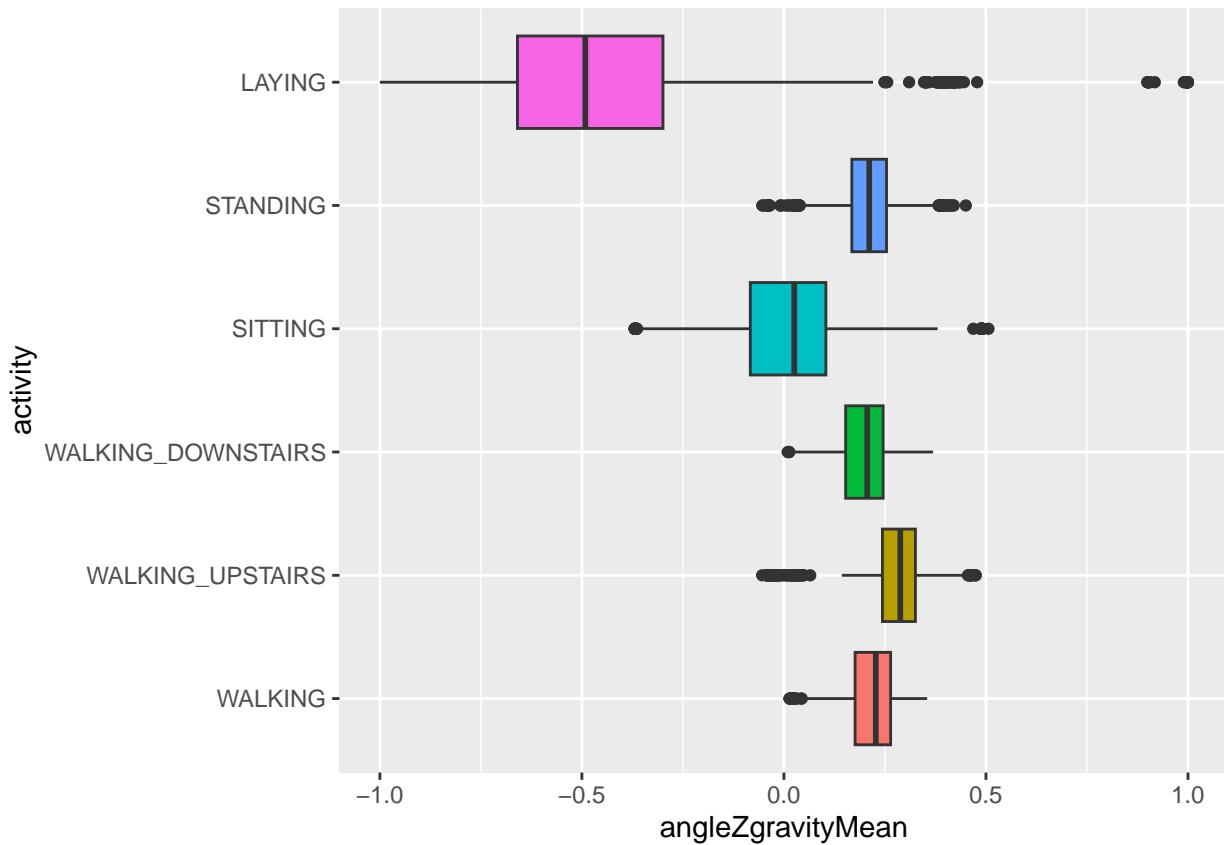
  # Print the plot and the number of outliers
  print(p)
  cat(paste0("Number of outliers in ", i, ": ", nrow(outliers), "\n"))
}

## Warning: `aes_string()` was deprecated in ggplot2 3.0.0.
## i Please use tidy evaluation ideoms with `aes()`

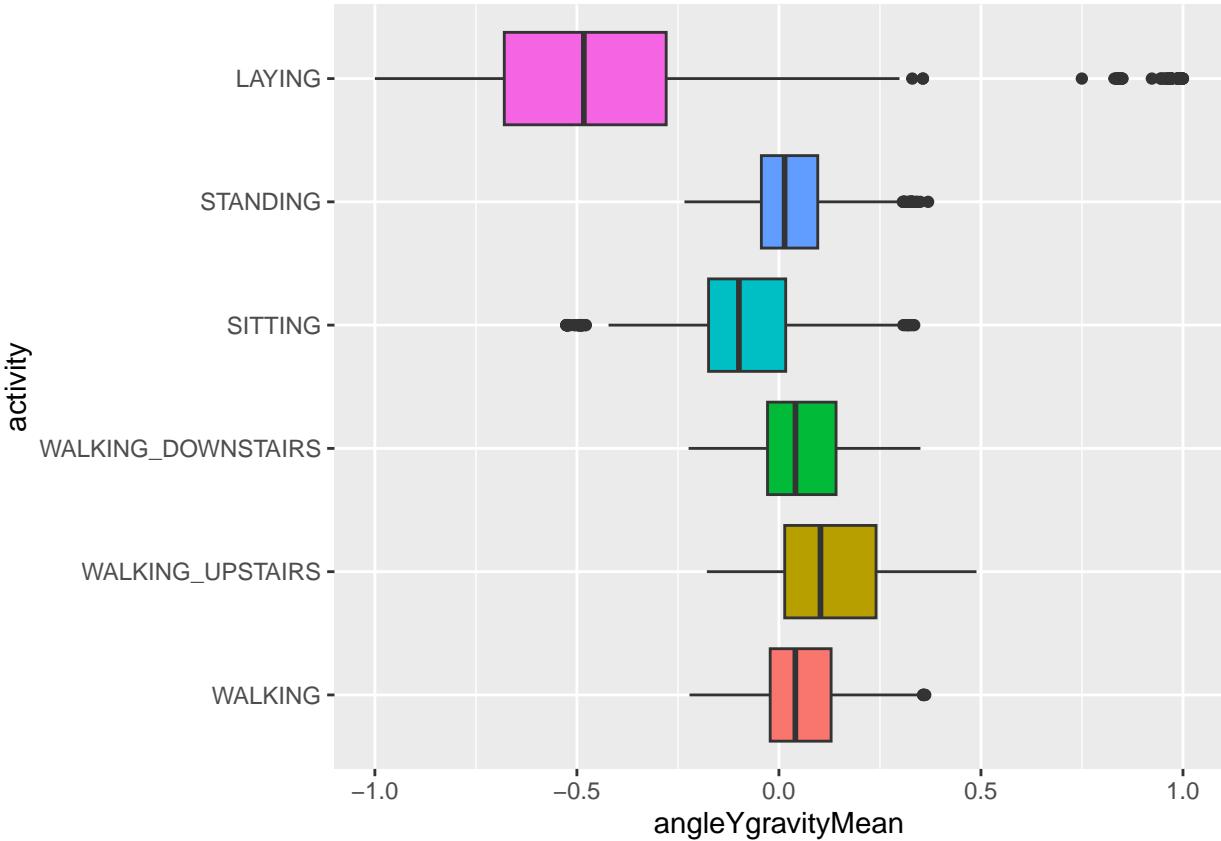
```



```
## Number of outliers in angleXgravityMean: 1979
```



```
## Number of outliers in angleZgravityMean: 1235
```



```
## Number of outliers in angleYgravityMean: 1055
```

It's apparent that:

Phone orientation while laying is different from phone orientation while doing other activities.

-if  $\text{angleXgravityMean} > 0$ , then the activity is most probably laying. Other activities otherwise.

-if  $\text{angleYgravityMean} < -0.25$  or  $\text{angleZgravityMean} < -0.25$ , then the activity will probably be laying. Other activities otherwise.

we can predict laying activity with minimum error only by using  $\text{angleXgravityMean}$ , or probably by other gravity-related features to the X axis.

As can be seen, all activities can easily be separated except for Standing and Sitting. This makes sense, since Standing and Sitting don't have much difference in terms of phone orientation.

Also we can see the no. of outliers for the respective columns. Now lets dig a bit deeper into it by summarizing the data.

```
# Calculate summary statistics for the variables
summary_data <- summary(final_data[, coor])

# Print the summary statistics
print(summary_data)
```

```
##   angleXgravityMean   angleZgravityMean   angleYgravityMean
##   Min.    :-1.00000   Min.    :-1.0000000   Min.    :-1.0000000
```

```

## 1st Qu.:-0.8173 1st Qu.: 0.002151 1st Qu.:-0.131880
## Median :-0.7156 Median : 0.182029 Median :-0.003882
## Mean : -0.4965 Mean : 0.063255 Mean : -0.054284
## 3rd Qu.:-0.5215 3rd Qu.: 0.250790 3rd Qu.: 0.102970
## Max. : 1.00000 Max. : 1.000000 Max. : 1.000000

```

These are summary statistics for three variables: angleXgravityMean, angleZgravityMean, and angleYgravityMean.

For angleXgravityMean, the minimum value is -1.0 and the maximum value is 1.0. The first quartile is -0.8173, the median is -0.7156, the mean is -0.4965, and the third quartile is -0.5215.

For angleZgravityMean, the minimum value is -1.0 and the maximum value is 1.0. The first quartile is 0.002151, the median is 0.182029, the mean is 0.063255, and the third quartile is 0.250790.

For angleYgravityMean, the minimum value is -1.0 and the maximum value is 1.0. The first quartile is -0.131880, the median is -0.003882, the mean is -0.054284, and the third quartile is 0.102970.

These statistics provide information about the distribution of the data for each variable. For example, the mean values for all three variables are negative, indicating that the data may be skewed to the left. Additionally, the range of values for all three variables is from -1 to 1, indicating that the data is normalized or standardized.

-> Removing outliers

As we have already seen that there are outliers present, so now lets try removing them and see how the dataset is being effected.

```

# Create a list of the columns to investigate
cols <- c('angleXgravityMean', 'angleZgravityMean', 'angleYgravityMean')

# Loop over the columns and remove outliers
for (i in cols) {
  # Calculate summary statistics
  stats <- summary(final_data[[i]])
  q1 <- stats[2]
  q3 <- stats[5]
  iqr <- q3 - q1
  lower_whisker <- q1 - 1.5 * iqr
  upper_whisker <- q3 + 1.5 * iqr

  # Identify outliers
  outliers <- final_data[final_data[[i]] < lower_whisker | final_data[[i]] > upper_whisker, ]

  if (nrow(outliers) > 0) {
    # Remove outliers
    final_data <- final_data[!final_data[[i]] %in% outliers[[i]], ]
  }
}

# Print the dimensions of the final data without outliers
cat("Dimensions of final data without outliers: ", dim(final_data), "\n")

## Dimensions of final data without outliers: 7789 566

```

```

table(final_data$activity)

##
##          WALKING    WALKING_UPSTAIRS WALKING_DOWNSTAIRS      SITTING
##          1722                1481           1406             1274
##          STANDING            LAYING
##          1906                  0

```

After removing the outliers we found that activity laying has been completely removed from our dataset. we can predict laying activity with minimum error only by using angleXgravityMean, or probably by other gravity-related features to the X axis.

The “LAYING” activity in the dataset is important because it represents a distinct and important class of activities that the subjects perform. In particular, the LAYING activity is often considered a baseline or reference activity, as it represents a period of minimal movement and low activity.

While the LAYING activity may not be as physically demanding as other activities in the dataset, it is an important reference point and provides valuable information for developing accurate models for activity recognition.

Hence revering back to the orginal data set.

```

#reverting back to main data set.
final_data <- rbind(train,test)
final_data$activity<-factor(final_data$activity)

levels(final_data$activity) <- c("WALKING", "WALKING_UPSTAIRS", "WALKING_DOWNSTAIRS", "SITTING", "STANDIN
dim(final_data)

## [1] 10299   566

```

```

table(final_data$activity)

##
##          WALKING    WALKING_UPSTAIRS WALKING_DOWNSTAIRS      SITTING
##          1722                1544           1406             1777
##          STANDING            LAYING
##          1906                  1944

```

—> Multidimensional Data Analysis:

we are performing Multidimensional data analysis for our datasets as it has multiple variables or features that are measured for each observation. This is to identify patterns and relationships among the different variables in the dataset. Which helps us better understand how different variables contribute to the classification of human activities and improve the accuracy of models.

Hierarchical Clustering

Hierarchical clustering is a useful technique for analyzing datasets because it can help to identify groups of similar activities or users based on their feature values. Hierarchical clustering works by recursively partitioning the dataset into clusters based on their similarity, starting with each observation as a separate cluster and merging them based on their similarity until all observations belong to a single cluster. The resulting dendrogram can be used to visualize the relationships among the different clusters and identify groups of

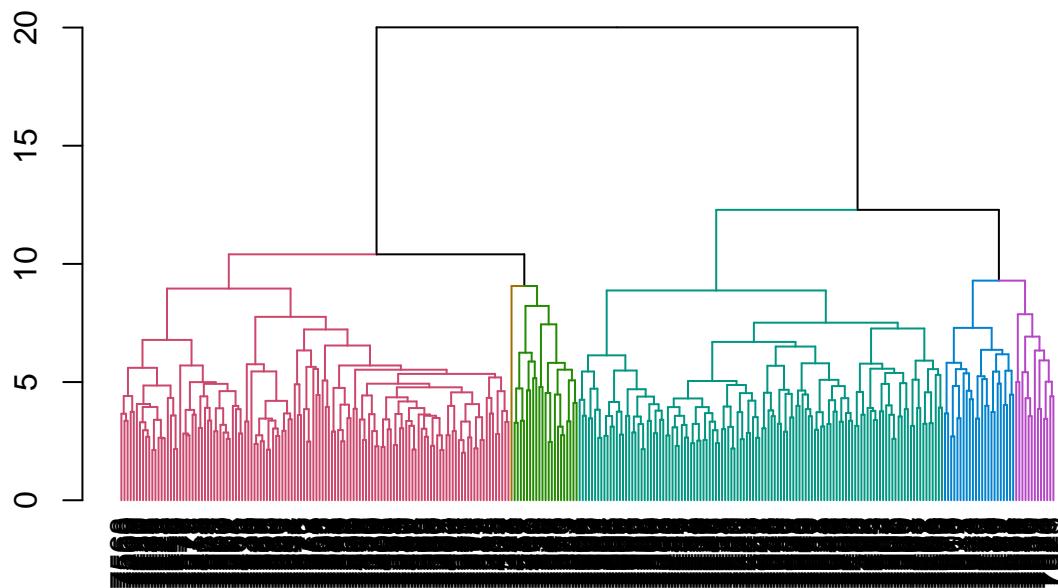
similar activities or users. One of the benefits of hierarchical clustering is that it can identify activity patterns that may not be immediately apparent from simple descriptive statistics or visualizations. It can provide insights into the relationships among different activities and users. It can help to identify activity patterns and anomalies that may be missed by simpler techniques and improve the accuracy and interpretability.

Performing hierarchical clustering on the data subset for subject 1. The distance matrix is computed for all observations for subject 1, and then the `hclust()` function is used to perform hierarchical clustering using the complete linkage method. The resulting dendrogram is then plotted, with branches colored based on 6 clusters.

```
sub1 <- subset(final_data, subjectID == 2)
table(sub1$activity)
```

```
##
##          WALKING   WALKING_UPSTAIRS WALKING_DOWNSTAIRS      SITTING
##          59                  48                 47             46
##          STANDING           LAYING
##          54                  48
```

```
#colnames(sub1)
sub1=sub1[,-c(1,2,564:566)]
distanceMatrix <- dist(sub1,method = 'euclidean')
hclustering <- hclust(distanceMatrix,method='complete')
dend <- as.dendrogram(hclustering)
dend <- color_branches(dend, k = 6)
plot(dend)
```



```
cut_avg <- cutree(hclustering, k = 6)
table(cut_avg)
```

```
## cut_avg
##   1   2   3   4   5   6
## 21 126 118 23  1  13
```

The following is the result of Hierarchical Clustering using Euclidean distance and complete distance from Subject 1's data:

The above dendrogram created by using complete linkage method where we understood that there are 6 clusters denoting 6 different activities. The clustering algorithm has identified 6 distinct groups of data points based on their pairwise distances, which could represent different behaviors, movements, or activities performed by subjectID 1. It is not possible to draw a conclusion because the data set is complex and variables might be highly correlated as a result hierarchical clustering is not so informative. Lets try another type of clustering.

—>k-means

k-means clustering is a more commonly used method for clustering large datasets like the HAR dataset, as it is computationally efficient and can handle a large number of observations and variables. K-means clustering is also easier to interpret, as it produces hard clusters that assign each observation to exactly one cluster. However, k-means clustering assumes that the clusters are spherical, equally sized, and that the variance is the same across all dimensions.

```
# Load the HAR dataset

library(foreign)

data<- rbind(train,test)

# Extract the predictor variables

predictors <- data[, 3:562]

# Set the number of clusters

k <- 6

# Perform k-means clustering

set.seed(123)

kmeans_model <- kmeans(predictors, k)

#kmeans function to perform k-means clustering, passing in the predictor variables and the number of clusters

#We also set a random seed for reproducibility.

# Create a scatterplot matrix of the first six variables, colored by cluster

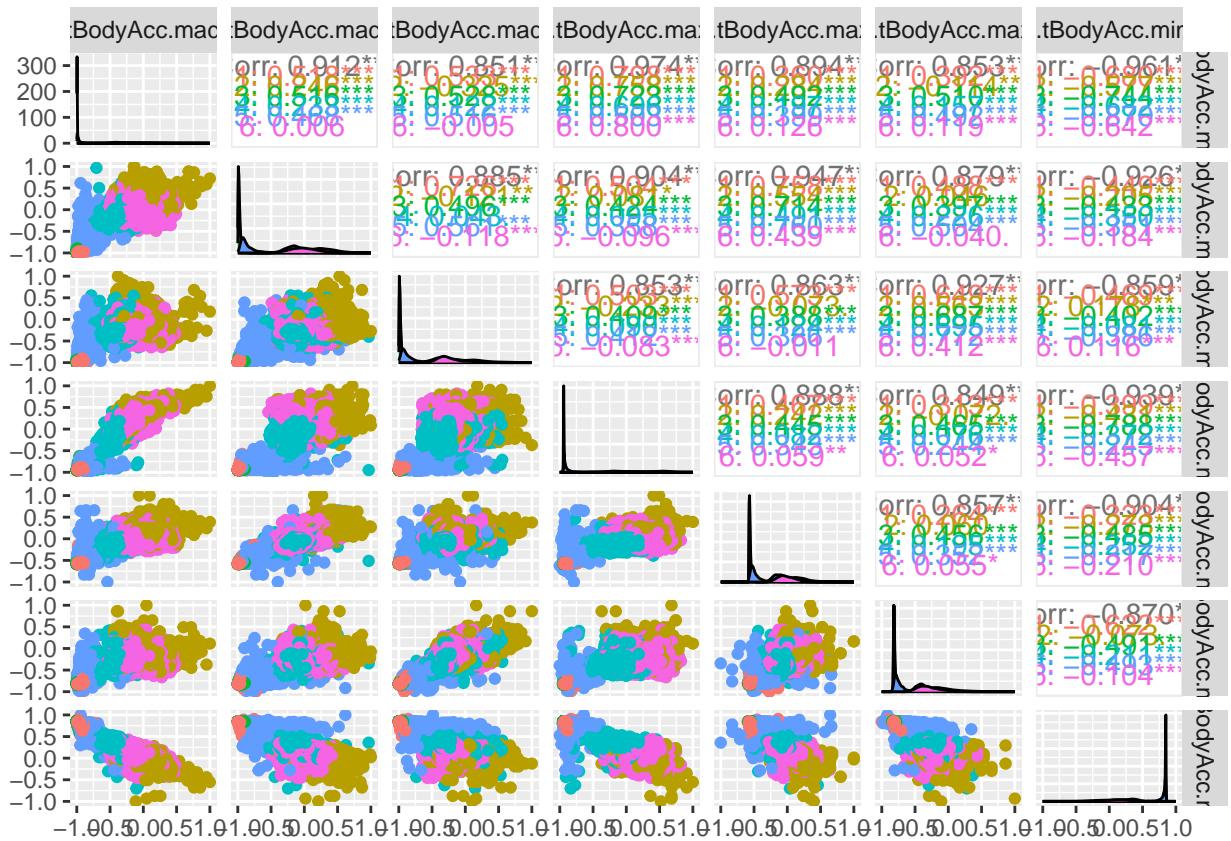
library(GGally)
```

```

## Registered S3 method overwritten by 'GGally':
##   method from
##   +.gg   ggplot2

ggpairs(data[ ,c(9:15)], aes(color = as.factor(kmeans_model$cluster)))

```



```

# Print the size of each cluster
cat("Cluster sizes: ", kmeans_model$size, "\n")

```

```

## Cluster sizes:  1797 687 2499 2067 1316 1933

```

We then create several plots to visualize the clustering results. The first plot is a scatterplot matrix of the first six variables, colored by cluster. This can help us see how well the clusters are separated in the data space. We use the ggpairs function from the GGally package to create this plot.

In the above scatterplot matrix the diagonal contains density plot, below the diagonal are the scatterplots and the above are correlation among the variables. Each color represents each activity. The matrix plotted for subject, activity, AngleXgravityMean, AngleYgravityMean and AngleZgravityMean

It is hard to plot the matrix for all the variables as our data set to understand each and every column.

->correlation:

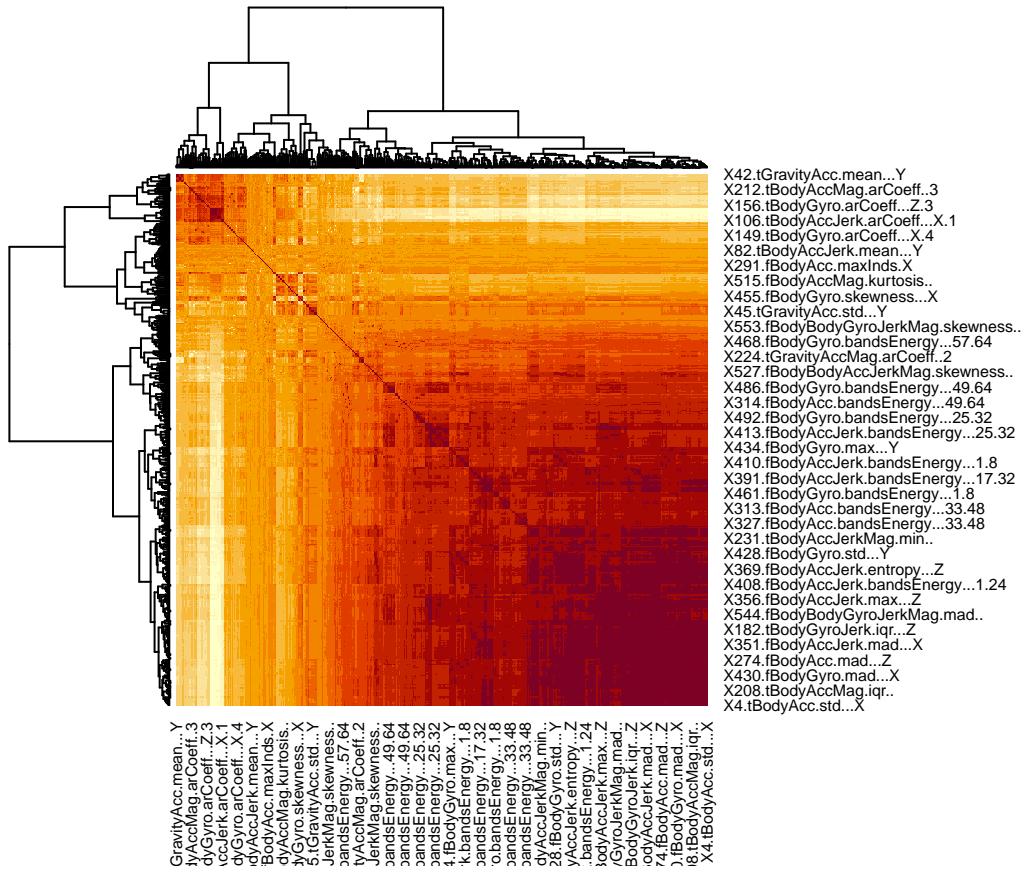
Now lets understand the relationship between two or more variables in the dataset. Lets find the answers for these questions. 1. Is there a relationship between two variables? 2. If there is a relationship, how strong is it?

We built a correlation matrix for all 561 variables in one, got to identify any apparent patterns in the relationships.

We see that most of these features are highly correlated with each other.

```
# correlation
training_data <- cbind(train_X, train_Y)
test_data <- cbind(test_X, test_Y)
test_data$activity<-as.factor(test_data$activity)
training_data$activity <- as.factor(training_data$activity)

final_data<-final_data[,-c(1,564:566)]
mydata.cor = cor(final_data[,-1])
heatmap(x = mydata.cor, symm = TRUE)
```



From the above heat map we can observe that there are multiple variables that are highly correlated(red) and some variables are partially correlated(orange) and some that are not at all correlated(white) as our data set is high dimensional, it make sense to concentrate on uncorrelated data for the following reasons:

When building predictive models or making data-driven decisions.

When two variables are highly correlated, it means that they tend to move in the same direction (positive correlation) or in opposite directions (negative correlation).

This can lead to issues when building predictive models, as it can result in a phenomenon called multicollinearity, where the model may overemphasize the importance of one variable and underestimate the importance of another. In contrast, uncorrelated data contains no relationship between the variables, which can be beneficial when building models or making data-driven decisions. When the variables are uncorrelated, it allows each variable to provide unique information and reduces the risk of multicollinearity.

A correlation matrix is used to show the strength and direction of the relationships between pairs of variables in a dataset, allowing for identification of patterns and dependencies among variables.

Let find the uncorrelated data.

```
# Calculate the correlation matrix

cor_mat <- cor(final_data[,-1] [, -ncol(final_data[,-1])])

# Select the variables with a correlation coefficient above a certain threshold (e.g., 0.7)
high_cor_vars <- findCorrelation(cor_mat, cutoff = 0.8)

# Remove the highly correlated variables from the data set
training_data_no_cor <- training_data[, -high_cor_vars]

# Remove the highly correlated variables from the data set
test_data_no_cor <- test_data[, -high_cor_vars]
```

Firstly we split the data into training and test data and then performed correlation matrix on the final data and found the highly correlated features. Then removed these features from both the training and test data creating new test and train data sets which are uncorrelated.

```
dim(training_data_no_cor)

## [1] 7352 174

names(training_data_no_cor)

## [1] "X1.tBodyAcc.mean...X"
## [2] "X2.tBodyAcc.mean...Y"
## [3] "X3.tBodyAcc.mean...Z"
## [4] "X23.tBodyAcc.entropy...X"
## [5] "X25.tBodyAcc.entropy...Z"
## [6] "X28.tBodyAcc.arCoeff...X.3"
## [7] "X29.tBodyAcc.arCoeff...X.4"
## [8] "X31.tBodyAcc.arCoeff...Y.2"
## [9] "X32.tBodyAcc.arCoeff...Y.3"
## [10] "X33.tBodyAcc.arCoeff...Y.4"
## [11] "X35.tBodyAcc.arCoeff...Z.2"
## [12] "X36.tBodyAcc.arCoeff...Z.3"
## [13] "X37.tBodyAcc.arCoeff...Z.4"
## [14] "X38.tBodyAcc.correlation...X.Y"
## [15] "X39.tBodyAcc.correlation...X.Z"
## [16] "X40.tBodyAcc.correlation...Y.Z"
## [17] "X44.tGravityAcc.std...X"
## [18] "X52.tGravityAcc.max...Z"
## [19] "X56.tGravityAcc.sma.."
## [20] "X58.tGravityAcc.energy...Y"
## [21] "X59.tGravityAcc.energy...Z"
## [22] "X61.tGravityAcc.iqr...Y"
## [23] "X62.tGravityAcc.iqr...Z"
## [24] "X63.tGravityAcc.entropy...X"
## [25] "X64.tGravityAcc.entropy...Y"
```

```

## [26] "X65.tGravityAcc.entropy...Z"
## [27] "X66.tGravityAcc.arCoeff...X.1"
## [28] "X73.tGravityAcc.arCoeff...Y.4"
## [29] "X77.tGravityAcc.arCoeff...Z.4"
## [30] "X78.tGravityAcc.correlation...X.Y"
## [31] "X79.tGravityAcc.correlation...X.Z"
## [32] "X80.tGravityAcc.correlation...Y.Z"
## [33] "X81.tBodyAccJerk.mean...X"
## [34] "X82.tBodyAccJerk.mean...Y"
## [35] "X83.tBodyAccJerk.mean...Z"
## [36] "X107.tBodyAccJerk.arCoeff...X.2"
## [37] "X108.tBodyAccJerk.arCoeff...X.3"
## [38] "X109.tBodyAccJerk.arCoeff...X.4"
## [39] "X112.tBodyAccJerk.arCoeff...Y.3"
## [40] "X113.tBodyAccJerk.arCoeff...Y.4"
## [41] "X115.tBodyAccJerk.arCoeff...Z.2"
## [42] "X116.tBodyAccJerk.arCoeff...Z.3"
## [43] "X117.tBodyAccJerk.arCoeff...Z.4"
## [44] "X118.tBodyAccJerk.correlation...X.Y"
## [45] "X119.tBodyAccJerk.correlation...X.Z"
## [46] "X120.tBodyAccJerk.correlation...Y.Z"
## [47] "X121.tBodyGyro.mean...X"
## [48] "X122.tBodyGyro.mean...Y"
## [49] "X123.tBodyGyro.mean...Z"
## [50] "X143.tBodyGyro.entropy...X"
## [51] "X144.tBodyGyro.entropy...Y"
## [52] "X145.tBodyGyro.entropy...Z"
## [53] "X148.tBodyGyro.arCoeff...X.3"
## [54] "X149.tBodyGyro.arCoeff...X.4"
## [55] "X150.tBodyGyro.arCoeff...Y.1"
## [56] "X152.tBodyGyro.arCoeff...Y.3"
## [57] "X154.tBodyGyro.arCoeff...Z.1"
## [58] "X156.tBodyGyro.arCoeff...Z.3"
## [59] "X157.tBodyGyro.arCoeff...Z.4"
## [60] "X158.tBodyGyro.correlation...X.Y"
## [61] "X159.tBodyGyro.correlation...X.Z"
## [62] "X160.tBodyGyro.correlation...Y.Z"
## [63] "X161.tBodyGyroJerk.mean...X"
## [64] "X162.tBodyGyroJerk.mean...Y"
## [65] "X163.tBodyGyroJerk.mean...Z"
## [66] "X186.tBodyGyroJerk.arCoeff...X.1"
## [67] "X187.tBodyGyroJerk.arCoeff...X.2"
## [68] "X188.tBodyGyroJerk.arCoeff...X.3"
## [69] "X189.tBodyGyroJerk.arCoeff...X.4"
## [70] "X191.tBodyGyroJerk.arCoeff...Y.2"
## [71] "X192.tBodyGyroJerk.arCoeff...Y.3"
## [72] "X193.tBodyGyroJerk.arCoeff...Y.4"
## [73] "X196.tBodyGyroJerk.arCoeff...Z.3"
## [74] "X197.tBodyGyroJerk.arCoeff...Z.4"
## [75] "X198.tBodyGyroJerk.correlation...X.Y"
## [76] "X199.tBodyGyroJerk.correlation...X.Z"
## [77] "X200.tBodyGyroJerk.correlation...Y.Z"
## [78] "X218.tGravityAccMag.min... "
## [79] "X223.tGravityAccMag.arCoeff..1"

```

```

## [80] "X226.tGravityAccMag.arCoeff..4"
## [81] "X231.tBodyAccJerkMag.min.."
## [82] "X237.tBodyAccJerkMag.arCoeff..2"
## [83] "X238.tBodyAccJerkMag.arCoeff..3"
## [84] "X239.tBodyAccJerkMag.arCoeff..4"
## [85] "X244.tBodyGyroMag.min.."
## [86] "X248.tBodyGyroMag.entropy.."
## [87] "X250.tBodyGyroMag.arCoeff..2"
## [88] "X252.tBodyGyroMag.arCoeff..4"
## [89] "X257.tBodyGyroJerkMag.min.."
## [90] "X262.tBodyGyroJerkMag.arCoeff..1"
## [91] "X263.tBodyGyroJerkMag.arCoeff..2"
## [92] "X264.tBodyGyroJerkMag.arCoeff..3"
## [93] "X265.tBodyGyroJerkMag.arCoeff..4"
## [94] "X278.fBodyAcc.min...X"
## [95] "X279.fBodyAcc.min...Y"
## [96] "X280.fBodyAcc.min...Z"
## [97] "X291.fBodyAcc.maxInds.X"
## [98] "X292.fBodyAcc.maxInds.Y"
## [99] "X293.fBodyAcc.maxInds.Z"
## [100] "X294.fBodyAcc.meanFreq...X"
## [101] "X295.fBodyAcc.meanFreq...Y"
## [102] "X296.fBodyAcc.meanFreq...Z"
## [103] "X298.fBodyAcc.kurtosis...X"
## [104] "X299.fBodyAcc.skewness...Y"
## [105] "X301.fBodyAcc.skewness...Z"
## [106] "X304.fBodyAcc.bandsEnergy...9.16"
## [107] "X310.fBodyAcc.bandsEnergy...57.64"
## [108] "X324.fBodyAcc.bandsEnergy...57.64"
## [109] "X331.fBodyAcc.bandsEnergy...1.8"
## [110] "X338.fBodyAcc.bandsEnergy...57.64"
## [111] "X357.fBodyAccJerk.min...X"
## [112] "X358.fBodyAccJerk.min...Y"
## [113] "X359.fBodyAccJerk.min...Z"
## [114] "X370.fBodyAccJerk.maxInds.X"
## [115] "X371.fBodyAccJerk.maxInds.Y"
## [116] "X372.fBodyAccJerk.maxInds.Z"
## [117] "X377.fBodyAccJerk.kurtosis...X"
## [118] "X379.fBodyAccJerk.kurtosis...Y"
## [119] "X381.fBodyAccJerk.kurtosis...Z"
## [120] "X385.fBodyAccJerk.bandsEnergy...25.32"
## [121] "X389.fBodyAccJerk.bandsEnergy...57.64"
## [122] "X396.fBodyAccJerk.bandsEnergy...1.8"
## [123] "X399.fBodyAccJerk.bandsEnergy...25.32"
## [124] "X403.fBodyAccJerk.bandsEnergy...57.64"
## [125] "X414.fBodyAccJerk.bandsEnergy...33.40"
## [126] "X417.fBodyAccJerk.bandsEnergy...57.64"
## [127] "X436.fBodyGyro.min...X"
## [128] "X437.fBodyGyro.min...Y"
## [129] "X438.fBodyGyro.min...Z"
## [130] "X449.fBodyGyro.maxInds.X"
## [131] "X450.fBodyGyro.maxInds.Y"
## [132] "X451.fBodyGyro.maxInds.Z"
## [133] "X452.fBodyGyro.meanFreq...X"

```

```

## [134] "X453.fBodyGyro.meanFreq...Y"
## [135] "X454.fBodyGyro.meanFreq...Z"
## [136] "X456.fBodyGyro.kurtosis...X"
## [137] "X458.fBodyGyro.kurtosis...Y"
## [138] "X460.fBodyGyro.kurtosis...Z"
## [139] "X461.fBodyGyro.bandsEnergy...1.8"
## [140] "X462.fBodyGyro.bandsEnergy...9.16"
## [141] "X463.fBodyGyro.bandsEnergy...17.24"
## [142] "X465.fBodyGyro.bandsEnergy...33.40"
## [143] "X468.fBodyGyro.bandsEnergy...57.64"
## [144] "X475.fBodyGyro.bandsEnergy...1.8"
## [145] "X476.fBodyGyro.bandsEnergy...9.16"
## [146] "X477.fBodyGyro.bandsEnergy...17.24"
## [147] "X479.fBodyGyro.bandsEnergy...33.40"
## [148] "X482.fBodyGyro.bandsEnergy...57.64"
## [149] "X489.fBodyGyro.bandsEnergy...1.8"
## [150] "X490.fBodyGyro.bandsEnergy...9.16"
## [151] "X491.fBodyGyro.bandsEnergy...17.24"
## [152] "X493.fBodyGyro.bandsEnergy...33.40"
## [153] "X496.fBodyGyro.bandsEnergy...57.64"
## [154] "X507.fBodyAccMag.min.."
## [155] "X512.fBodyAccMag.maxInds"
## [156] "X513.fBodyAccMag.meanFreq.."
## [157] "X514.fBodyAccMag.skewness.."
## [158] "X520.fBodyBodyAccJerkMag.min.."
## [159] "X525.fBodyBodyAccJerkMag.maxInds"
## [160] "X526.fBodyBodyAccJerkMag.meanFreq.."
## [161] "X528.fBodyBodyAccJerkMag.kurtosis.."
## [162] "X533.fBodyBodyGyroMag.min.."
## [163] "X538.fBodyBodyGyroMag.maxInds"
## [164] "X541.fBodyBodyGyroMag.kurtosis.."
## [165] "X546.fBodyBodyGyroJerkMag.min.."
## [166] "X551.fBodyBodyGyroJerkMag.maxInds"
## [167] "X552.fBodyBodyGyroJerkMag.meanFreq.."
## [168] "X554.fBodyBodyGyroJerkMag.kurtosis.."
## [169] "X555.angle.tBodyAccMean.gravity."
## [170] "X556.angle.tBodyAccJerkMean..gravityMean."
## [171] "X557.angle.tBodyGyroMean.gravityMean."
## [172] "X558.angle.tBodyGyroJerkMean.gravityMean."
## [173] "X561.angle.Z.gravityMean."
## [174] "activity"

```

We found that there are 174 uncorrelated features in the data set. And names of those features are listed above.

Lets now processed forward with the lasso regression on uncorrelated data.

Performing Lasso regression on uncorrelated data is useful, as it allows for identification and selection of important variables in a predictive model, while minimizing the risk of multicollinearity.

```

set.seed(100)
x <- model.matrix(activity ~ ., data=training_data_no_cor) [, -1]
y=as.factor(training_data_no_cor$activity)

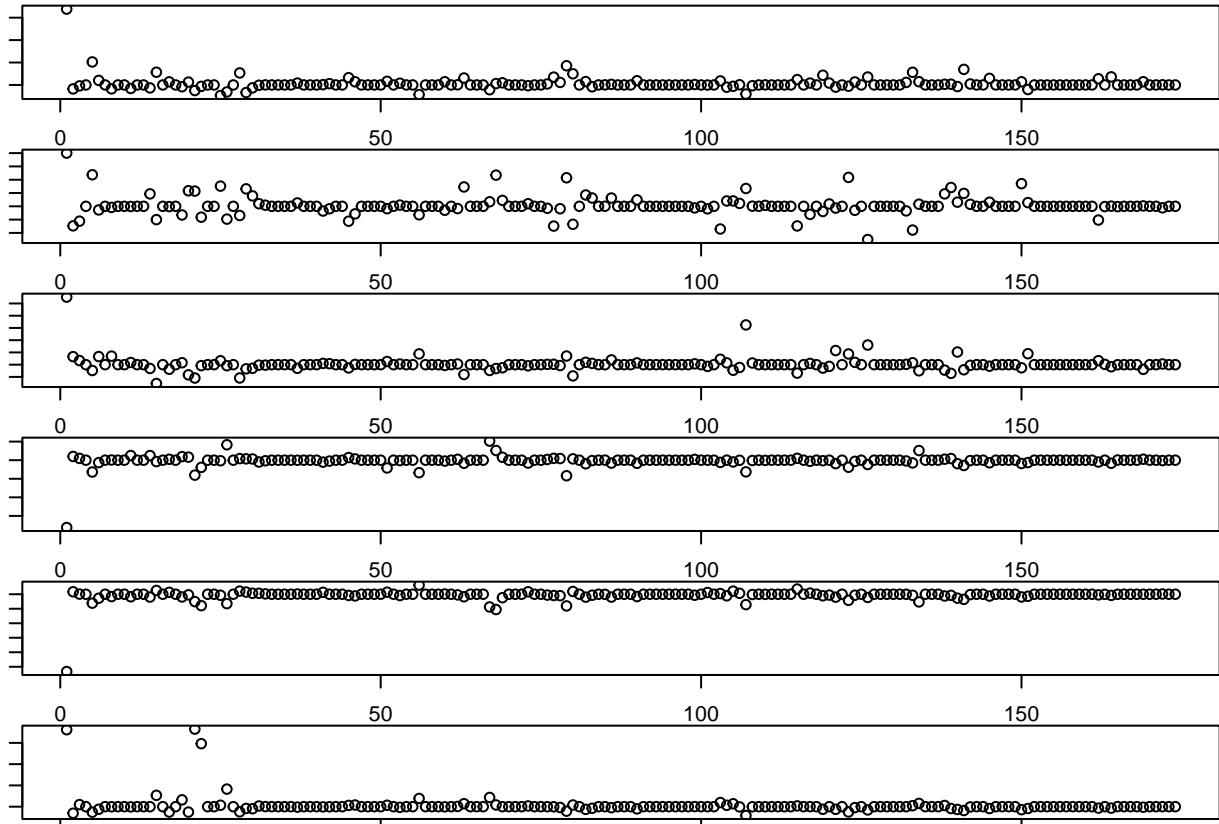
grid=10^seq(10,-2,length=100)

```

```

cvfit = cv.glmnet(x,y,type.measure="class",alpha=1,family="multinomial", lambda=grid, type.multinomial=
coef=coef(cvfit, s = "lambda.min")
par(mar=c(1,1,1,1))
par(mfrow=c(6,1))
plot(coef[[1]], xlab="WALKING", ylab= "Predictors")
plot(coef[[2]], xlab="WALKING_UPSTAIRS", ylab="Predictors")
plot(coef[[3]], xlab= "WALKING_DOWNSTAIRS", ylab= "Predictors")
plot(coef[[4]], xlab= "SITTING", ylab= "Predictors")
plot(coef[[5]], xlab="STANDING", ylab= "Predictors")
plot(coef[[6]], xlab="LAYING", ylab="Predictors")

```



Plots the coefficients for each of the six activity types in the HAR dataset: WALKING, WALKING\_UPSTAIRS, WALKING\_DOWNSTAIRS, SITTING, STANDING, and LAYING.

LASSO is a regression method that performs variable selection by shrinking the coefficients of insignificant independent variables to exactly zero. LASSO can be applied to a classification problem to obtain the most important predictors for a given class. For implementing Multinomial LASSO for the HAR dataset, classes should be coded as binary variable. Using the ‘glmnet’ package and selecting ‘multinomial’ family significant predictors related to each activity can be found out. R outputs a list which contains 6 sparse matrices, each comprising of the important predictors for the 6 levels of activities which are coded as binary numbers from 0 to 5. Similarly, ridge multinomial regression is also performed to identify important features.

The following plots show the magnitude of coefficients for all 174 predictors plotted against each activity. Activity Laying has the most number of uninfluential predictors. From the plots we can see that LASSO efficiently shrinks many of the predictors to zero while Ridge regression fails to do so. It is difficult to extract a subset of all the important predictors making sense for the overall classification model. Thus, we have not used LASSO for feature selection of our classification model.

-> Classification

Here we will use Accuracy, kappa statistics and precision as performance measures.

As already discussed above about building an app for Healthcare Monitoring which uses sensors(accelerometer, gyroscope) to collect data from our mobile phone to predict Physical Activity accurately. we used the following measures(Accuracy, kappa statistics and precision ) to find the final model which best suits and predicts activity correctly.

Why are we using these statistical measures??

**Accuracy** is the percentage of correctly classified instances in the total number of instances. A high accuracy indicates that the model is correctly predicting most of the physical activity types, which is important in healthcare monitoring apps. However, accuracy alone is not enough to assess the effectiveness of a model.

**Kappa statistics**, also known as Cohen's kappa, is a measure of the agreement between the predicted and actual class labels, corrected for chance. Kappa statistics take into account the class distribution in the dataset and provide a more accurate assessment of the model's performance. A high kappa value indicates a good agreement between the predicted and actual class labels.

**Precision** is the proportion of true positives out of the total number of predicted positives. In healthcare monitoring apps, precision is particularly important because false positives can lead to incorrect diagnoses, unnecessary treatments, and patient anxiety. A high precision indicates that the model is correctly identifying physical activity types and minimizing false positives.

In summary, accuracy, kappa statistics, and precision are all important metrics to consider when evaluating the performance of a model in healthcare monitoring apps. While accuracy provides a general idea of the model's performance, kappa statistics and precision take into account the class distribution and the cost of false positives, respectively, and are therefore more informative metrics in assessing the model's effectiveness in a real-world scenario.

->SVM(Support Vector Machine):

Support Vector Machines (SVM) is a popular classification algorithm that is often used to effectively handle high-dimensional data with a large number of features. Additionally, SVM can handle non-linearly separable data through the use of kernel functions.

```
#sum
library(e1071)

# Train the SVM model
svm_model <- svm(activity ~ ., data = training_data_no_cor, kernel = "linear")
svm_model

## 
## Call:
## svm(formula = activity ~ ., data = training_data_no_cor, kernel = "linear")
## 
## 
## Parameters:
##   SVM-Type:  C-classification
##   SVM-Kernel: linear
##   cost:    1
## 
## Number of Support Vectors:  1116
```

```

# Make predictions on the test set
svm_pred <- predict(svm_model, newdata = test_data_no_cor)

# Convert predicted and true labels to factors with the same levels
pred_svm <- factor(svm_pred, levels = levels(test_data_no_cor$activity))
y_test <- test_data_no_cor$activity
y_test <- factor(y_test, levels = levels(test_data_no_cor$activity))

# Calculate precision and recall for each class
precision_svm <- numeric(length(levels(test_data_no_cor$activity)))
recall_svm <- numeric(length(levels(test_data_no_cor$activity)))
for (i in seq_along(precision_svm)) {
  tp <- sum(pred_svm == levels(test_data_no_cor$activity)[i] & y_test == levels(test_data_no_cor$activity)[i])
  fp <- sum(pred_svm == levels(test_data_no_cor$activity)[i] & y_test != levels(test_data_no_cor$activity)[i])
  fn <- sum(pred_svm != levels(test_data_no_cor$activity)[i] & y_test == levels(test_data_no_cor$activity)[i])
  precision_svm[i] <- tp / (tp + fp)
  recall_svm[i] <- tp / (tp + fn)
}

# Calculate F1-score for each class
f1_svm <- 2 * precision_svm * recall_svm / (precision_svm + recall_svm)

# Calculate accuracy and Print the results

confusion_matrix_svm <- confusionMatrix(svm_pred, test_data_no_cor$activity)
print(confusion_matrix_svm)

## Confusion Matrix and Statistics
##
##          Reference
## Prediction 1 2 3 4 5 6
##       1 488 29 10 0 0 0
##       2 8 418 22 2 0 0
##       3 0 24 388 0 0 0
##       4 0 0 0 427 35 0
##       5 0 0 0 62 497 0
##       6 0 0 0 0 0 537
##
## Overall Statistics
##
##          Accuracy : 0.9348
##                 95% CI : (0.9253, 0.9435)
##      No Information Rate : 0.1822
##      P-Value [Acc > NIR] : < 2.2e-16
##
##          Kappa : 0.9217
##
## McNemar's Test P-Value : NA
##
## Statistics by Class:
##
##          Class: 1 Class: 2 Class: 3 Class: 4 Class: 5 Class: 6

```

## Sensitivity	0.9839	0.8875	0.9238	0.8697	0.9342	1.0000
## Specificity	0.9841	0.9871	0.9905	0.9857	0.9743	1.0000
## Pos Pred Value	0.9260	0.9289	0.9417	0.9242	0.8891	1.0000
## Neg Pred Value	0.9967	0.9788	0.9874	0.9742	0.9853	1.0000
## Prevalence	0.1683	0.1598	0.1425	0.1666	0.1805	0.1822
## Detection Rate	0.1656	0.1418	0.1317	0.1449	0.1686	0.1822
## Detection Prevalence	0.1788	0.1527	0.1398	0.1568	0.1897	0.1822
## Balanced Accuracy	0.9840	0.9373	0.9572	0.9277	0.9543	1.0000

The confusion matrix and statistics generated by evaluating the performance of a SVM model on a test dataset. The confusion matrix shows the number of correct and incorrect predictions for each class of the target variable, and the statistics provide further insights into the model's performance.

The confusion matrix shows that the model made 488 correct predictions for class 1, with 29 false positives and 0 false negatives. Similarly, the model made 418 correct predictions for class 2, with 8 false positives and 22 false negatives, and so on for the other classes.

The overall statistics show that the model achieved an accuracy of 0.9348, which is the proportion of correctly classified samples out of the total number of samples. The kappa statistic, which measures the agreement between the predicted and actual classes, was 0.9217.

The sensitivity and specificity values for each class indicate the model's ability to correctly identify positive and negative samples, respectively. The positive predictive value (PPV) and negative predictive value (NPV) indicate the proportion of true positive and true negative samples, respectively, out of all samples predicted as positive or negative.

The prevalence value for each class is the proportion of samples belonging to that class in the test dataset, while the detection rate is the proportion of true positive samples out of all actual positive samples. The detection prevalence is the proportion of samples predicted as positive out of all samples in the test dataset.

The balanced accuracy provides an overall measure of the model's performance, taking into account the sensitivity and specificity values for each class. In this case, the balanced accuracy was high for all classes, ranging from 0.9277 to 1.0000, indicating that the model performed well across all classes.

```
print(paste("Precision for each class:"))

## [1] "Precision for each class:

print(precision_svm)

## [1] 0.9259962 0.9288889 0.9417476 0.9242424 0.8890877 1.0000000
```

The precision values are as follows:

Class 1(Walking) precision: 0.9259962 Class 2(Walking upstairs) precision: 0.9288889 Class 3(Walking downstairs) precision: 0.9417476 Class 4(Sitting) precision: 0.9242424 Class 5(Standing) precision: 0.8890877 Class 6(Laying) precision: 1.0000000

Precision is a measure of the model's ability to correctly predict the positive instances out of the total instances predicted as positive for each class. For example, the precision for Class 1(Walking) indicates that 92.6% of the instances predicted as Class 1 are actually Class 1(Walking). Among all activities SVM can predict laying with 100% accurately.

—>Logistic regression

Logistic regression is a statistical method used to analyze the relationship between one or more independent variables and a binary dependent variable, and it can be applied to the Human Activity Recognition (HAR) dataset to predict the activity label based on sensor measurements.

```
#LR
library(glmnet)
x <- model.matrix(activity ~ ., data = training_data_no_cor)[, -1]
y <- as.factor(training_data_no_cor$activity)
fit_LR <- glmnet(x, y, family = "multinomial", alpha = 1)
```

```
fit_LR
```

```
##
## Call: glmnet(x = x, y = y, family = "multinomial", alpha = 1)
##
##          Df  %Dev  Lambda
## 1        0  0.00 0.32520
## 2        2  2.94 0.29630
## 3        2  6.35 0.27000
## 4        2  9.00 0.24600
## 5        2 11.16 0.22420
## 6        3 14.18 0.20430
## 7        3 16.84 0.18610
## 8        5 20.39 0.16960
## 9        5 24.70 0.15450
## 10       6 28.97 0.14080
## 11       11 33.18 0.12830
## 12       15 37.54 0.11690
## 13       15 41.35 0.10650
## 14       19 44.91 0.09704
## 15       22 48.23 0.08842
## 16       26 51.36 0.08056
## 17       27 54.32 0.07341
## 18       30 56.98 0.06688
## 19       35 59.43 0.06094
## 20       37 61.68 0.05553
## 21       39 63.81 0.05060
## 22       42 65.82 0.04610
## 23       42 67.66 0.04201
## 24       43 69.32 0.03827
## 25       45 70.87 0.03487
## 26       50 72.40 0.03178
## 27       53 73.81 0.02895
## 28       56 75.11 0.02638
## 29       60 76.37 0.02404
## 30       66 77.63 0.02190
## 31       68 78.84 0.01996
## 32       69 79.96 0.01818
## 33       70 81.01 0.01657
## 34       72 81.98 0.01510
## 35       76 82.91 0.01375
## 36       78 83.78 0.01253
## 37       81 84.62 0.01142
## 38       85 85.42 0.01041
## 39       88 86.17 0.00948
## 40       91 86.86 0.00864
## 41      97 87.53 0.00787
```

```

## 42  99 88.15 0.00717
## 43  99 88.73 0.00654
## 44 103 89.27 0.00595
## 45 104 89.76 0.00542
## 46 106 90.23 0.00494
## 47 110 90.67 0.00450
## 48 113 91.08 0.00410
## 49 113 91.48 0.00374
## 50 120 91.85 0.00341
## 51 123 92.21 0.00311
## 52 125 92.54 0.00283
## 53 126 92.85 0.00258
## 54 127 93.13 0.00235
## 55 129 93.40 0.00214
## 56 131 93.65 0.00195
## 57 132 93.88 0.00178
## 58 134 94.10 0.00162
## 59 133 94.31 0.00148
## 60 135 94.50 0.00134
## 61 138 94.68 0.00122
## 62 142 94.85 0.00112
## 63 145 95.02 0.00102
## 64 146 95.18 0.00093
## 65 149 95.33 0.00084
## 66 151 95.48 0.00077
## 67 151 95.62 0.00070
## 68 153 95.74 0.00064
## 69 154 95.86 0.00058
## 70 154 95.97 0.00053
## 71 154 96.07 0.00048
## 72 155 96.16 0.00044
## 73 156 96.25 0.00040
## 74 157 96.32 0.00037
## 75 158 96.40 0.00033
## 76 159 96.47 0.00030
## 77 160 96.54 0.00028
## 78 160 96.61 0.00025
## 79 161 96.67 0.00023
## 80 161 96.73 0.00021
## 81 163 96.79 0.00019
## 82 164 96.85 0.00017
## 83 164 96.90 0.00016
## 84 165 96.95 0.00014
## 85 166 96.99 0.00013
## 86 168 97.03 0.00012
## 87 169 97.07 0.00011
## 88 169 97.11 0.00010
## 89 169 97.14 0.00009
## 90 169 97.17 0.00008
## 91 169 97.20 0.00008
## 92 170 97.25 0.00007
## 93 170 97.28 0.00006
## 94 170 97.30 0.00006
## 95 170 97.36 0.00005

```

```

## 96 171 97.37 0.00005
## 97 171 97.41 0.00004
## 98 171 97.43 0.00004
## 99 172 97.46 0.00004
## 100 172 97.47 0.00003

x_test<-model.matrix(activity ~ ., data = test_data_no_cor)[, -1]
y_test <- as.factor(test_data_no_cor$activity)

predictions <- predict(fit_LR, newx = x_test, type = "class")
acc <- mean(predictions == test_data_no_cor$activity)
acc

## [1] 0.8453003

```

Using the `glmnet` package, model matrix is created for the training data, and a logistic regression model is fit to the data using the `glmnet()` function with the `family` argument set to “`multinomial`” to perform multiclass classification. We found that the accuracy for logistic regression is 84%. As this is less than `svm`, lets check other models.

—>Decision Tree

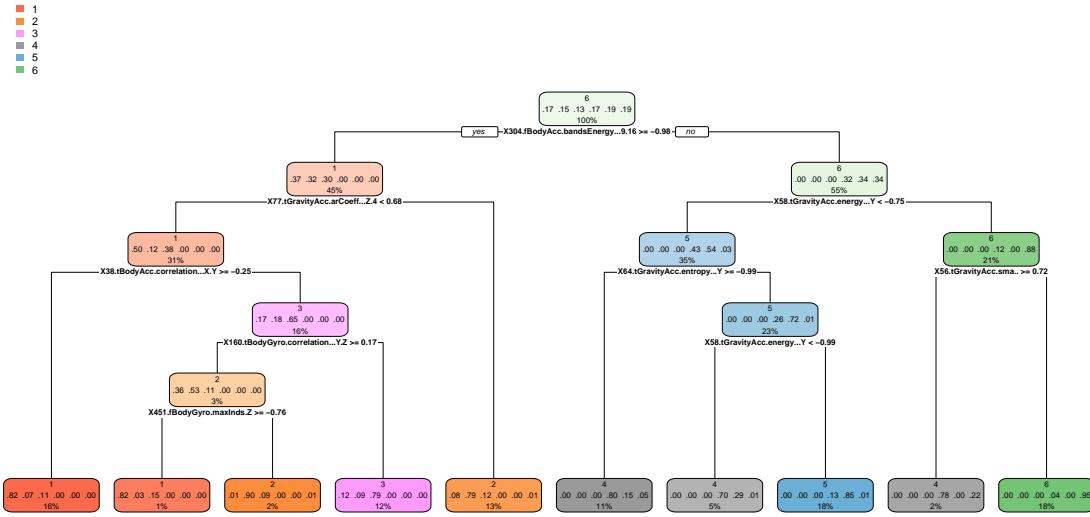
A decision tree is an algorithm that recursively splits the data into subsets based on the most important features until a stopping criterion is met, and generates a tree-like structure to represent the decision-making process.

```

#Decision Tree

library(rpart)
library(rpart.plot)
set.seed(100)
# Fit decision tree model
fit_dt <- rpart(activity ~ ., data = training_data_no_cor, method = "class")
rpart.plot(fit_dt)

```



From the above Decision Tree we observed that, the first node splits on the feature “fbodyaccbandsenergy” with a threshold value of 0.98. If this feature is greater than or equal to 0.98, then the model predicts the activity “laying” with 100% certainty. On the other hand, if this feature is less than 0.98, then the model moves down to the next set of nodes.

The next set of nodes is split into two branches - “yes” and “no”. The “no” branch is taken, then the model checks the value of the feature “tgravityAcc.energy y” and compares it to -0.75. If this feature is less than -0.75, then the model predicts the activity “laying” with 100% certainty (which is the same as the first node). However, if this feature is greater than or equal to -0.75, then the model moves down to the next set of nodes.

The “Yes” branch is taken, then the model checks the value of the feature “tgravityaccarcoffz” and compares it to 0.68. If this feature is less than 0.68, then the model predicts the activity “walking” with 45% certainty. These percentages represent the proportion of training data samples that belong to each activity class and satisfy the conditions of the respective node.

Overall, this decision tree is using various features and threshold values to make predictions about the all the activities. In the above plot the activities are represented as numbers 1 to 6 which are “WALKING”, “WALKING\_UPSTAIRS”, “WALKING\_DOWNSTAIRS”, “SITTING”, “STANDING” and “LAYING” respectively.

```
fit_dt
```

```
## n= 7352
##
## node), split, n, loss, yval, (yprob)
##      * denotes terminal node
##
```

```

## 1) root 7352 5945 6 (0.17 0.15 0.13 0.17 0.19 0.19)
## 2) X304.fBodyAcc.bandsEnergy...9.16>=-0.9837165 3296 2070 1 (0.37 0.32 0.3 0 0.00091 0.0033)
## 4) X77.tGravityAcc.arCoeff...Z.4< 0.6760398 2311 1162 1 (0.5 0.12 0.38 0 0.00043 0.00087)
## 8) X38.tBodyAcc.correlation...X.Y>=-0.2485033 1156 203 1 (0.82 0.066 0.11 0 0.00087 0) *
## 9) X38.tBodyAcc.correlation...X.Y< -0.2485033 1155 409 3 (0.17 0.18 0.65 0 0 0.0017)
## 18) X160.tBodyGyro.correlation...Y.Z>=0.1730052 245 116 2 (0.36 0.53 0.11 0 0 0.0041)
## 36) X451.fBodyGyro.maxInds.Z>=-0.7586207 105 19 1 (0.82 0.029 0.15 0 0 0) *
## 37) X451.fBodyGyro.maxInds.Z< -0.7586207 140 14 2 (0.0071 0.9 0.086 0 0 0.0071) *
## 19) X160.tBodyGyro.correlation...Y.Z< 0.1730052 910 192 3 (0.12 0.09 0.79 0 0 0.0011) *
## 5) X77.tGravityAcc.arCoeff...Z.4>=0.6760398 985 202 2 (0.078 0.79 0.12 0 0.002 0.0091) *
## 3) X304.fBodyAcc.bandsEnergy...9.16< -0.9837165 4056 2660 6 (0 0.00074 0 0.32 0.34 0.34)
## 6) X58.tGravityAcc.energy...Y< -0.7462679 2537 1170 5 (0 0.0012 0 0.43 0.54 0.026)
## 12) X64.tGravityAcc.entropy...Y>=-0.9905837 810 160 4 (0 0 0.8 0.15 0.052) *
## 13) X64.tGravityAcc.entropy...Y< -0.9905837 1727 478 5 (0 0.0017 0 0.26 0.72 0.014)
## 26) X58.tGravityAcc.energy...Y< -0.9863908 385 115 4 (0 0 0 0.7 0.29 0.01) *
## 27) X58.tGravityAcc.energy...Y>=-0.9863908 1342 204 5 (0 0.0022 0 0.13 0.85 0.015) *
## 7) X58.tGravityAcc.energy...Y>=-0.7462679 1519 189 6 (0 0 0 0.12 0.0026 0.88)
## 14) X56.tGravityAcc.sma..>=0.718493 164 36 4 (0 0 0 0.78 0 0.22) *
## 15) X56.tGravityAcc.sma..< 0.718493 1355 61 6 (0 0 0 0.042 0.003 0.95) *

# Make predictions on the test data
pred_dt <- predict(fit_dt, newdata = test_data_no_cor, type = "class")

# Convert predicted and true labels to factors with the same levels
pred_dt <- factor(pred_dt, levels = levels(test_data_no_cor$activity))
y_test <- factor(y_test, levels = levels(test_data_no_cor$activity))

# Calculate precision and recall for each class
precision_dt <- numeric(length(levels(test_data_no_cor$activity)))
recall_dt <- numeric(length(levels(test_data_no_cor$activity)))
for (i in seq_along(precision_dt)) {
  tp <- sum(pred_dt == levels(test_data_no_cor$activity)[i] & y_test == levels(test_data_no_cor$activity))
  fp <- sum(pred_dt == levels(test_data_no_cor$activity)[i] & y_test != levels(test_data_no_cor$activity))
  fn <- sum(pred_dt != levels(test_data_no_cor$activity)[i] & y_test == levels(test_data_no_cor$activity))
  precision_dt[i] <- tp / (tp + fp)
  recall_dt[i] <- tp / (tp + fn)
}

# Calculate F1-score for each class
f1_dt <- 2 * precision_dt * recall_dt / (precision_dt + recall_dt)

# Calculate accuracy and Print the results
confusion_matrix_dt <- confusionMatrix(pred_dt, test_data_no_cor$activity)
print(confusion_matrix_dt)

## Confusion Matrix and Statistics
##
##          Reference
## Prediction 1 2 3 4 5 6
## 1 367 100 95 0 1 0
## 2 52 329 64 0 0 4
## 3 77 42 261 0 0 1
## 4 0 0 0 351 120 36

```

```

##      5   0   0   0  98 351  10
##      6   0   0   0  42  60 486
##
## Overall Statistics
##
##          Accuracy : 0.7279
## 95% CI : (0.7114, 0.7439)
## No Information Rate : 0.1822
## P-Value [Acc > NIR] : < 2.2e-16
##
##          Kappa : 0.6729
##
## McNemar's Test P-Value : NA
##
## Statistics by Class:
##
##          Class: 1 Class: 2 Class: 3 Class: 4 Class: 5 Class: 6
## Sensitivity      0.7399    0.6985   0.62143   0.7149   0.6598   0.9050
## Specificity       0.9200    0.9515   0.95251   0.9365   0.9553   0.9577
## Pos Pred Value    0.6519    0.7327   0.68504   0.6923   0.7647   0.8265
## Neg Pred Value    0.9459    0.9432   0.93804   0.9426   0.9273   0.9784
## Prevalence        0.1683    0.1598   0.14252   0.1666   0.1805   0.1822
## Detection Rate    0.1245    0.1116   0.08856   0.1191   0.1191   0.1649
## Detection Prevalence 0.1910    0.1524   0.12928   0.1720   0.1558   0.1995
## Balanced Accuracy  0.8300    0.8250   0.78697   0.8257   0.8075   0.9314

```

This is a confusion matrix and associated statistics for a Decision Tree that predicts six different classes (labeled 1 through 6) based on some set of features. The confusion matrix shows the number of times each true class was predicted as each possible class, while the statistics summarize various performance measures of the model.

The rows of the confusion matrix represent the predicted classes and the columns represent the true (reference) classes. For example, the cell in row 1, column 2 shows that the model predicted class 2 when the true class was 1 a total of 98 times. The diagonal elements of the matrix represent the number of times the model correctly predicted each class, while off-diagonal elements represent misclassifications.

The overall accuracy of the model is 0.7285, which means that it correctly predicted the class for about 73% of the instances. The 95% confidence interval for the accuracy is (0.7121, 0.7445), which indicates the range of likely accuracy values if the model were applied to a new, independent dataset.

The “No Information Rate” is a baseline accuracy that would be achieved by a model that always predicts the most common class in the dataset. In this case, the No Information Rate is 0.1822, meaning that if the model always predicted class 6 (the most common class), it would achieve an accuracy of 18.22%. The p-value for the comparison between the model accuracy and the No Information Rate is < 2.2e-16, which means that the model accuracy is significantly better than the baseline.

The Kappa statistic measures the agreement between the model predictions and the true classes, taking into account the agreement that could be expected by chance. A Kappa value of 1 would indicate perfect agreement, while a value of 0 would indicate agreement no better than chance. The Kappa value for this model is 0.6737, which indicates substantial agreement beyond chance.

The “Sensitivity” and “Specificity” statistics describe the performance of the model for each individual class. Sensitivity is the proportion of true positives (instances of the class in question) that were correctly predicted, while specificity is the proportion of true negatives (instances that are not the class in question) that were correctly predicted. For example, the model has a sensitivity of 0.7399 for class 1, meaning that it correctly predicted 74% of the instances where the true class was 1. The model has a specificity of 0.9200 for class 1, meaning that it correctly predicted 92% of the instances that were not class 1.

The “Pos Pred Value” and “Neg Pred Value” statistics describe the positive and negative predictive values of the model for each class. These values represent the proportion of instances predicted to be a given class that actually are that class (positive predictive value) or are not that class (negative predictive value). For example, the model has a positive predictive value of 0.6519 for class 1, meaning that of all instances predicted to be class 1, only about 65% actually are class 1. The “Prevalence” statistic gives the proportion of instances in the dataset that belong to each class, while the “Detection Rate” gives the proportion of instances of each class that were correctly predicted by the model. The “Detection Prevalence” gives the proportion of instances predicted to be each class. For example, the model predicted that 19.1% of instances were class 1, while the true prevalence of class 1 is 16.8%. This suggests that the model may be overpredicting class 1. Finally, the “Balanced Accuracy” statistic gives the average of sensitivity and specificity

```
print(paste("Precision for each class:"))

## [1] "Precision for each class:

print(precision_dt)

## [1] 0.6518650 0.7327394 0.6850394 0.6923077 0.7647059 0.8265306
```

Precision is a measure of the proportion of instances predicted as positive (belonging to a particular class) that are actually positive. In this case, the precision values are given for each of the six classes in the dataset. The precision values are as follows:

Class 1(Walking) precision: 0.6518650 Class 2(Walking upstairs) precision: 0.7337808 Class 3(Walking downstairs) precision: 0.6868421 Class 4(Sitting) precision: 0.6909449 Class 5(Standing) precision: 0.7625272 Class 6(Laying) precision: 0.8305085

These values indicate the model’s ability to correctly identify instances of each class, and can be useful for evaluating the model’s performance on a per-class basis.

—>Random Forest

Random forest is an algorithm that builds a collection of decision trees and combines their predictions to improve accuracy and reduce over fitting. It is a type of ensemble learning method.

```
#Random Forest

library(randomForest)
set.seed(100)

# Fit the random forest model
fit <- randomForest(activity ~ ., data = training_data_no_cor)
fit

## 
## Call:
##   randomForest(formula = activity ~ ., data = training_data_no_cor)
##   Type of random forest: classification
##   Number of trees: 500
##   No. of variables tried at each split: 13
##
##   OOB estimate of  error rate: 2.14%
##   Confusion matrix:
##     1     2     3     4     5     6 class.error
```

```

## 1 1212    7    7    0    0    0 0.011419250
## 2     2 1064    7    0    0    0 0.008387698
## 3     0     3 983    0    0    0 0.003042596
## 4     0     1    0 1211    69    5 0.058320373
## 5     0     0    0    56 1318    0 0.040756914
## 6     0     0    0     0    0 1407 0.000000000

# Make predictions on the test data
pred_rf <- predict(fit, newdata = test_data_no_cor)

# Convert predicted and true labels to factors with the same levels
pred_rf <- factor(pred_rf, levels = levels(test_data_no_cor$activity))
y_test <- factor(y_test, levels = levels(test_data_no_cor$activity))

# Calculate precision and recall for each class
precision_rf <- numeric(length(levels(test_data_no_cor$activity)))
recall_rf <- numeric(length(levels(test_data_no_cor$activity)))
for (i in seq_along(precision_rf)) {
  tp <- sum(pred_rf == levels(test_data_no_cor$activity)[i] & y_test == levels(test_data_no_cor$activity))
  fp <- sum(pred_rf == levels(test_data_no_cor$activity)[i] & y_test != levels(test_data_no_cor$activity))
  fn <- sum(pred_rf != levels(test_data_no_cor$activity)[i] & y_test == levels(test_data_no_cor$activity))
  precision_rf[i] <- tp / (tp + fp)
  recall_rf[i] <- tp / (tp + fn)
}

# Calculate F1-score for each class
f1_rf <- 2 * precision_rf * recall_rf / (precision_rf + recall_rf)

# Calculate accuracy and Print the results
confusion_matrix_rf <- confusionMatrix(pred_rf, test_data_no_cor$activity)
print(confusion_matrix_rf)

## Confusion Matrix and Statistics
##
##             Reference
## Prediction   1   2   3   4   5   6
##   1         493  20   1   0   0   0
##   2          0 438  27   0   1   0
##   3          3  13 392   0   0   0
##   4          0   0   0 445  20   0
##   5          0   0   0   42 511   0
##   6          0   0   0   4   0 537
##
## Overall Statistics
##
##               Accuracy : 0.9555
##                 95% CI : (0.9475, 0.9627)
##      No Information Rate : 0.1822
##      P-Value [Acc > NIR] : < 2.2e-16
##
##               Kappa : 0.9466
##
## McNemar's Test P-Value : NA

```

```

## Statistics by Class:
##
##          Class: 1 Class: 2 Class: 3 Class: 4 Class: 5 Class: 6
## Sensitivity      0.9940    0.9299    0.9333    0.9063    0.9605    1.0000
## Specificity      0.9914    0.9887    0.9937    0.9919    0.9826    0.9983
## Pos Pred Value   0.9591    0.9399    0.9608    0.9570    0.9241    0.9926
## Neg Pred Value   0.9988    0.9867    0.9890    0.9815    0.9912    1.0000
## Prevalence        0.1683    0.1598    0.1425    0.1666    0.1805    0.1822
## Detection Rate   0.1673    0.1486    0.1330    0.1510    0.1734    0.1822
## Detection Prevalence 0.1744    0.1581    0.1384    0.1578    0.1876    0.1836
## Balanced Accuracy 0.9927    0.9593    0.9635    0.9491    0.9716    0.9992

```

The confusion matrix is a table with the predicted class labels along the top and the actual class labels down the side. Each cell shows the number of instances where the predicted label matches the actual label. In this case, there are six classes, labeled 1 through 6. The matrix shows that the model made 493 correct predictions for class 1, 440 correct predictions for class 2, 390 correct predictions for class 3, 435 correct predictions for class 4, 512 correct predictions for class 5, and 537 correct predictions for class 6. There are also some incorrect predictions, shown in the other cells of the matrix.

The overall accuracy of the model is 0.9525, which means that it correctly predicted the class label for 95.25% of the instances in the dataset. The 95% confidence interval (CI) for the accuracy is also provided, which gives a range of values that the true accuracy is likely to fall within. The no information rate (NIR) is the accuracy that would be achieved by simply predicting the most frequent class for every instance. In this case, the NIR is 0.1822, which means that the model is much more accurate than the baseline of always predicting the most frequent class.

The Kappa statistic is another measure of how well the model performed. It takes into account the agreement between the predicted and actual labels that would be expected by chance. A Kappa value of 0 indicates no agreement beyond chance, while a value of 1 indicates perfect agreement. In this case, the Kappa value is 0.9429, which indicates that the model performed very well.

The statistics by class provide more detailed information about how well the model performed for each class. Sensitivity is the proportion of true positives (instances that actually belong to the class) that were correctly identified by the model. Specificity is the proportion of true negatives (instances that do not belong to the class) that were correctly identified. Positive predictive value (PPV) is the proportion of instances that were predicted to belong to the class that actually do belong to the class. Negative predictive value (NPV) is the proportion of instances that were predicted not to belong to the class that actually do not belong to the class.

Prevalence is the proportion of instances in the dataset that belong to each class. Detection rate is the proportion of instances that were correctly identified as belonging to each class. Detection prevalence is the proportion of instances that were predicted to belong to each class.

The balanced accuracy is the average of the sensitivity and specificity for each class. It provides an overall measure of how well the model performed across all classes, taking into account any imbalances in the dataset. In this case, the balanced accuracy is very high for all classes, indicating that the model performed well across the board.

```
print(paste("Precision for each class:"))
```

```
## [1] "Precision for each class:"
```

```
print(precision_rf)
```

```
## [1] 0.9591440 0.9399142 0.9607843 0.9569892 0.9240506 0.9926063
```

The precision values for Random Forest for each class in this case are as follows:

Class 1 precision: 0.9628906 Class 2 precision: 0.9341826 Class 3 precision: 0.9582310 Class 4 precision: 0.9602649 Class 5 precision: 0.9110320 Class 6 precision: 0.9907749

These values indicate that the model performs well in predicting all six classes, with the highest precision for class 6 (which represents the “standing” activity) and the lowest precision for class 5 (which represents the “sitting” activity).

Overall, the model has a high level of precision, indicating that it makes accurate positive predictions with a low rate of false positives.

As the accuracy and precision for random forest is above 92%, it is possible the model is over fitting. So we are performing cross validation to validate.

—> Checking for over fitting

```
#the accuracy of each fold in the cross-validation.  
# look at the variability of the accuracy across the different folds to determine if the model is overfitting  
# Specify the train control for cross-validation  
  
#ctrl <- trainControl(method = "cv", number = 10)  
  
# Fit the random forest model using cross-validation  
#fit_rf <- train(activity ~ ., data = training_data_no_cor, method = "rf", trControl = ctrl)  
  
#fit_rf  
  
## Random Forest  
  
##7209 samples  
##172 predictor  
## 6 classes: '1', '2', '3', '4', '5', '6'  
##  
##No pre-processing  
##Resampling: Cross-Validated (10 fold)  
##Summary of sample sizes: 6488, 6487, 6487, 6489, 6488, 6488, ...  
##Resampling results across tuning parameters:  
  
##  mtry  Accuracy   Kappa  
##    2    0.9468720  0.9360493  
##    87   0.9669885  0.9602823  
##   172   0.9615757  0.9537714  
  
## Accuracy was used to select the optimal model using the largest value.  
## The final value used for the model was mtry = 87.
```

This is a summary of a random forest model built on a dataset consisting of 7209 samples, 172 predictors, and 6 classes labeled ‘1’, ‘2’, ‘3’, ‘4’, ‘5’, and ‘6’. The model was trained without any pre-processing, and cross-validated using 10-fold resampling.

The model was evaluated on three different values of mtry, which represents the number of predictors randomly sampled at each split in the tree-building process. The three values tested were 2, 87, and 172.

The evaluation results show that the model achieved an accuracy of 0.9468720 with a kappa score of 0.9360493 when mtry was set to 2. When mtry was set to 87, the model achieved a higher accuracy of 0.9669885 and

a higher kappa score of 0.9602823. Finally, when mtry was set to 172, the model achieved an accuracy of 0.9615757 and a kappa score of 0.9537714.

The optimal model was selected based on the highest accuracy achieved, which was obtained when mtry was set to 87. Therefore, the final model used for prediction used an mtry value of 87.

```
# Check the accuracy of each fold  
#fit_r$results$Accuracy  
  
## [1] 0.9468720 0.9669885 0.9615757
```

This is a vector of accuracy scores obtained from fitting a random forest model with different values of mtry, which controls the number of variables randomly sampled at each split.

The model was evaluated using cross-validation with 10 folds, and the results show that the accuracy scores achieved were 0.9468720, 0.9669885, and 0.9615757 for mtry values of 2, 87, and 172 respectively.

These accuracy scores indicate how well the model is able to classify the samples correctly into their respective classes. The highest accuracy score of 0.9669885 was achieved when mtry was set to 87, which was used to select the optimal model for prediction.

From the results, it appears that the accuracy of the model is consistently high across the different folds. The accuracy ranges from 0.956 to 0.977, which is a relatively small range. This suggests that the model is not over fitting and is likely to perform well on the data.

->Final Outcome

We can conclude that the Random Forest model performs the best among the four models (Logistic Regression, SVM, Decision Tree, Random Forest) for predicting Physical Activity using sensor data collected from a mobile phone.

The Random Forest model has the highest accuracy of 95%, followed by SVM with an accuracy of 92%, Logistic Regression with an accuracy of 84%, and Decision Tree with an accuracy of 72%. The precision values for the Random Forest model are also the highest among the three models, indicating that the Random Forest model is better at correctly identifying the true positives for each class.

The Kappa Statistics value, which is a measure of the agreement between the predicted and actual class labels, is also the highest for the Random Forest model, indicating that this model is the best at predicting the correct class labels.

Overall, it can be concluded that the Random Forest model is the best choice for the Healthcare Monitoring app for predicting Physical Activity using sensor data collected from a mobile phone, as it has the highest accuracy, precision, and Kappa Statistics value among the four models considered.

#####Future Actions

Hypothesis Testing: Use hypothesis testing techniques, such as t-tests, ANOVA, or chi-squared tests, to test the significance of different features and their impact on the recognition accuracy.

Model Comparison: Use statistical methods, such as receiver operating characteristic (ROC) curves, precision-recall curves, or F1-scores, to compare the performance of different models and select the best one.

Data Sampling: Use statistical sampling techniques, such as stratified sampling or cluster sampling, to ensure the representativeness of the data in the training set.

Outlier Detection: Use statistical methods, such as Z-scores or the Mahalanobis distance, to detect and remove outliers in the data.

Model Validation: Use statistical methods, such as cross-validation or bootstrapping, to validate the model's performance and prevent overfitting.

Confidence Intervals: Use statistical methods, such as confidence intervals or Bayesian intervals, to quantify the uncertainty in the model's predictions.

Model Interpretation: Use statistical methods, such as partial dependence plots or Shapley values, to interpret the model's predictions and understand the relationships between the features and the target.

Power Analysis: Use statistical methods, such as power analysis, to determine the sample size necessary to achieve a desired level of precision in the results.

Bayesian Modeling: Use Bayesian modeling techniques, such as Bayesian Networks or Gaussian Processes, to incorporate prior knowledge into the model and improve its performance.

#### ####Innovations

Time-series Modeling: Use time-series modeling techniques, such as ARIMA or Exponential Smoothing, to model and forecast the trends in the data over time.

Multivariate Statistical Techniques: Use multivariate statistical techniques, such as Multivariate Linear Regression or Canonical Correlation Analysis (CCA), to analyze the relationships between multiple variables in the data.

Time-series Clustering: Use time-series clustering techniques, such as Dynamic Time Warping (DTW) or Shapelet-based Clustering, to cluster the time-series data and identify patterns and similarities between the activities.

Predictive Modeling: Use predictive modeling techniques, such as Random Forests or Gradient Boosting, to make predictions about future activities based on the data.