

Classe Predictor

Nazar Mammedov

2025-01-12

Introduction

Create training and test partitions. We create internal train and test partitions from the data provided in “pml-training.csv” because we want to be able to test our model.

```
library(caret)
```

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

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

```
library(ggplot2)
library(e1071)
library(dplyr)
```

```
##
```

```
## Attaching package: 'dplyr'
```

```
## The following objects are masked from 'package:stats':
```

```
##
```

```
##      filter, lag
```

```
## The following objects are masked from 'package:base':
```

```
##
```

```
##      intersect, setdiff, setequal, union
```

```
data = read.csv("pml-training.csv")
external_test = read.csv("pml-testing.csv")
```

Create a new dataset with proper columns both for training and external test. We drop columns with a lot of NAs because they are not going to be useful anyway, and also drop some variables such as timestamps.

```
# Drop columns with any NA values
```

```
data_no_na <- data %>% select_if(~ !any(is.na(.)))
```

```
test_no_na <- external_test %>% select_if(~ !any(is.na(.)))
```

```
#Select only numeric columns
```

```
data_numeric <- data_no_na %>% select_if(is.numeric) %>% bind_cols(data_no_na %>% select(classe))
```

```

test_numeric <- test_no_na %>% select_if(is.numeric)

# drop first 4 columns X, timestamps, and num_window
# because I don't think they are good predictors anyway
new_data <- data_numeric[, -c(1:4)]
new_data$classe <- as.factor(new_data$classe)
new_test_data <- test_numeric[, -c(1:4)]
new_test_data <- new_test_data %>% select(-problem_id)

```

Now given this data set. We train the model with SVM linear. It gives .7856 accuracy level which is acceptable for our purposes.

```

# Split the data into training and testing sets (70/30 split)
set.seed(123)
inTrain <- createDataPartition(new_data$classe, p = 0.7, list = FALSE)
train_data <- new_data[inTrain, ]
test_data <- new_data[-inTrain, ]

# Build a classification model (e.g., SVM)
model_all <- svm(classe ~ ., data = train_data, kernel = "linear")

# Make predictions and evaluate the model
predictions_all <- predict(model_all, test_data)
cm <- confusionMatrix(predictions_all, test_data$classe)
cm

```

Confusion Matrix and Statistics

```

##
##           Reference
## Prediction    A    B    C    D    E
##           A 1557  149  109   63   60
##           B   28  837   83   42  155
##           C   34   65  796  112   85
##           D   43   18   23  703   52
##           E   12   70   15   44  730
##

```

Overall Statistics

```

##
##           Accuracy : 0.7856
##           95% CI : (0.7748, 0.796)
##           No Information Rate : 0.2845
##           P-Value [Acc > NIR] : < 2.2e-16
##

```

```

##
##           Kappa : 0.7271
##

```

McNemar's Test P-Value : < 2.2e-16

##

Statistics by Class:

```

##
##           Class: A Class: B Class: C Class: D Class: E
## Sensitivity          0.9301   0.7349   0.7758   0.7293   0.6747
## Specificity          0.9095   0.9351   0.9391   0.9724   0.9706
## Pos Pred Value       0.8034   0.7310   0.7289   0.8379   0.8381

```

```
## Neg Pred Value      0.9704  0.9363  0.9520  0.9483  0.9298
## Prevalence          0.2845  0.1935  0.1743  0.1638  0.1839
## Detection Rate      0.2646  0.1422  0.1353  0.1195  0.1240
## Detection Prevalence 0.3293  0.1946  0.1856  0.1426  0.1480
## Balanced Accuracy    0.9198  0.8350  0.8575  0.8508  0.8227
```

Out of sample error rate

Out of sample error rate is (1-Accuracy rate), but can be calculated as below. Its value is 0.2144.

```
conf_matrix <- cm$table
misclassifications <- sum(conf_matrix) - sum(diag(conf_matrix))
total_observations <- sum(conf_matrix)
misclassification_rate <- misclassifications / total_observations
print(misclassification_rate)
```

```
## [1] 0.2144435
```

Predicting 20 cases

Now we predict “classe” in the external test data which was provided in “pml-testing.csv”. As predicted this gives about 80% of accuracy, which was also the expected Quiz result.

```
predictions_test <- predict(model_all, new_test_data)
print(predictions_test)
```

```
##  1  2  3  4  5  6  7  8  9 10 11 12 13 14 15 16 17 18 19 20
##  C  A  B  C  A  E  D  B  A  A  C  A  B  A  E  E  A  B  B  B
## Levels: A B C D E
```