

HarvardX: PH125.9x: Data Science: Capstone: Predicting the Chance of a Patient Falling Using Data from China

A Capstone Project for the Professional Certificate in Data Science offered by
Harvard University (HarvardX) via EdX

John King'athia Karuitha

Friday December 16, 2022

Contents

Abstract	2
Introduction	3
Objective of the Exercise	3
Summary of Results	3
The Data	3
Training set/ Test set split	4
Exploratory Data Analysis- Training set	5
Data Structure	5
Missing Data	6
Other Summary Statistics	10
Data Visualization	11
Correlation among the variables	14
Class Balance/ Imbalance	16
Method	16
Principal Components Analysis and Handling Class Imbalance	16
Visualizing the Principal Components	17
Running the ML models	29
Classification tree	29
The Random Forest Model	32
K-Nearest Neighbours (KNN)	33
Extreme Gradient Boosting (XGBoost)	33
Multinomial Logit Model	35
Ensemble	36
Model Evaluation	38
Conclusion	38
References	38

```
doParallel::registerDoParallel()
```

Abstract

In this project, I use data from China to predict whether or not a patient will fall. In hospitals, falls by patients can be fatal and it is paramount to have an idea which patient is likely or unlikely to fall so as to take appropriate interventions. The dependent variable **activity** consists of a range of patient activities among them falling. The independent variables capture a range of physiological conditions such as heart rate and blood flow. The independent variables are highly correlated and hence, I start by performing principal components analysis (PCA) on the data. I then run 5 models- classification tree, random forest, K-Nearest Neighbours (KNN), Extreme gradient boosting (XGBoost), and multinomial logit. Finally, I assemble these models into one predictive algorithm (the ensemble). The random forest model offers the best specificity (0.8859) while extreme gradient boosting has the highest sensitivity (0.5437). The ensembled algorithm offers the best balanced accuracy (0.7048).

Introduction

In hospitals, predicting that a patient will fall goes a long way in ensuring their well being. Seriously ill and elderly patients could easily die if they attempt to take a walk but fall. In this exercise, I use a redacted dataset from China to predict the chance that a patient will fall¹. I sourced the data from Kaggle. I downloaded the data and loaded it in my github account. From the outset, it important to note that it is much more expensive to predict that a patient will NOT fall who ends up falling. It is much more acceptable to predict a patient will fall and they do NOT fall. From this perspective, maximizing specificity of the model is paramount in the choice of the right model.

Note that I partly adopt the Tidymodels framework in building the machine learning model. Tidymodels is a suit of packages for building machine learning models developed by R-Studio. Max Kuhn, the developer of `caret`, the other popular machine learning platform in R is the lead scientist in developing Tidymodels. Like the tidyverse, Tidymodels provides a consistent API for predictive modelling. Tidymodels has a smoother, more intuitive work-flow compared to caret and saves the data analyst of having to keep repeating data preprocessing steps.

The exercise proceeds as follows. In the next section, I download and load the required R packages and then present a summary of the final results. I then load and describe the data. I then explain the method of analysis, and then do the analysis. I then discuss the results and conclude.

Objective of the Exercise

The objective of the exercise is to build a machine learning model that maximizes the chances of telling whether or not a patient is likely to fall.

Summary of Results

The results show that the random forest model does best in predicting patients do not fall who actually do not fall (specificity). The random forest model offers the best specificity (0.8859) while extreme gradient boosting has the highest sensitivity (0.5437). The ensemble algorithm offers the best balanced accuracy (0.7048). Overall, we can rely on the models to predict which patients are at risk of falling and taking measures to support them.

The Data

In this section, I load the data into a file called fall. The data has seven variables. The dependent variable is `category` and has six levels.

1. Standing (coded 0)
2. Walking (coded 1)
3. Sitting (coded 2)
4. Falling (coded 3)
5. Cramps (coded 4)
6. Running (coded 5)

The rest of the variables are independent.

1. Time: monitoring time.
2. SL: sugar level.
3. EEG: Electroencephalogram (EEG) is a test that detects electrical activity in the brain.
4. BP: blood pressure.
5. HR: heart beat rate.
6. Circulation: blood circulation.

¹The data is a redacted version of the one used in a study by Özdemir & Barshan (2017). However, the original dataset had over 300 features whilst the redacted version has 7 features

In total, the data has 16382 rows of data and 7 columns of data.

```
#####
# SECTION 3
#####
# Load the data ----

## Get the URL for the datafile
## The original source of the data is kaggle.com, specifically https://www.kaggle.com/pitasr/falldata
## However, it is hard to download data straight from kaggle without supplying login details
## hence, I loaded the data into my github account in the url below.

url <- "https://raw.githubusercontent.com/Karuitha/Final_Project_HarvardX/main/falldetection.csv"

## Download the dataset.

download.file(url = url, destfile = "fall.csv", method = "curl")

## load the dataset into R

fall <- read.csv("fall.csv") %>%

## Clean the column names by removing capital letters, spaces and special characters
janitor::clean_names()

## The dependent variable is activity, classified as follows
## Fall detection data set of Chinese hospitals of old age patients.

### 0- Standing
### 1- Walking
### 2- Sitting
### 3- Falling
### 4- Cramps
### 5- Running

## Independent variables

### time: monitoring time
### sl: sugar level
### eeg: eeg monitoring rate- electroencephalogram (EEG) is a test that detects electrical activity in ...
### bp blood pressure
### hr: heart beat rate
### circulation: blood circulation

# Convert the dependent variable into a factor with category falling as the base
fall$activity <- factor(fall$activity, levels = c(3, 0, 1, 2, 4, 5))
```

Training set/ Test set split

I split the data into a training set and a test set. Here I use the function `initial_split` from `tidymodels` where I provide the data, the proportion of the data that goes into the training set and the dependent variable. After generating the index, I use the functions `training` and `testing` to specify the training and testing sets.

```
set.seed(123, sample.kind = "Rounding")
```

```

# Specify the index to split data into training and testing set
index <- initial_split(fall, prop = 0.7, strata = activity)

# Sopecify the training set
fall_train <- training(index)

# Dimensions of the training set
dim(fall_train) # 11470 observations of 7 variables

## [1] 11465      7

# Get the testing set
fall_test <- testing(index)

# Dimensions of the testing set
dim(fall_test) # 4912 observations of 7 variables

## [1] 4917      7

```

Note that any additional analysis will involve only the training set with the testing set reserved for evaluation of the final model. Next, I explore the training data

Exploratory Data Analysis- Training set

Data Structure

I first examine the structure of the data. Note that except for the dependent variable that is a factor with six levels, all the other variables are numeric. The training dataset has 11465 observations of 7 variables. The dependent variable is **activity**, a factor variable with 6 levels as follows.

0. Standing
1. Walking
2. Sitting
3. Falling
4. Cramps
5. Running

The level of interest in this case is to forecast whether or not a patient will fall. When weighing the options, it is better to predict that a patient will not fall when in fact they do not fall. While its also good to predict which patient will fall, the former has more weight. For this reason, I evaluate the model mainly using **specificity** and **balanced accuracy** rather than **sensitivity**.

The dependent variables are as follows.

1. time: monitoring time
2. sl: sugar level.
3. eeg: eeg monitoring rate- electroencephalogram (EEG) is a test that detects electrical activity in the brain.
4. bp: blood pressure.
5. hr: heart beat rate.
6. circulation: blood circulation.

I summarise the data below.

```

## Overview of the training data.
## Structure of the training data
str(fall_train)

## 'data.frame':    11465 obs. of  7 variables:

```

```

## $ activity   : Factor w/ 6 levels "3","0","1","2",...: 2 2 2 2 2 2 2 2 2 ...
## $ time      : num  10256 6628 7837 10788 18609 ...
## $ sl        : num  66668 6491 5545 103850 300553 ...
## $ eeg       : num  -6050 -1923 -2661 -4705 -12391 ...
## $ bp        : int  97 23 41 85 73 264 61 18 63 92 ...
## $ hr        : int  207 100 126 237 464 465 390 111 299 365 ...
## $ circluation: int  3048 485 371 4804 10829 8741 6741 345 7334 3633 ...
## First 6 observations of the training dataset
head(fall_train) %>% knitr::kable(caption = "First Six Observations of the Training Set")

```

Table 1: First Six Observations of the Training Set

	activity	time	sl	eeg	bp	hr	circluation
13	0	10255.70	66668.00	-6050.00	97	207	3048
18	0	6628.08	6490.83	-1923.04	23	100	485
19	0	7836.78	5545.00	-2661.00	41	126	371
20	0	10787.80	103850.00	-4705.00	85	237	4804
27	0	18609.30	300553.00	-12391.00	73	464	10829
29	0	23099.10	280276.00	-13215.00	264	465	8741

```

## last 6 observations of the training dataset
tail(fall_train) %>% knitr::kable(caption = "Last Six Observations of the Training Set")

```

Table 2: Last Six Observations of the Training Set

	activity	time	sl	eeg	bp	hr	circluation
16325	5	9363.74	37928.80	-3795.490	38	180	2054
16334	5	13684.60	52970.90	-5976.250	53	288	2334
16346	5	10137.40	34148.60	-3518.520	95	176	1922
16357	5	4795.02	3404.89	-1322.750	41	71	249
16366	5	13250.40	51768.90	-6069.660	52	288	2334
16370	5	3825.38	1220.87	-768.261	37	51	102

```

## Number of rows in the training dataset
nrow(fall_train)

```

```

## [1] 11465
## Number of columns in the training dataset
ncol(fall_train)

```

```

## [1] 7

```

Missing Data

As the summary below shows, the training data has no missing data points. Hence, I explore the data using visualizations next using the `missmap` function from the `Amelia` package.

```

# ****
# Exploratory data analysis: missing data
## Check for missing data.
sapply(fall_train, is.na) %>%

```

```

## Get the colsums of the logical dataframe.
## The colsums represent missing data for each column.
colSums() %>%

## Make a tibble of column names and missing values.
dplyr::tibble(variables = names(fall), missing = .) %>%

## Arrange missing values in descending order of missingness
dplyr::arrange(desc(missing)) %>%

## Get the top 7 (number of columns)
## Our tibble has no missing data.
head(7) %>%

## make a nice table
knitr::kable(caption = "Missing Data in the Training Set")

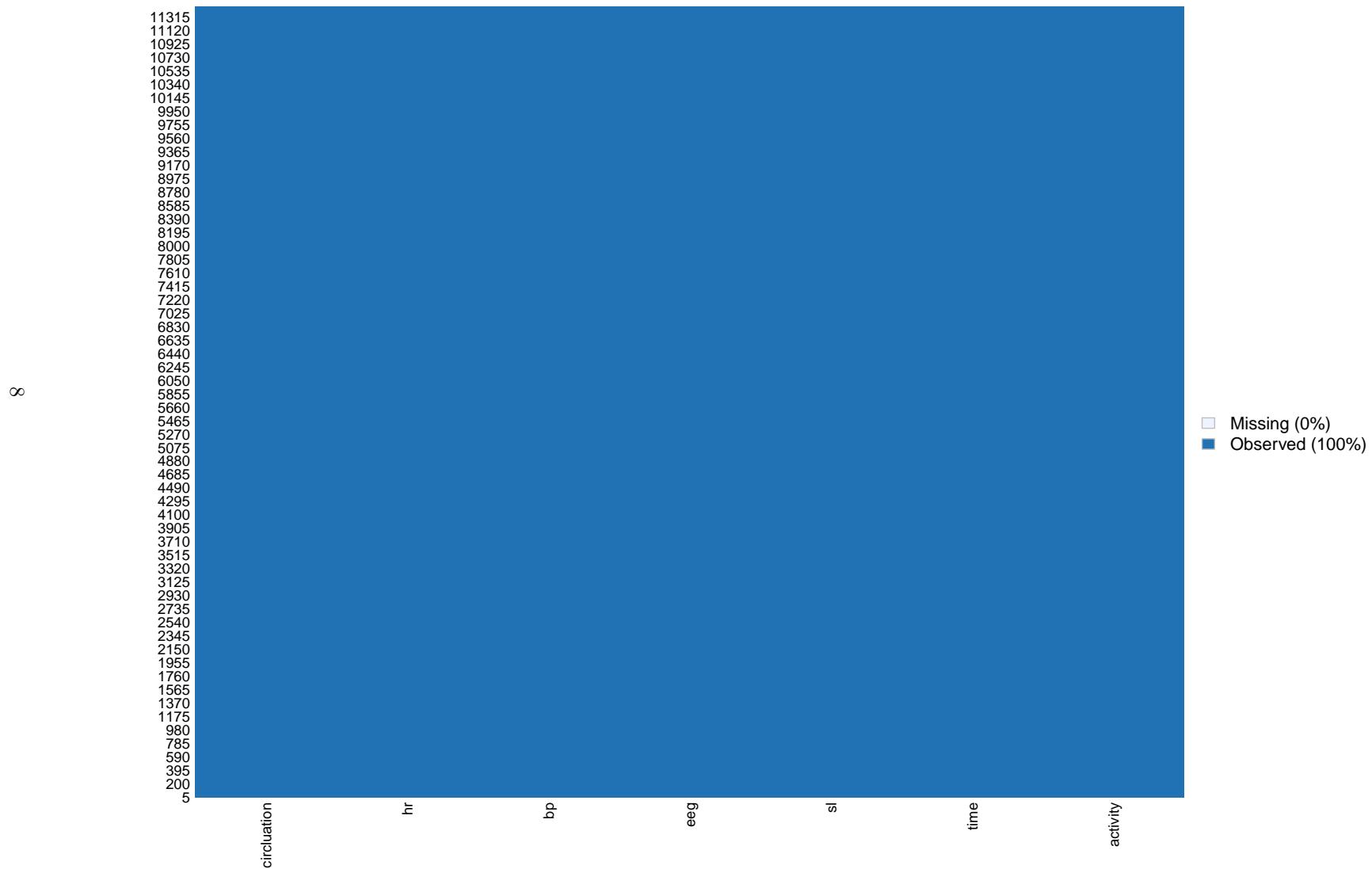
```

Table 3: Missing Data in the Training Set

variables	missing
activity	0
time	0
sl	0
eeg	0
bp	0
hr	0
circluation	0

```
## Visualizing missingness of data
Amelia::missmap(fall_train, main = "Figure 1: Missingness Map- Training Set")
```

Figure 1: Missingness Map– Training Set



```
## Again there is no missing data
```

Other Summary Statistics

The independent variable is categorical with six categories. The category of interest is 3 (falling). While the other categories are important, maximizing the prediction for falling remains the most important objective as the other activities are not as harmful.

```
## Summary statistics for the dependent variables
summary(fall_train$activity) %>% knitr::kable(caption = "Summary of Dependent Variable")
```

Table 4: Summary of Dependent Variable

	x
3	2511
0	3228
1	355
2	1751
4	2447
5	1173

Table 5 below shows the summary statistics for the independent variables (features). The data shows that there is a wide variation in the dataset with the standard deviation of the variables ranging from 48.38601 to 130029.85851.

```
#####
# Exploratory data analysis: summary statistics
fall_train %>%

## Deselect the activity column
select(-activity) %>%

## Make a table of summary statistics
skimr::skim() %>%

## Remove some uninformative columns
dplyr::select(-contains(c("missing", "complete", "hist", "skim_type"))) %>%

## Rename remaining columns
dplyr::rename(Variable = skim_variable,
             Mean = numeric.mean, SD = numeric.sd,
             Min = numeric.p0, Q1 = numeric.p25,
             Median = numeric.p50,
             Q3 = numeric.p75, Max = numeric.p100) %>%

## make a nice table
knitr::kable(caption = "Summary Statistics for the `Fall` dataset", align = "l")
```

Table 5: Summary Statistics for the Fall dataset

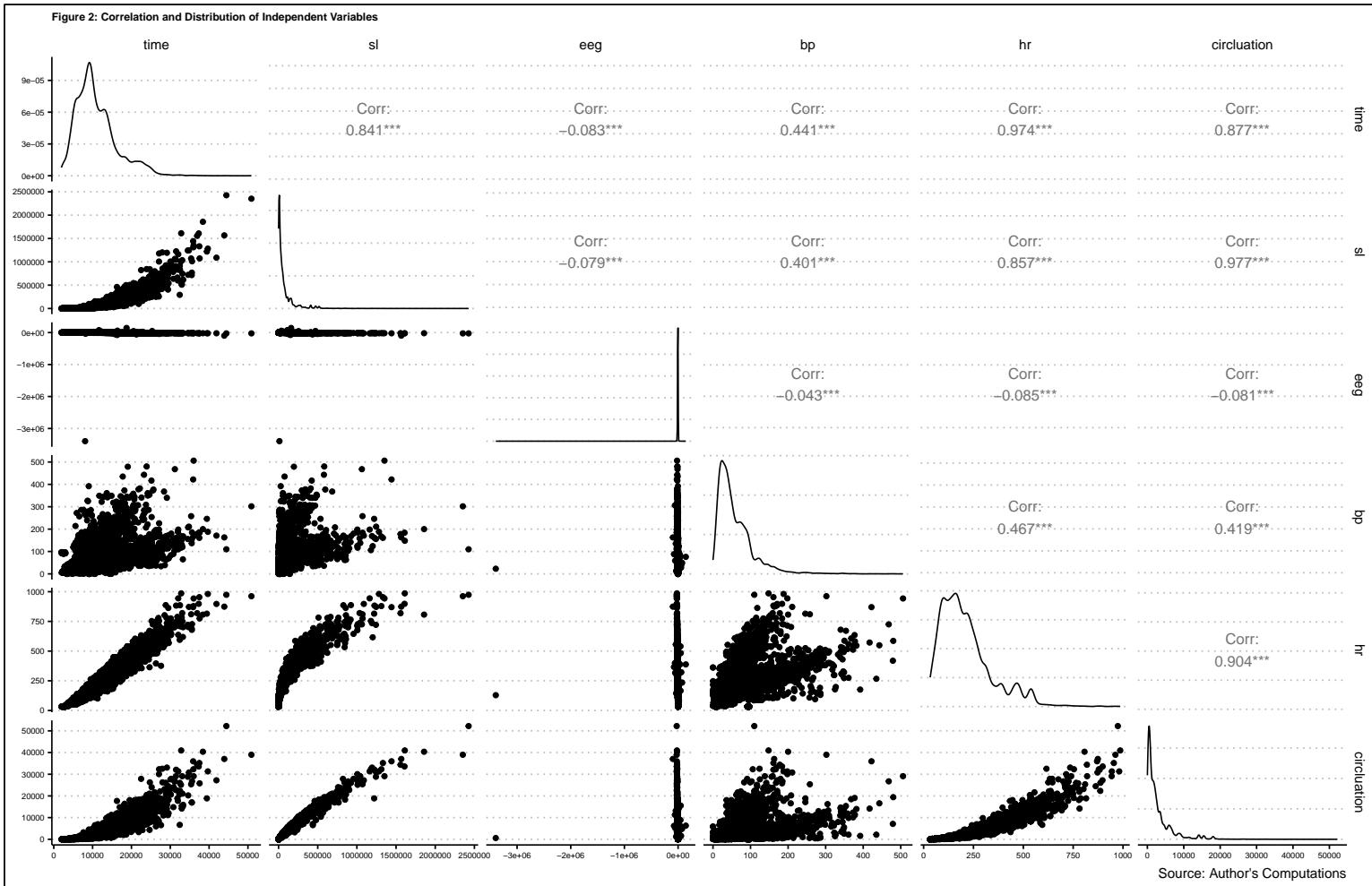
Variable	Mean	SD	Min	Q1	Median	Q3	Max
time	10963.86416	5279.52148	1.95423e+03	7269.79	9779.28	13503.90	50895.5
sl	75951.65790	129410.22133	4.22242e+01	10069.00	32234.50	80896.80	2426140.0
eeg	-4562.64414	31952.36356	-	-5646.00	-3372.05	-2170.96	144000.0
			3.39680e+06				
bp	58.43044	48.30781	0.00000e+00	26.00	44.00	78.00	506.0
hr	212.29551	130.28488	3.30000e+01	119.00	180.00	271.00	986.0
circulation	2913.92708	3850.79797	5.00000e+00	587.00	1626.00	3539.00	52210.0

Data Visualization

In this section, I visualize the training dataset. First, I visualize the correlation matrix that shows extremely high correlation between the independent variables, for instance between `circulation` and `heart rate`. In its current form, running models using collinear data is likely to yield unstable coefficients and hence unstable predictions. For this reason, I will run principal components analysis on the independent variables.

```
# Exploratory data analysis: data visualization
fall_train[,-1] %>% GGally::ggpairs() +  
  
## Add title and caption  
labs(title = "Figure 2: Correlation and Distribution of Independent Variables",  
    caption = "Source: Author's Computations") +  
  
## Add themes and adjust the font size plot title  
ggthemes::theme_clean() + theme(plot.title = element_text(size = 8)) +  
  
## Adjust font sizes of axis text  
theme(axis.text = element_text(size = 6))
```

Figure 2: Correlation and Distribution of Independent Variables



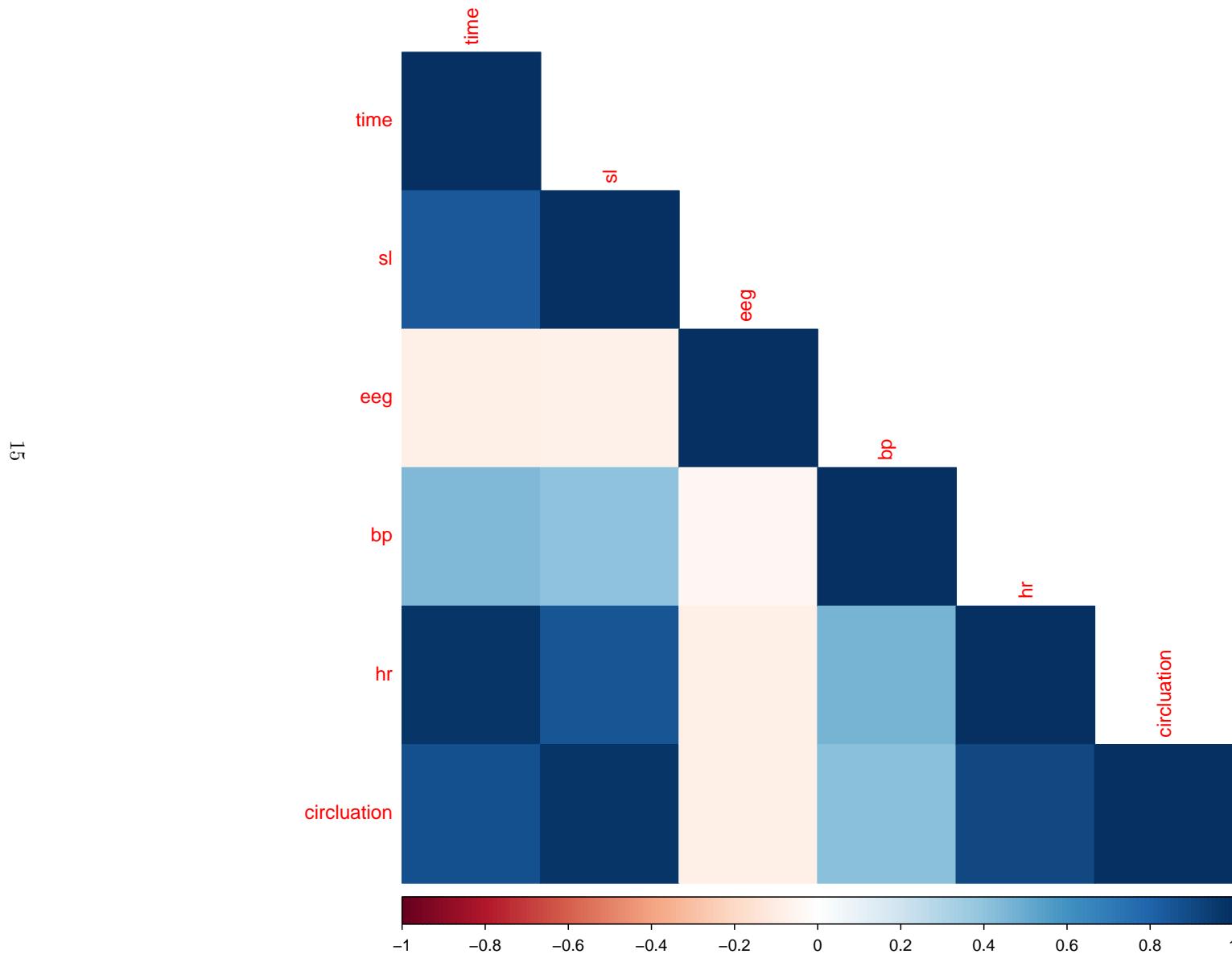
Correlation among the variables

Figure 2 shows the high degree of correlation between the independent variables which could make a model have unstable coefficients. I intend to correct the multicollinearity by running principal components analysis.

```
## Visualize the correlation
## Run correlation on independent variables and call the corrplot
cor(fall_train[,-1]) %>% corrplot::corrplot(method = "color",

      ## Specify the corrplot type and add title
      type = "lower", title = "Figure 3: A Visual of the Correlation Matrix - Training Set")
```

Figure 5: A visual of the Correlation Matrix – Training Set



Class Balance/ Imbalance

Here I check for data balance and find that class 1 has very low prevalence that is likely to affect predictive accuracy on the test set. I intend to correct for the low prevalence by up-sampling the classes that have very low prevalence in order to create a balanced dataset.

```
## Check for possible class imbalance on the dependent variable
## Make a table of class counts
table(fall_train$activity)

##
##      3      0      1      2      4      5
## 2511 3228 355 1751 2447 1173

## Make a proportionate table of class counts
prop.table(table(fall_train$activity))

##
##            3            0            1            2            4            5
## 0.2190144 0.2815526 0.0309638 0.1527257 0.2134322 0.1023114
# Class 1 has a problem of low prevalence.
```

Method

I adopt the following strategy;

- First, given that the data has high collinearity, I run a principal components analysis (PCA) and generate uncorrelated variables.
- I deal with the problem of imbalanced data by up-sampling the under-represented classes.
- I run seven machine learning models and then make an ensemble.
- I compare the performance of the models and choose the most optimal.

Principal Components Analysis and Handling Class Imbalance

In this section, I do principal components analysis (PCA) and balance the datasets. Note that I have applied the transformation to deal with extreme values that exist in our data.

```
## Create a PCA recipe
pca_recipe <- recipe(activity ~ ., data = fall_train)

## Do the PCA analysis
pca_trans <- pca_recipe %>%

## Ensure all predictors have a mean of zero
step_center(all_predictors()) %>%

## Ensure all predictors have a standard deviation of one
step_scale(all_predictors()) %>%

## Run the principal components analysis
step_pca(all_predictors()) %>%

## Apply adjustment to deal with outliers
step_spatialsign(all_predictors()) %>%
```

```

## Adjust data to deal with missing values
step_upsample(activity, over_ratio = 1) %>%
  ## Apply all the steps above and generate new dataset.
  prep()

```

Visualizing the Principal Components

Here, I extract the standard deviations, compute the variance and the cumulative variance for the PCs. The data shows that the first principal component accounts for about 66% of the variability while the second PC accounts for 17% of the variation.

```

## Check the names of the dependent variables
names(pca_trans)

## [1] "var_info"      "term_info"      "steps"        "template"
## [5] "levels"        "retained"       "requirements" "tr_info"
## [9] "orig_lvls"     "last_term_info"

## Access the standard deviations
sdev <- pca_trans$steps[[3]]$res$sdev

# View the standard deviations output
sdev

## [1] 1.9941791 0.9952037 0.8673943 0.4869939 0.1728212 0.1158289

## The contribution of PCA to the total variation
variance_explained <- (sdev ^ 2) / sum(sdev ^ 2)

# The variance explained by each principal component
variance_explained

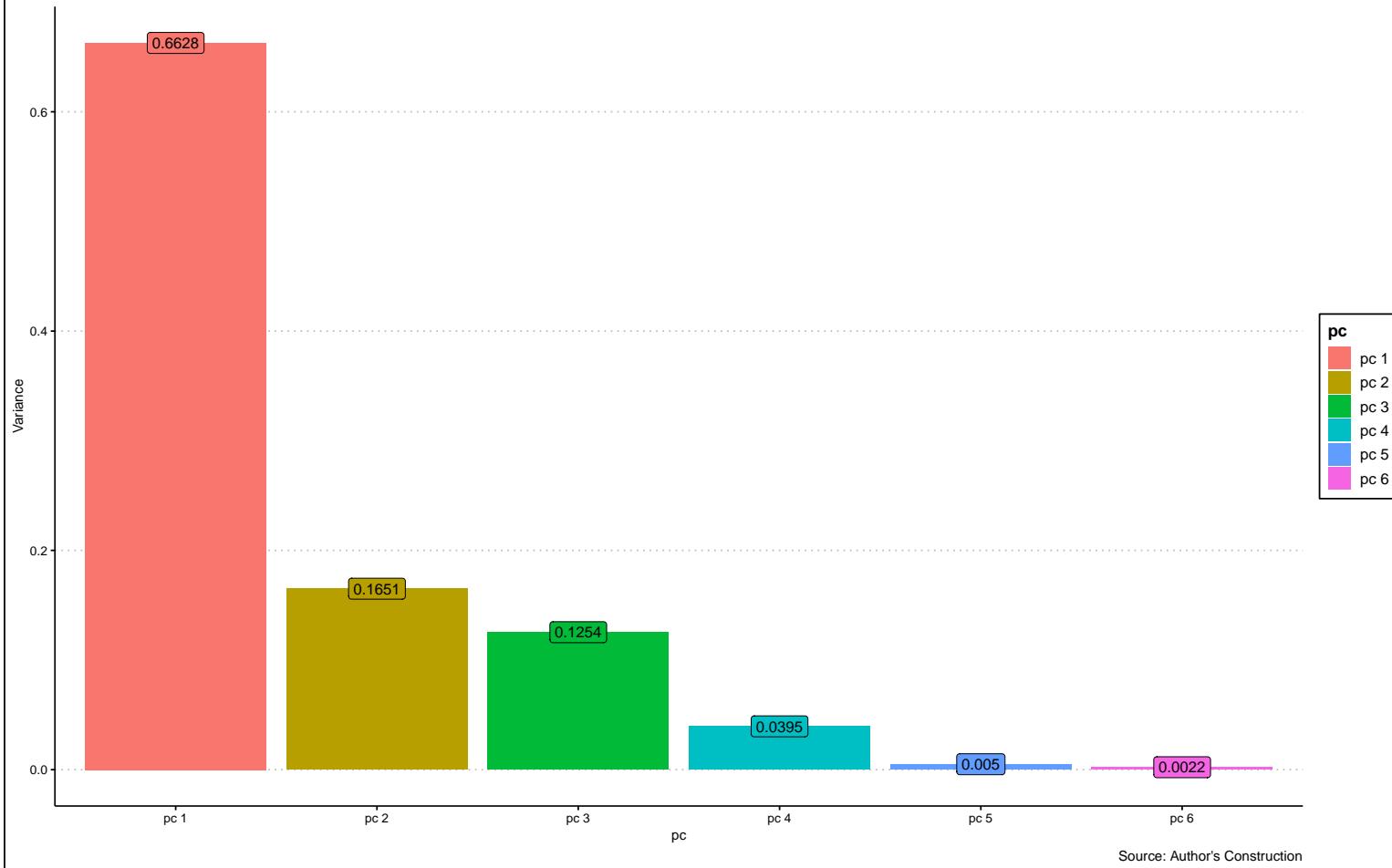
## [1] 0.662791684 0.165071746 0.125395480 0.039527175 0.004977860 0.002236054

```

Here, I make the scree plot showing the contribution of each principal component to the overall variability in the data.

```
## Plot a scree plot
## Create dataframe of principal components and variance explained by each PC>
data.frame(pc = paste("pc", 1:length(variance_explained)), variance = variance_explained) %>%
  # Call ggplot and supply the axes
  ggplot(aes(x = pc, y = variance, fill = pc)) +
  ## Add the geoms- geom_col and geom_label
  geom_col() + geom_label(show.legend = FALSE, mapping = aes(label = round(variance, 4))) +
  # Add a title
  ggtitle("Figure 4: Skree Plot - Contribution of Each PC to Total Variability") +
  # Remove title
  theme(legend.position = "none") +
  # Add a pleasant theme
  ggthemes::theme_clean() +
  # Add labels and caption
  labs(y = "Variance", caption = "Source: Author's Construction")
```

Figure 4: Skree Plot – Contribution of Each PC to Total Variability



In this section I make the plot for the cumulative variance. both the summary and the graph show that the first four principal components account for over 99% of the variability.

```
## Plot a scree plot
## The cumulative variance captured by the PCAs
variance_explained_cum <- cumsum(sdev ^ 2) / sum(sdev ^ 2)

## View the cumulative variance explained
variance_explained_cum

## [1] 0.6627917 0.8278634 0.9532589 0.9927861 0.9977639 1.0000000
```

```
## Create a dataframe of principal components and cumulative variance
data.frame(pc = paste("pc", 1:length(variance_explained_cum)),  

           variance = variance_explained_cum) %>%  

  
## Specify the axes
ggplot(aes(x = pc, y = variance, group = pc)) +  

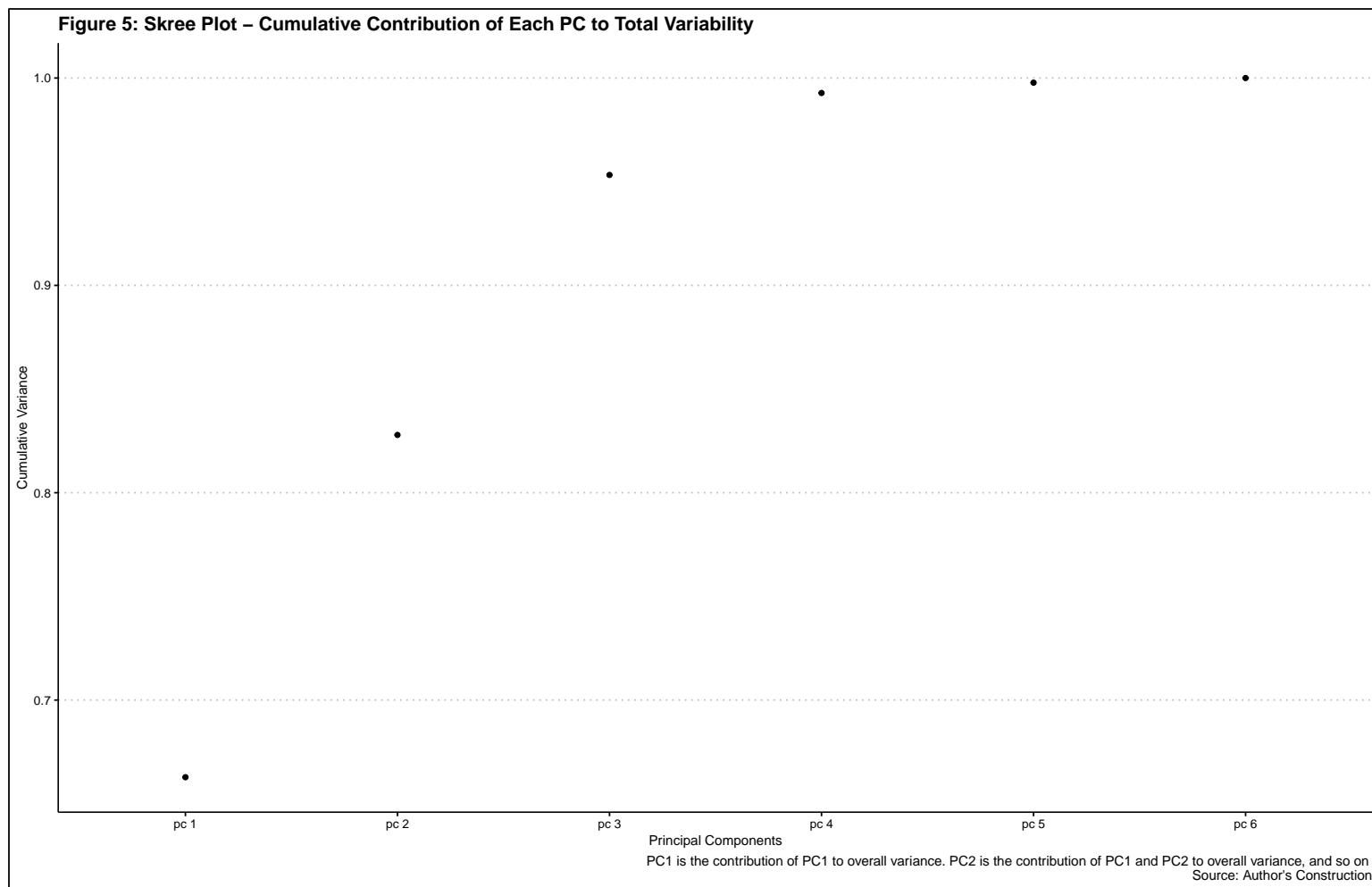
  
# Add a geom
geom_point() +  

  
# Add a title
ggtitle("Figure 5: Skree Plot - Cumulative Contribution of Each PC to Total Variability") +  

  
# remove legend
theme(legend.position = "none") +  

  
# Add a nice theme
ggthemes::theme_clean() +  

  
# Add axes labels and title
labs(y = "Cumulative Variance", x = "Principal Components", caption = "PC1 is the contribution of PC1 to overall variance. PC2 is the co
```



Next, I extract the summary statistics. Note that the new dataset for PCAs has 19368 rows of data. Note that the classes are now evenly balanced with each class having 3226 observations. Also, the summary shows that the variability between the variables has reduced markedly.

```
## Extract the transformed dataset
fall_train <- pca_trans %>% juice()

## The number of rows in the training dataset containing pcas
pca_trans %>% juice() %>% nrow()

## [1] 19368

## Checking class balances by dependent variable
table(pca_trans %>%

      ## Extract transformed data
      juice() %>%

      ## Select dependent variable to check for balance in classes
      select(activity)) %>%

      ## Make a nice table and add title
      knitr::kable(caption = "Distribution of Classes")
```

Table 6: Distribution of Classes

activity	Freq
3	3228
0	3228
1	3228
2	3228
4	3228
5	3228

```
## Summary of the PCAs
pca_trans %>%

## Extract transformed data
juice() %>%

# Select all variables except activity
select(-activity) %>%

## Summarize the data
skim_without_charts() %>%

## Select some variables in the resulting table
select(-skim_type, -n_missing, -complete_rate) %>%

## Make a nice table
knitr::kable(caption = "Summary Statistics for PCAs in the training set")
```

Table 7: Summary Statistics for PCAs in the training set

skim_variable	numeric.mean	numeric.sd	numeric.p0	numeric.p25	numeric.p50	numeric.p75	numeric.p100
PC1	-0.3311002	0.7324191	-	-0.9443104	-0.7607755	0.3784137	0.9987812
			0.9993410				
PC2	-0.0029457	0.0443862	-	-0.0193300	-0.0003134	0.0142165	0.9989081
			0.9170863				
PC3	-0.0461149	0.4929493	-	-0.3216948	-0.1837881	0.1533594	0.9992720
			0.9983935				
PC4	0.0519831	0.2987749	-	-0.1862040	-0.0179253	0.2256539	0.9970300
			0.5785683				
PC5	0.0016704	0.1220847	-	-0.0428751	0.0095705	0.0502035	0.9220046
			0.8874046				

In this section, I examine the distribution of the values of each principal components to see the extent to which outliers exist.

```
## Checking for extreme values
pca_trans %>% juice() %>%

  ## Convert the data to tidy format
  pivot_longer(-activity, names_to = "pc", values_to = "value") %>%

  ## Filter for values that meet threshold of 7.5 to filter out extreme values
  filter(value > 7.5)

## # A tibble: 0 x 3
## # ... with 3 variables: activity <fct>, pc <chr>, value <dbl>
## Plot histogram for PCAs
pca_trans %>% juice() %>%

  ## Convert data to tidy format
  pivot_longer(-activity, names_to = "pc", values_to = "value") %>%

  ## Filter for values within limits to avoid extreme values
  filter(value <= 7.5 & value >= -7.5) %>%

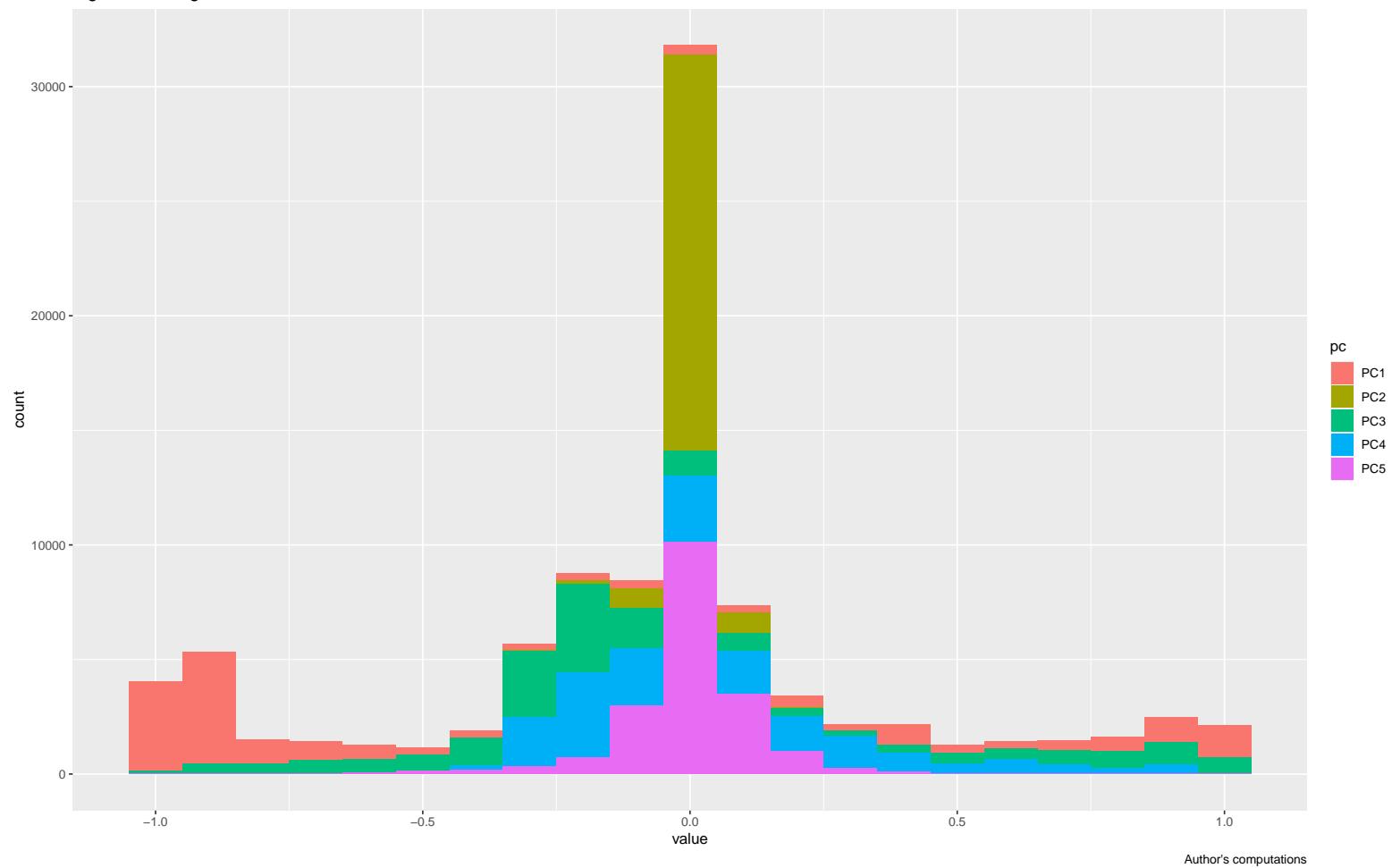
  ## Plot the data by supplying axes
  ggplot(mapping = aes(x = value, fill = pc),

         color = "black") +

  ## Add the geom and specify binwidth
  geom_histogram(binwidth = 0.1) +

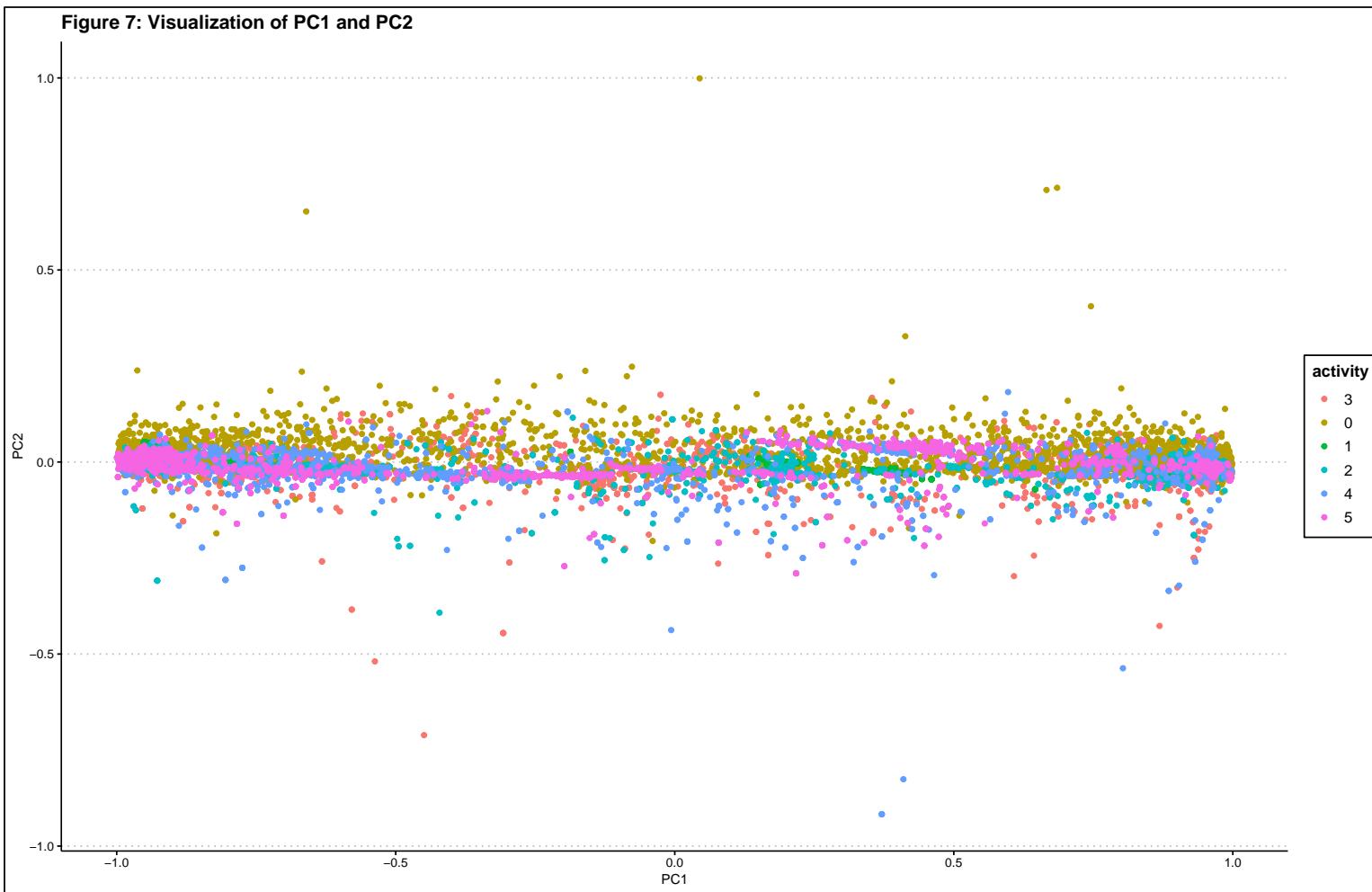
  ## Add titles and labels
  labs(title = "Figure 6: Histogram of PCs", caption = "Author's computations")
```

Figure 6: Histogram of PCs



Next, I visualize the first 2 PCs

```
pca_trans %>%  
  
  ## Extract transformed variables - the principal components  
  juice() %>%  
  
  ## Plot the first two principal components  
  ggplot(mapping = aes(x = PC1, y = PC2, color = activity)) +  
  
  ## Add a geom and a pleasant theme  
  geom_point() + ggthemes::theme_clean() +  
  
  ## Add labels and titles  
  labs(title = "Figure 7: Visualization of PC1 and PC2")
```



Next, I apply the same transformation I have made on the training set to the testing set. Here, I use a function `bake`.

```
## Transform the testing data similar to the training data
fall_test <- pca_trans %>%
  ## Use of bake function to generate testing test transformed exactly like the training set
  bake(fall_test)
```

Running the ML models

I now run machine learning models and evaluate each of them in the following order.

1. Classification tree.
2. Random Forest.
3. K-Nearest Neighbours.
4. Extreme Gradient Boosting.
5. Generalized Multinomial Logit Model.
6. Ensemble of all the models above.

Note that in evaluating the models, I will focus mainly on the specificity given that it is more expensive to predict no-fall in a patient who falls than it is to predict a fall in a patient who does not fall. Hence, a model that predicts that a patient will not fall with a higher precision is better than one that predicts a fall with equal level of precision. However, I will also, side by side, consider overall accuracy and specificity in case of ties. Also, given that our model has more than two classes, I consider the specificity of the main class `fall`. While the other classes are important, our case will be specific to predicting which patients fall or do not fall. The baseline accuracy on the test set, guessing the most frequent outcome, is 0.219036.

In all cases, I run a 10-fold cross validation and set a random seed as follows.

```
## Set seed to be used in the models
seeds <- set.seed(123, sample.kind = "Rounding")

## Set up cross validation parameters
control <- trainControl(method = "repeatedcv",
                         repeats = 10,
                         seeds = seeds)
```

Classification tree

In this section, I run the classification tree using the code chunk below. The tree model has overall accuracy 0.5138 against a no information rate (NIR) of 0.2814. The model has a specificity of 0.86453 and a sensitivity of 0.44866. Note that the optimal complexity parameter is near zero.

```
# The classification tree model
tree <- caret::train(activity ~ .,
                      data = fall_train,
                      # specify engine to use
                      method = "rpart",
                      # Set up cross validation
                      trControl = control,
                      # Set up metric to use
                      metric = "Accuracy",
                      # Tuning parameters
                      tuneGrid = expand.grid(cp = seq(0, 0.05, 0.01)))

# make predictions on the test set
```

```

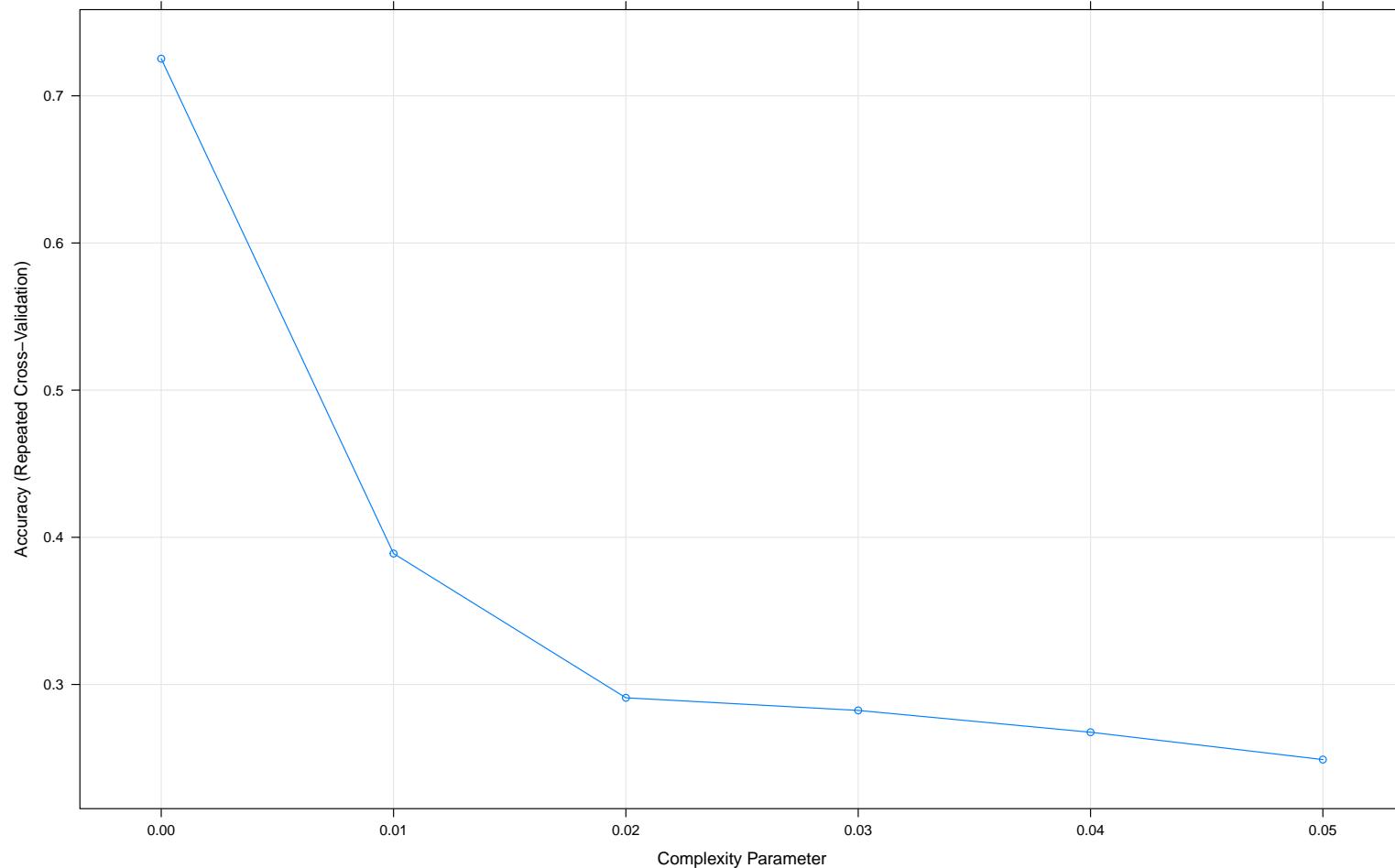
tree_prediction <- predict(tree, newdata = fall_test)

# Generate confusion matrix on the test set
confusionMatrix(tree_prediction, fall_test$activity)

## Confusion Matrix and Statistics
##
##          Reference
## Prediction 3 0 1 2 4 5
##           3 471 140 1 113 189 53
##           0 114 880 3 53 86 49
##           1 21 22 107 71 24 12
##           2 169 83 26 417 68 32
##           4 199 134 5 53 467 111
##           5 103 121 5 44 213 258
##
## Overall Statistics
##
##          Accuracy : 0.5288
## 95% CI : (0.5147, 0.5428)
## No Information Rate : 0.2807
## P-Value [Acc > NIR] : < 2.2e-16
##
##          Kappa : 0.4148
##
## McNemar's Test P-Value : < 2.2e-16
##
## Statistics by Class:
##
##          Class: 3 Class: 0 Class: 1 Class: 2 Class: 4 Class: 5
## Sensitivity      0.43733  0.6377  0.72789  0.55526  0.44604  0.50097
## Specificity       0.87083  0.9138  0.96855  0.90927  0.87028  0.88960
## Pos Pred Value    0.48707  0.7426  0.41634  0.52453  0.48194  0.34677
## Neg Pred Value    0.84658  0.8660  0.99142  0.91897  0.85309  0.93841
## Prevalence        0.21904  0.2807  0.02990  0.15274  0.21293  0.10474
## Detection Rate    0.09579  0.1790  0.02176  0.08481  0.09498  0.05247
## Detection Prevalence 0.19666  0.2410  0.05227  0.16168  0.19707  0.15131
## Balanced Accuracy  0.65408  0.7757  0.84822  0.73226  0.65816  0.69528

```

```
plot(tree)
```



```
#rpart.plot(tree$finalModel)
```

The Random Forest Model

The random forest model does way better than the tree in terms of accuracy with an overall accuracy of 0.6179 against a no information rate of 0.2814, a specificity level of 0.8859, and a sensitivity level of 0.5171.

```
## The random forest model
## Set up the tuning parameters
tunegrid <- expand.grid(.mtry= sqrt(ncol(fall_train)))

## Set up the random forest model
rf_default <- caret::train(activity ~ ., data = fall_train,

    # Set up engine, cross validation and tuning parameters
    method = "rf", tunegrid = tunegrid, trControl = control)

## make predictions on the test set
rf_prediction <- predict(rf_default, newdata = fall_test)

## generate confusion matrix
confusionMatrix(rf_prediction, fall_test$activity)

## Confusion Matrix and Statistics
##
##          Reference
## Prediction   3     0     1     2     4     5
##   3  593    70     1    99   174    36
##   0   94  1119     4    59    75    52
##   1    9     5   105    41    14     5
##   2  162    43    31   491    49    23
##   4  171    87     3    44   597   123
##   5   48    56     3    17   138   276
##
## Overall Statistics
##
##              Accuracy : 0.6469
##                  95% CI : (0.6334, 0.6603)
##      No Information Rate : 0.2807
##      P-Value [Acc > NIR] : < 2.2e-16
##
##              Kappa : 0.5558
##
## McNemar's Test P-Value : 0.0002631
##
## Statistics by Class:
##
##          Class: 3 Class: 0 Class: 1 Class: 2 Class: 4 Class: 5
## Sensitivity      0.5506    0.8109   0.71429   0.65379   0.5702   0.53592
## Specificity       0.9010    0.9197   0.98449   0.92607   0.8894   0.94048
## Pos Pred Value    0.6095    0.7976   0.58659   0.61452   0.5824   0.51301
## Neg Pred Value    0.8773    0.9257   0.99114   0.93686   0.8844   0.94542
## Prevalence        0.2190    0.2807   0.02990   0.15274   0.2129   0.10474
## Detection Rate    0.1206    0.2276   0.02135   0.09986   0.1214   0.05613
## Detection Prevalence 0.1979    0.2853   0.03640   0.16250   0.2085   0.10942
## Balanced Accuracy  0.7258    0.8653   0.84939   0.78993   0.7298   0.73820
```

K-Nearest Neighbours (KNN)

Although the KNN model does not perform as well as the random forest model, it is way better than the tree model. The overall accuracy for the KNN is 0.5189 against a no information rate of 0.2134. The sensitivity is 0.5105 while the specificity is 0.8587. Finally, the balanced accuracy is quite good at 0.6846.

```
## K-Nearest Neighbours (KNN)
## Set up the model and cross validation
knn <- train(activity ~ ., data = fall_train, method = "knn",
             trControl = control)

# Generate the confusion matrix
confusionMatrix(fall_test$activity, predict(knn, newdata = fall_test))

## Confusion Matrix and Statistics
##
##          Reference
## Prediction 3 0 1 2 4 5
##           3 545 44 20 191 176 101
##           0 150 674 29 113 200 214
##           1 1 0 126 17 2 1
##           2 109 12 88 460 41 41
##           4 179 35 37 85 498 213
##           5 50 25 18 20 106 296
##
## Overall Statistics
##
##          Accuracy : 0.5286
##                 95% CI : (0.5145, 0.5426)
##      No Information Rate : 0.2103
##      P-Value [Acc > NIR] : < 2.2e-16
##
##          Kappa : 0.4227
##
## McNemar's Test P-Value : < 2.2e-16
##
## Statistics by Class:
##
##          Class: 3 Class: 0 Class: 1 Class: 2 Class: 4 Class: 5
## Sensitivity          0.5271  0.8532  0.39623  0.51919  0.4868  0.3418
## Specificity          0.8630  0.8289  0.99543  0.92781  0.8590  0.9459
## Pos Pred Value       0.5060  0.4884  0.85714  0.61252  0.4756  0.5748
## Neg Pred Value       0.8727  0.9672  0.95975  0.89774  0.8643  0.8705
## Prevalence           0.2103  0.1607  0.06467  0.18019  0.2081  0.1761
## Detection Rate       0.1108  0.1371  0.02563  0.09355  0.1013  0.0602
## Detection Prevalence 0.2190  0.2807  0.02990  0.15274  0.2129  0.1047
## Balanced Accuracy    0.6950  0.8410  0.69583  0.72350  0.6729  0.6439
```

Extreme Gradient Boosting (XGBoost)

In this section, I run the extreme gradient boosting (XGBoost) model. XGBoost refers to the engineering goal to push the limit of computations resources for boosted tree algorithms.

```

# The XGboost model
## Set up the tuning parameters
tune_grid <- expand.grid(nrounds = 200,
                         max_depth = 5,
                         eta = 0.05,
                         gamma = 0.01,
                         colsample_bytree = 0.75,
                         min_child_weight = 0,
                         subsample = 0.5)

## Set up model, cross validation and tuning parameters
xgb_model <- train(activity ~., data = fall_train, method = "xgbTree",
                     trControl = control,
                     tuneGrid = tune_grid,
                     tuneLength = 10)

## generate confusion matrix for the test set predictions
confusionMatrix(fall_test$activity, predict(xgb_model, newdata = fall_test))

## Confusion Matrix and Statistics
##
##             Reference
## Prediction   3    0    1    2    4    5
##   3 566   91   18  160  172   70
##   0 125   972   17   74   95   97
##   1    0   3 122   19    2    1
##   2 117   54   79  414   57   30
##   4 209   71   42   91  464  170
##   5   51   38   22   20  113  271
##
## Overall Statistics
##
##                 Accuracy : 0.5713
##                 95% CI : (0.5573, 0.5852)
##     No Information Rate : 0.2499
##     P-Value [Acc > NIR] : < 2.2e-16
##
##                 Kappa : 0.4665
##
## McNemar's Test P-Value : < 2.2e-16
##
## Statistics by Class:
##
##             Class: 3 Class: 0 Class: 1 Class: 2 Class: 4 Class: 5
## Sensitivity          0.5300    0.7909    0.40667   0.5321   0.51384   0.42410
## Specificity          0.8672    0.8894    0.99459   0.9186   0.85476   0.94296
## Pos Pred Value       0.5255    0.7043    0.82993   0.5513   0.44317   0.52621
## Neg Pred Value       0.8693    0.9273    0.96268   0.9126   0.88656   0.91640
## Prevalence           0.2172    0.2499    0.06101   0.1582   0.18365   0.12996
## Detection Rate       0.1151    0.1977    0.02481   0.0842   0.09437   0.05511
## Detection Prevalence 0.2190    0.2807    0.02990   0.1527   0.21293   0.10474
## Balanced Accuracy    0.6986    0.8401    0.70063   0.7254   0.68430   0.68353

```

Multinomial Logit Model

In the multinomial logit model below, the overall accuracy is very poor- at 0.2828 against a no information rate of 0.2814. However the model has good specificity at 0.8630 and a very low sensitivity of 0.2970. Consequently, the balanced accuracy is also low at 0.57995. Overall this model lacks good predictive power going by the balanced accuracy and sensitivity relative to the other models. However, it has specificity that is reasonable.

```
## Set up the multinomial logit model
multinom <- multinom(activity ~ ., data = fall_train)

## # weights:  42 (30 variable)
## initial value 34702.797400
## iter 10 value 33136.804144
## iter 20 value 32317.779276
## iter 30 value 31659.552143
## final value 31583.698807
## converged

## Predict the multinomial logit model on the test set
multinom_predict <- predict(multinom, newdata = fall_test)

## Get confusion matrix
confusionMatrix(multinom_predict, fall_test$activity)

## Confusion Matrix and Statistics
##
##          Reference
## Prediction 3 0 1 2 4 5
##       3 319 102 5 92 202 49
##       0 152 771 17 71 113 52
##       1 317 280 92 382 292 145
##       2 88 50 7 65 77 44
##       4 103 54 1 60 148 83
##       5 98 123 25 81 215 142
##
## Overall Statistics
##
##           Accuracy : 0.3126
##           95% CI : (0.2996, 0.3258)
##           No Information Rate : 0.2807
##           P-Value [Acc > NIR] : 4.425e-07
##
##           Kappa : 0.1866
##
##   Mcnemar's Test P-Value : < 2.2e-16
##
## Statistics by Class:
##
##           Class: 3 Class: 0 Class: 1 Class: 2 Class: 4 Class: 5
## Sensitivity      0.29619  0.5587  0.62585  0.08655  0.14136  0.27573
## Specificity      0.88281  0.8855  0.70314  0.93615  0.92222  0.87687
## Pos Pred Value    0.41482  0.6556  0.06101  0.19637  0.32962  0.20760
## Neg Pred Value    0.81726  0.8372  0.98387  0.85041  0.79879  0.91188
## Prevalence        0.21904  0.2807  0.02990  0.15274  0.21293  0.10474
## Detection Rate    0.06488  0.1568  0.01871  0.01322  0.03010  0.02888
## Detection Prevalence 0.15640  0.2392  0.30669  0.06732  0.09132  0.13911
```

```
## Balanced Accuracy      0.58950  0.7221  0.66450  0.51135  0.53179  0.57630
```

Ensemble

Ensemble 1: With multinomial logit In this section, I assemble all the models and use a voting method to build an ensemble. The prediction is the value that receives the most votes from each of the models. In the ensemble, the overall accuracy is 0.5928, a specificity of 0.8804, and a sensitivity of 0.5291. The model does a good job in predicting who will not fall but rather poorly in predicting who will fall. Part of the reason for the low sensitivity maybe the inclusion of the multinomial logit model that had extremely low sensitivity. I remove the multinomial logit model from the ensemble and build a new ensemble next.

```
## make a dataframe with predictions on the test on all the models
ensemble <- tibble(tree = predict(tree, newdata = fall_test), rf = predict(rf_default, newdata = fall_t

## Create a new colum with outcome being the most popular outcome for the models
ensemble$ensemble_all <- apply(ensemble, 1, function(x) {
  ux <- unique(x)
  ux[which.max(tabulate(match(x, ux)))]
})

## Convert the ensembled column into a factor
ensemble$ensemble_all <- factor(ensemble$ensemble_all, levels = levels(ensemble$rf))

## Compute the confusion matrix on the model
confusionMatrix(ensemble$ensemble_all, fall_test$activity)

## Confusion Matrix and Statistics
##
##          Reference
## Prediction   3    0    1    2    4    5
##           3 586  94   1 105 192  41
##           0  87 1035   3  47  67  45
##           1  21  12 118  76  28  17
##           2 164  55  21 455  69  22
##           4 155  92   4  45 512 109
##           5  64  92   0  23 179 281
##
##          Overall Statistics
##
##                Accuracy : 0.6075
##                95% CI : (0.5937, 0.6212)
##                No Information Rate : 0.2807
##                P-Value [Acc > NIR] : < 2.2e-16
##
##                Kappa : 0.5105
##
##    Mcnemar's Test P-Value : < 2.2e-16
##
##    Statistics by Class:
##
##                      Class: 3 Class: 0 Class: 1 Class: 2 Class: 4 Class: 5
##    Sensitivity          0.5441  0.7500  0.80272  0.60586  0.4890  0.54563
##    Specificity          0.8872  0.9296  0.96771  0.92055  0.8953  0.91867
##    Pos Pred Value       0.5751  0.8061  0.43382  0.57888  0.5583  0.43975
```

```

## Neg Pred Value      0.8740  0.9050  0.99376  0.92835  0.8663  0.94530
## Prevalence        0.2190  0.2807  0.02990  0.15274  0.2129  0.10474
## Detection Rate   0.1192  0.2105  0.02400  0.09254  0.1041  0.05715
## Detection Prevalence 0.2072  0.2611  0.05532  0.15985  0.1865  0.12996
## Balanced Accuracy 0.7157  0.8398  0.88522  0.76320  0.6922  0.73215

```

Ensemble 2: Without Multinomial logit The multinomial logit performs poorly and hence I remove it in making the second ensemble. However, both the sensitivity and specificity barely change with this tweak given that the other models subsume the mistakes made by the multinomial logit model.

```

## make a dataframe with predictions on the test on all the models
ensemble2 <- tibble(tree = predict(tree, newdata = fall_test), rf = predict(rf_default, newdata = fall_)

## Create a new column with outcome being the most popular outcome for the models
ensemble2$ensemble_all <- apply(ensemble2, 1, function(x) {
  ux <- unique(x)
  ux[which.max(tabulate(match(x, ux)))]
})

## Convert the ensembled column into a factor
ensemble2$ensemble_all <- factor(ensemble2$ensemble_all, levels = levels(ensemble2$rf))

## Compute the confusion matrix on the model
confusionMatrix(ensemble$ensemble_all, fall_test$activity)

## Confusion Matrix and Statistics
##
##             Reference
## Prediction    3     0     1     2     4     5
##       3  586    94    1  105  192   41
##       0   87  1035    3    47   67   45
##       1   21    12  118    76   28   17
##       2  164    55    21  455   69   22
##       4  155    92     4    45  512  109
##       5   64    92     0    23  179  281
##
## Overall Statistics
##
##           Accuracy : 0.6075
##           95% CI : (0.5937, 0.6212)
##   No Information Rate : 0.2807
##   P-Value [Acc > NIR] : < 2.2e-16
##
##           Kappa : 0.5105
##
## McNemar's Test P-Value : < 2.2e-16
##
## Statistics by Class:
##
##           Class: 3 Class: 0 Class: 1 Class: 2 Class: 4 Class: 5
## Sensitivity      0.5441    0.7500    0.80272   0.60586   0.4890  0.54563
## Specificity      0.8872    0.9296    0.96771   0.92055   0.8953  0.91867
## Pos Pred Value   0.5751    0.8061    0.43382   0.57888   0.5583  0.43975
## Neg Pred Value   0.8740    0.9050    0.99376   0.92835   0.8663  0.94530

```

```

## Prevalence      0.2190  0.2807  0.02990  0.15274  0.2129  0.10474
## Detection Rate 0.1192  0.2105  0.02400  0.09254  0.1041  0.05715
## Detection Prevalence 0.2072  0.2611  0.05532  0.15985  0.1865  0.12996
## Balanced Accuracy 0.7157  0.8398  0.88522  0.76320  0.6922  0.73215

```

Model Evaluation

I present the overall results in the table 8 below. Overall, the random forest model does best in terms of specificity and overall accuracy followed by the ensemble without the multinomial logit model. In terms of sensitivity, the extreme gradient boosting model performs best followed by the ensemble without the multinomial logit model. Ensemble with Multinomial logit model has the highest balanced accuracy. The results are in the table below.

```

## make a table of results
tribble(~ Model, ~ Specificity, ~ Sensitivity, ~ BalancedAccuracy, ~ Accuracy, ~ NoInformationRate, "Cl
      "Random Forest Model", "0.8859", "0.5171", "0.7015", "0.6179", "0.2814", "K Nearest Neighbours"

```

Table 8: Results of the Machine Learning Models

Model	Specificity	Sensitivity	BalancedAccuracy	Accuracy	NoInformationRate
Classification Tree	0.8645	0.4487	0.6566	0.5138	0.2814
Random Forest Model	0.8859	0.5171	0.7015	0.6179	0.2814
K Nearest Neighbours	0.8587	0.5105	0.6846	0.5189	0.2134
Extreme Gradient Boosting	0.8563	0.5437	0.7000	0.5450	0.2484
Multinomial Logit	0.8630	0.2970	0.5800	0.2828	0.2814
Ensemble with Multinomial	0.8815	0.5282	0.7049	0.5930	0.2814
Ensemble without Muntinomial	0.8804	0.5291	0.7048	0.5928	0.2814

Conclusion

In this project, I have built models to predict whether or not a patient falls. I have developed several models- the classification tree, the random forest model, the K-nearest neighbour model, extreme gradient boosting, multinomial logit models and two ensembles- one with the multinomial logit model and one without the multinomial logit model. In evaluating the models I use specificity- the accuracy in predicting that a patient will not fall when in fact, they do not fall as it less expensive to predict that a patient will fall and they do not fall. Overall, the random forest model does the best in terms of specificity, while extreme gradient boosting has the best sensitivity. The challenges I have encountered in this exercise include computing power where the models take too long to run. The computational power has affected my ability to explore more models and fine tune the parameters for better results.

References

- Anava, O., & Levy, K. (2016). k*-nearest neighbours: From global to local. In Advances in neural information processing systems (pp. 4916-4924).
- Chen, Z., & Fan, W. D. (2019). A multinomial logit model of pedestrian-vehicle crash severity in North Carolina. International journal of transportation science and technology, 8(1), 43-52.
- Özdemir, A. T., & Barshan, B. (2014). Detecting falls with wearable sensors using machine learning techniques. Sensors, 14(6), 10691-10708.
- Shaikhina, T., Lowe, D., Daga, S., Briggs, D., Higgins, R., & Khovanova, N. (2019). Decision tree and random forest models for outcome prediction in antibody incompatible kidney transplantation. Biomedical Signal Processing and Control, 52, 456-462.