

Classifying patients by analyzing the biomechanical features of orthopedic patients

Sadia Boksh

01/03/2021

Introduction

In this project we will use the data set from Kaggle to study the biochemical features of orthopedic patients and classify the patients based on the features. Using this data set, we will train some machine learning models to classify patients as belonging to one out of three categories: Normal, Disk Hernia or Spondylolisthesis. We will perform the model fitting on scaled raw data. We will choose the best performing model by analyzing their accuracy.

Methods

Data Exploration

The data set used in this project can be found in <https://www.kaggle.com/uciml/biomechanical-features-of-orthopedic-patients>. This data set has 310 rows and 7 columns. There are 6 features and 1 response variable. There are no missing values in the data set.

```
head(df)
```

```
##      pelvic_incidence pelvic_tilt lumbar_lordosis_angle sacral_slope pelvic_radius
## 1          63.02782    22.552586          39.60912      40.47523         98.67292
## 2          39.05695    10.060991          25.01538      28.99596        114.40543
## 3          68.83202    22.218482          50.09219      46.61354        105.98514
## 4          69.29701    24.652878          44.31124      44.64413        101.86850
## 5          49.71286     9.652075          28.31741      40.06078        108.16872
## 6          40.25020    13.921907          25.12495      26.32829        130.32787
##      degree_spondylolisthesis class
## 1          -0.254400 Hernia
## 2           4.564259 Hernia
## 3          -3.530317 Hernia
## 4          11.211523 Hernia
## 5           7.918501 Hernia
## 6           2.230652 Hernia
```

This data set has three categories in response variable. The patients are to be classified into these three categories:

```
## [1] "Hernia"          "Spondylolisthesis" "Normal"
```

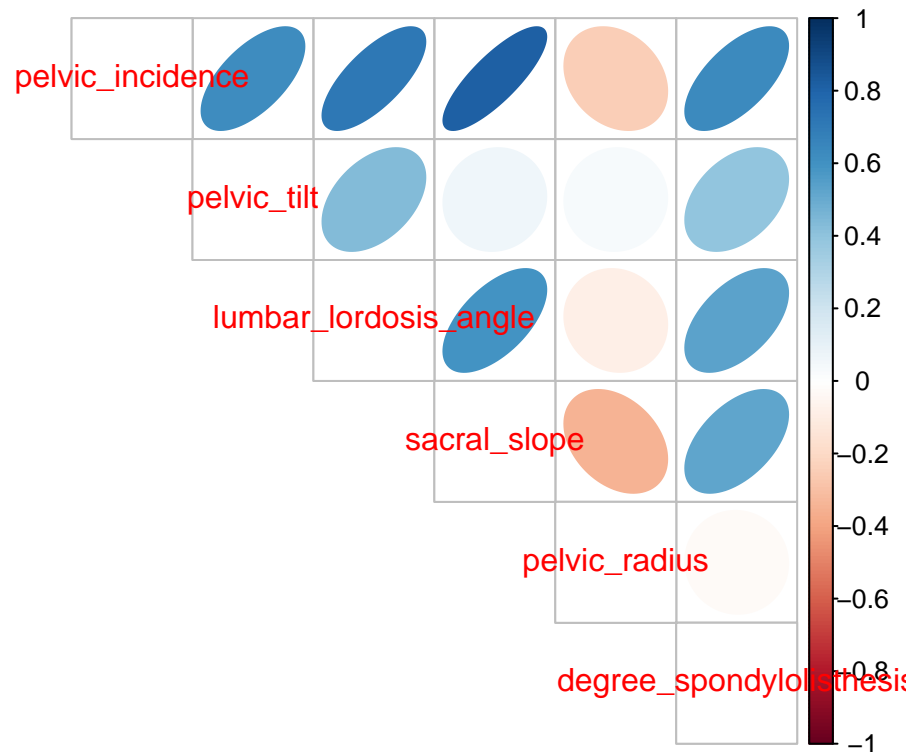
Below is the table that shows the proportion of patients in each class:

```
##           Class      Prop
## 1         Hernia 0.1935484
## 2         Normal 0.3225806
## 3 Spondylolisthesis 0.4838710
```

Plots

Below is the correlation plot:

```
##           pelvic_incidence pelvic_tilt lumbar_lordosis_angle
## pelvic_incidence           1.0000000  0.62919877          0.71728236
## pelvic_tilt                0.6291988  1.00000000          0.43276386
## lumbar_lordosis_angle       0.7172824  0.43276386          1.00000000
## sacral_slope                0.8149600  0.06234529          0.59838689
## pelvic_radius               -0.2474672  0.03266781         -0.08034361
## degree_spondylolisthesis    0.6387427  0.39786228          0.53366701
##
##           sacral_slope pelvic_radius degree_spondylolisthesis
## pelvic_incidence    0.81495999 -0.24746721          0.63874275
## pelvic_tilt          0.06234529  0.03266781          0.39786228
## lumbar_lordosis_angle 0.59838689 -0.08034361          0.53366701
## sacral_slope         1.00000000 -0.34212835          0.52355746
## pelvic_radius        -0.34212835  1.00000000         -0.02606501
## degree_spondylolisthesis 0.52355746 -0.02606501          1.00000000
```



From the plot above we can see pelvic_incidence is highly correlated with sacral_slope.

Principle Component Analysis

We will apply PCA to explore the variable importance of each feature. Using the summary function we can see the variability explained by each PC:

```
#transform to a matrix
x <- df[, 1:6] %>% as.matrix()

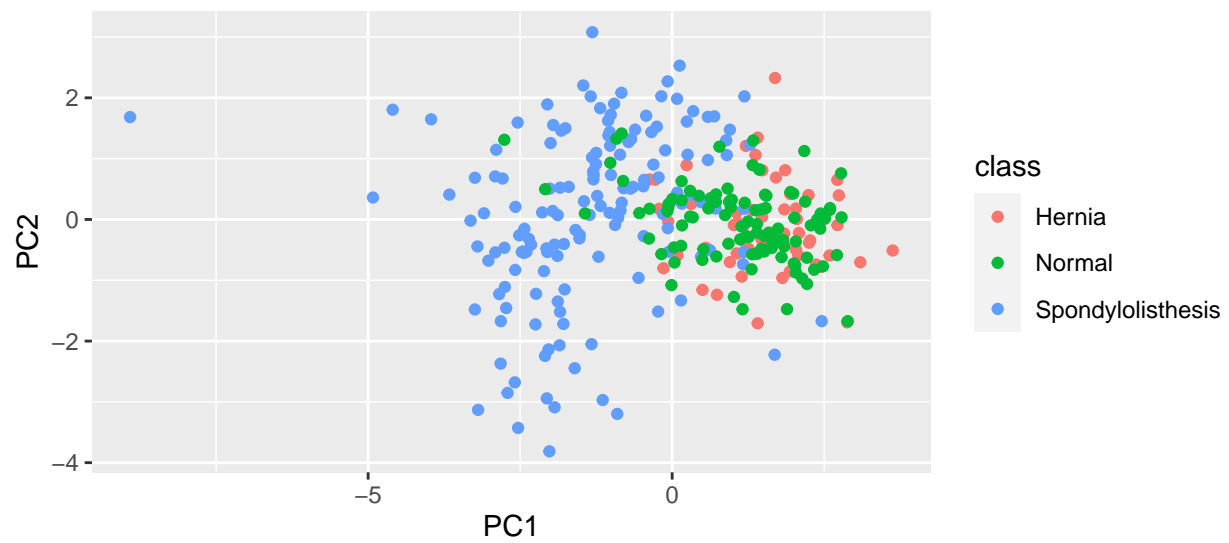
# scale and center the feature matrix
x_centered <- sweep(x, 2, colMeans(x))
scaled_X <- sweep(x_centered, 2, colSds(x), FUN = "/")

# principal components
pca <- prcomp(scaled_X)
summary(pca)$importance
```

```
##                PC1      PC2      PC3      PC4      PC5
## Standard deviation    1.801605 1.09297 0.872405 0.6874067 0.5709795
## Proportion of Variance 0.540960 0.19910 0.126850 0.0787500 0.0543400
## Cumulative Proportion 0.540960 0.74006 0.866910 0.9456600 1.0000000
##
##                PC6
## Standard deviation    1.935122e-10
## Proportion of Variance 0.000000e+00
## Cumulative Proportion 1.000000e+00
```

We can plot the first two PCS to see how they explain the variability:

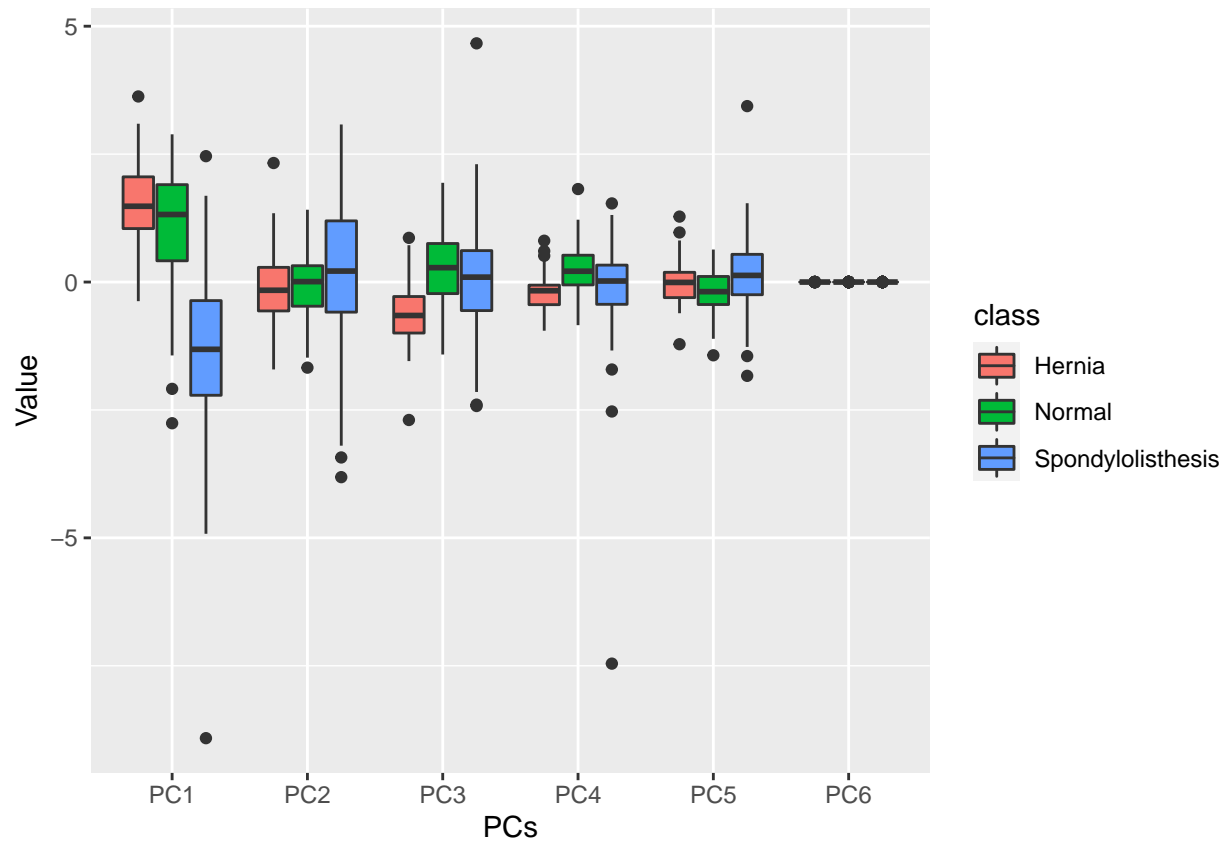
```
data.frame(pca$x[,1:2], class=df$class) %>%
ggplot(aes(PC1,PC2, col = class))+
geom_point() +
coord_fixed(ratio = 1)
```



We can see PC1 and PC2 has separated the patients into two categories: Spondylolisthesis and non Spondylolisthesis. Lower PC1 explains Spondylolisthesis and higher PC1 explains either Normal or Hernia.

We can also plot the first 10 PCs:

```
data.frame(pca$x[,1:6], class=df$class) %>% gather(PCs,Value, -class) %>%
ggplot(aes(PCs,Value, fill = class))+
geom_boxplot()
```



From the plot above we can see PC1 is not overlapping with other PCs.

Modelling

Now We will fit LDA, KNN and Random forest, SVM Linear models to the scaled data set and compare their accuracy. First we will split the scaled data set to 80% train set and 20% test set.

LDA

```
set.seed(5, sample.kind = "Rounding")

train_lda <- train(train_x, train_y, method = "lda")
pred_lda <- predict(train_lda, test_x)
acc_lda <- confusionMatrix(pred_lda, test_y)$overall['Accuracy']
acc_lda
```

```
## Accuracy
## 0.8064516
```

K Nearest Neighbours

For KNN, I am using tuning parameter k from 15 to 40 and the default cross validation is performed by taking 25 bootstrap samples comprised of 25% of the observations

```
set.seed(7, sample.kind = "Rounding")

train_knn <- train(train_x, train_y, method = "knn", tuneGrid = data.frame(k=c(15:40,2)))
pred_knn <- predict(train_knn, test_x)
acc_knn <- confusionMatrix(pred_knn,test_y)$overall['Accuracy']

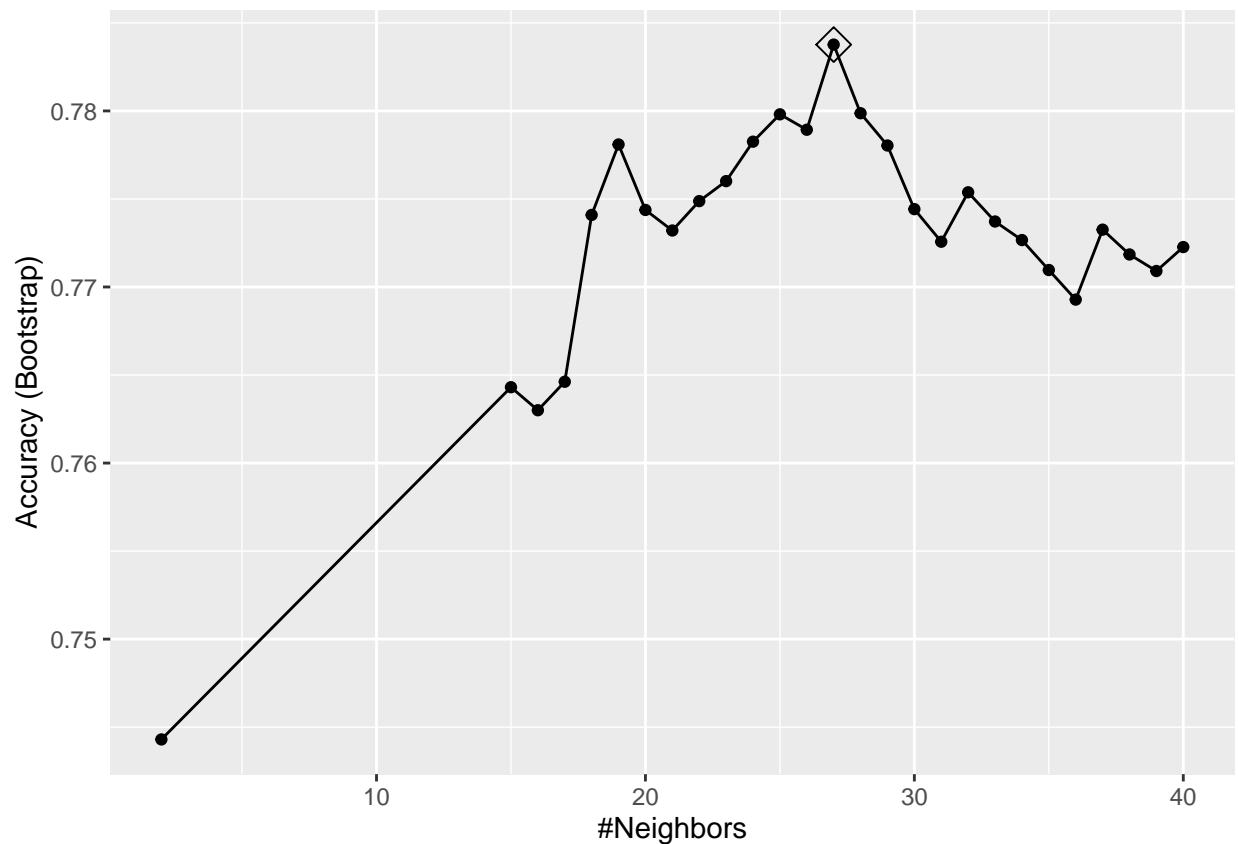
acc_knn
```

```
## Accuracy
## 0.8064516
```

```
train_knn$bestTune
```

```
##      k
## 14 27
```

```
ggplot(train_knn, highlight = TRUE)
```



SVM Linear Model

For SVM Linear model, I have used tuning parameter C from 1 to 10 and 10 fold cross validation.

```

set.seed(20, sample.kind = "Rounding")
train_control <- trainControl(method="repeatedcv", number=10, repeats=3)

train_svm <- train(train_x, train_y, method = "svmLinear", tuneGrid = data.frame(C=c(1:10,2)), trControl=train_control)
pred_svm <- predict(train_svm, test_x)
acc_svm <- confusionMatrix(pred_svm, test_y)$overall['Accuracy']

acc_svm

```

```

## Accuracy
## 0.8709677

```

```

train_svm$bestTune

```

```

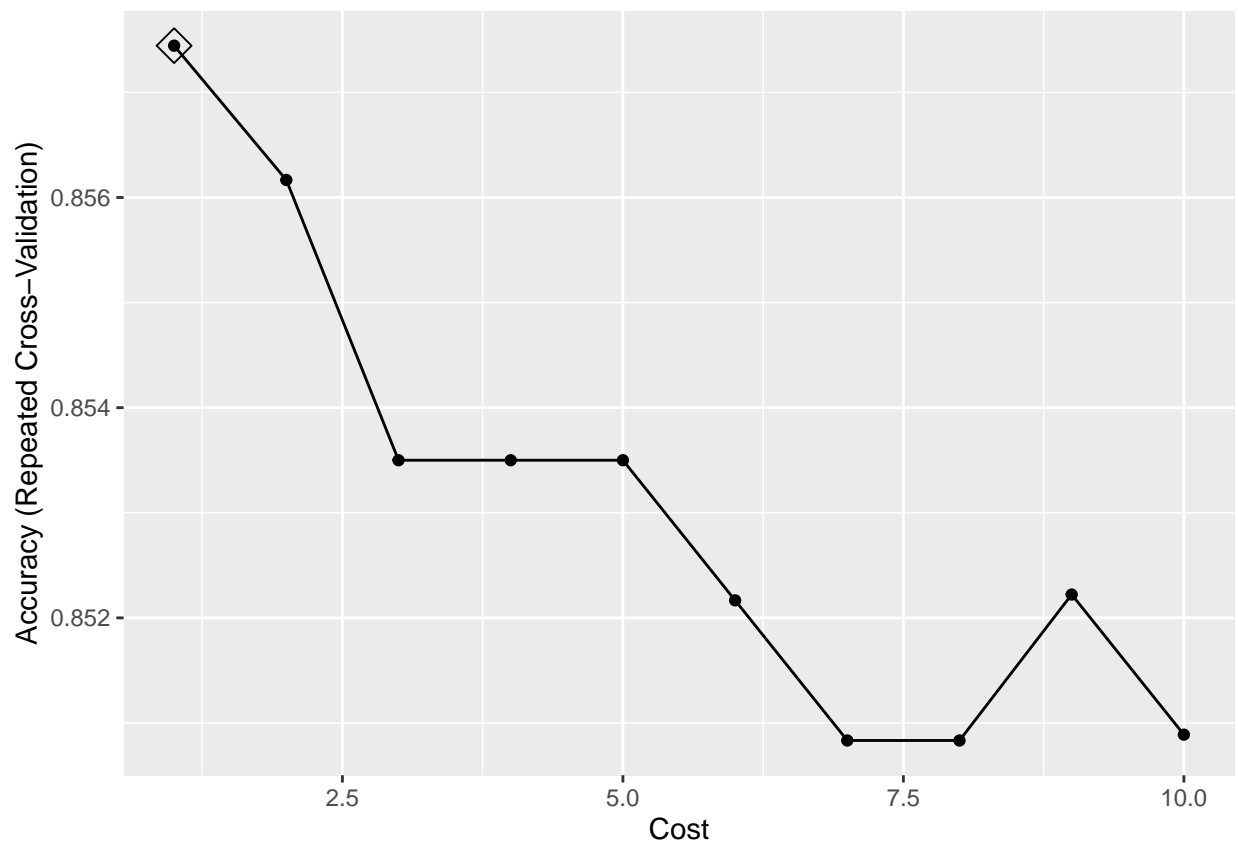
## C
## 1 1

```

```

ggplot(train_svm, highlight = TRUE)

```



Random Forest

For Random forest, tune grid parameter is mtry with values from 3 to 13.

```

set.seed(9, sample.kind = "Rounding")

train_rf <- train(train_x, train_y, method = "rf", tuneGrid = data.frame(mtry=c(3, 5, 7, 9, 11, 13)), in
pred_rf <- predict(train_rf, test_x)
acc_rf <- confusionMatrix(pred_rf, test_y)$overall['Accuracy']
acc_rf

```

```

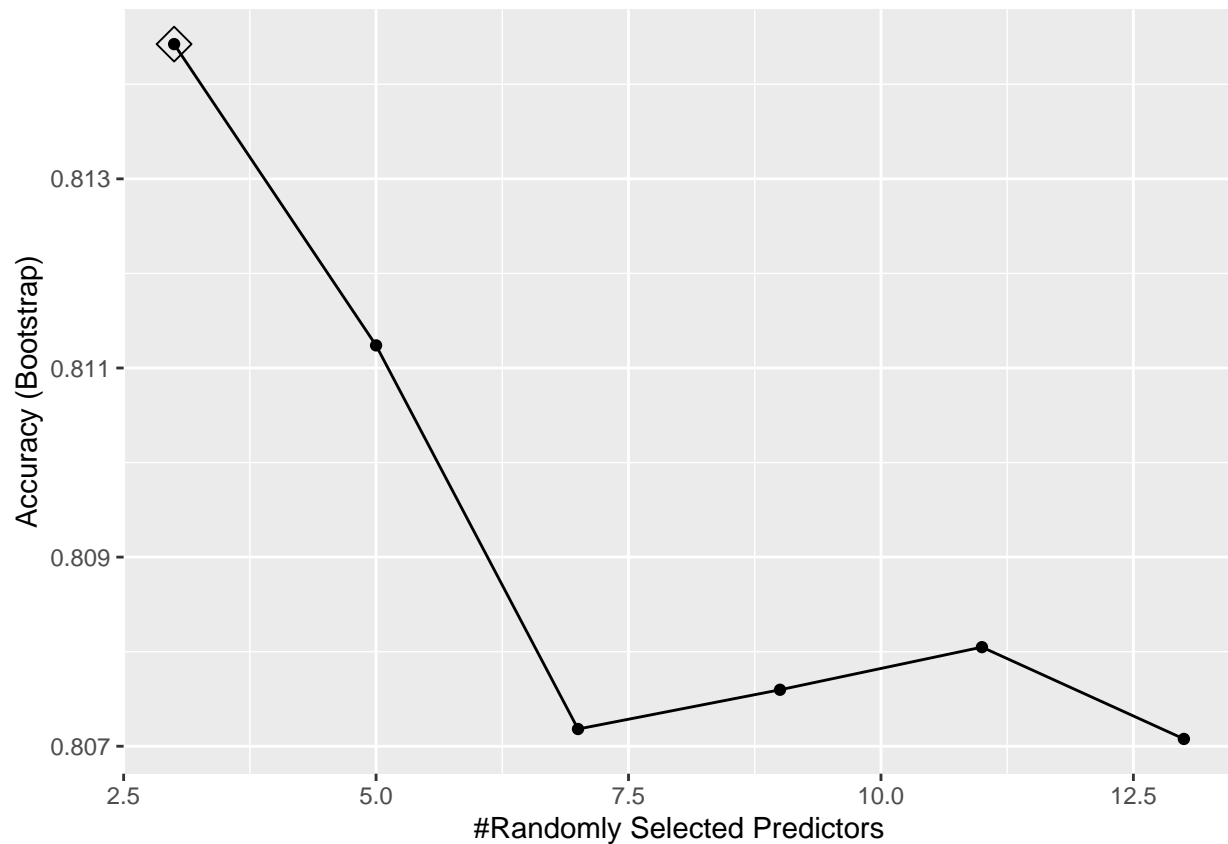
## Accuracy
## 0.8548387

```

```

ggplot(train_rf, highlight = TRUE)

```



```

## rf variable importance
##
##   variables are sorted by maximum importance across the classes
##
##           Hernia Normal Spondylolisthesis
## degree_spondylolisthesis 53.835 76.668      100.000
## pelvic_radius             6.613 39.122       12.049
## sacral_slope              29.818  2.236        9.504
## pelvic_tilt                6.539 22.305       11.707
## pelvic_incidence          10.589 14.748       16.734
## lumbar_lordosis_angle     15.302  0.000        8.624

```


Results

Now we can compare the results of different models and their accuracy.

Below is the accuracy table summary:

| ## | Method | Accuracy |
|------|------------|-----------|
| ## 1 | lda | 0.8064516 |
| ## 2 | knn | 0.8064516 |
| ## 3 | rf | 0.8548387 |
| ## 4 | svm_linear | 0.8709677 |

Conclusion

In summary, this analysis shows it is possible to classify the orthopedic patients by analyzing their biochemical features. SVM Linear is the highest performing model with accuracy around 87%. Future work can be done to improve the accuracy above 87%.