

# Course Project for Practical Machine Learning

This project reads in the training data from people who did the weight lifting exercise and try to predict if people are doing exercise correctly in the test data.

1. Load the training data. As we believe that information like the data index and the time of recording is not relevant to if people are doing actions correctly, we discard them. Additionally, there are a lot of features with mostly NA or blank values, we discard them as well. So we manually select the remaining 54 features which are non-NA to use below.

```
rm(list=ls())  
library(caret)
```

```
## Loading required package: lattice  
## Loading required package: ggplot2
```

```
data <- read.csv('pml-training.csv')  
  
# select the non-NA features  
colIndex <- c(7,8,9,10,11,37:49,60:68,84:86,102,113:124,140,151:160)  
data <- data[,colIndex]
```

2. For training the model and cross validation, we randomly split the data into 70% training and 30% test data.

```
# split the training data into training and test for cross validation  
set.seed(98123)  
inTrain <- createDataPartition(y=data$classe, p=0.7, list=FALSE)  
trainData <- data[inTrain,]  
testData <- data[-inTrain,]  
  
dims <- dim(trainData)  
col <- dims[2]
```

3. Perform PCA to reduce the dimension of the training data and apply PCA model on

```
# do PCA to reduce the dimension
preProc <- preProcess(trainData[, -col], method="pca", thresh=0.9)
trainPC <- predict(preProc, trainData[, -col])
testPC <- predict(preProc, testData[, -col])

# add the class label to train and test data
dims <- dim(trainPC)
colPC <- dims[2]

newTrainPC <- cbind(trainPC, trainData[, col])
names(newTrainPC)[colPC+1] <- "classe"
newTestPC <- cbind(testPC, testData[, col])
names(newTestPC)[colPC+1] <- "classe"
```

4. Prepare the test data according to the preprocess steps done for training data:

```
# prepare validation data according to training data
validationData <- read.csv('pml-testing.csv')
validationData <- validationData[, colIndex]
validationPC <- predict(preProc, validationData[, -col])
```

5. Train the model on training data and validate the model on test data.

```
modelFit <- train(classe ~., data=newTrainPC, method="rf")
```

```
## Loading required package: randomForest
## randomForest 4.6-10
## Type rfNews() to see new features/changes/bug fixes.
```

```
predTrain <- predict(modelFit, newTrainPC)
predtest <- predict(modelFit, newTestPC)

confTrain <- confusionMatrix(trainData[, col], predTrain)
conftest <- confusionMatrix(testData[, col], predtest)
```

The accuracy of the model on the training dataset is:

```
train_accuracy <- as.numeric(confTrain$overall["Accuracy"])
```

The in sample error is:

```
in_sample_error <- (1-train_accuracy)
in_sample_error
```

```
## [1] 0
```

The accuracy of the model on the test dataset is:

```
test_accuracy <- as.numeric(conftest$overall["Accuracy"])
```

The out of sample error is expected at:

```
out_of_sample_error <- (1-test_accuracy)
out_of_sample_error
```

```
## [1] 0.01852167
```

So we can see that the in sample error is lower than the expected out of sample error, which is normal for prediction models.

6. Finally we use the model to predict the testcases which is the validation data:

```
predValidation <- predict(modelFit, validationPC)
predValidation
```

```
## [1] B A A A A E D B A A B C B A E E A B B B
## Levels: A B C D E
```

From the submission results, our model got 19 out of 20 test cases correct, yielding a 95% accuracy on validation data. So the real out of sample error is bigger than the expectation.