

Predicting the type of movement based on accelerator data

```
## Warning: package 'caret' was built under R version 3.1.2
```

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

```
## Warning: package 'AppliedPredictiveModeling' was built under R version  
## 3.1.2
```

```
## Warning: package 'kernlab' was built under R version 3.1.2
```

```
## Warning: package 'e1071' was built under R version 3.1.2
```

```
## Warning: package 'randomForest' was built under R version 3.1.2
```

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

executive summary

This report presents the build-up of a model based on a training set of several person exercising to predict the type of movement they are doing based on accelerometer data. The data used in this analysis is from groupware project on Human Activity Recognition [1]. The outcome is the column “classe” which is a factor of 5 levels, from A to E corresponding to a type of activity. The model chosen for this prediction was a random forest prediction without preprocessing of the data except some simple data cleaning. The accuracy of the model was around 99% on a subset of the training set for testing.

Data exploration and data cleaning

The objective is to predict the “classe” parameter based on the other metrics. Looking at the data set, we have 160 columns, any of them are statistically summary of more detailed parametrics taking during a time widow.

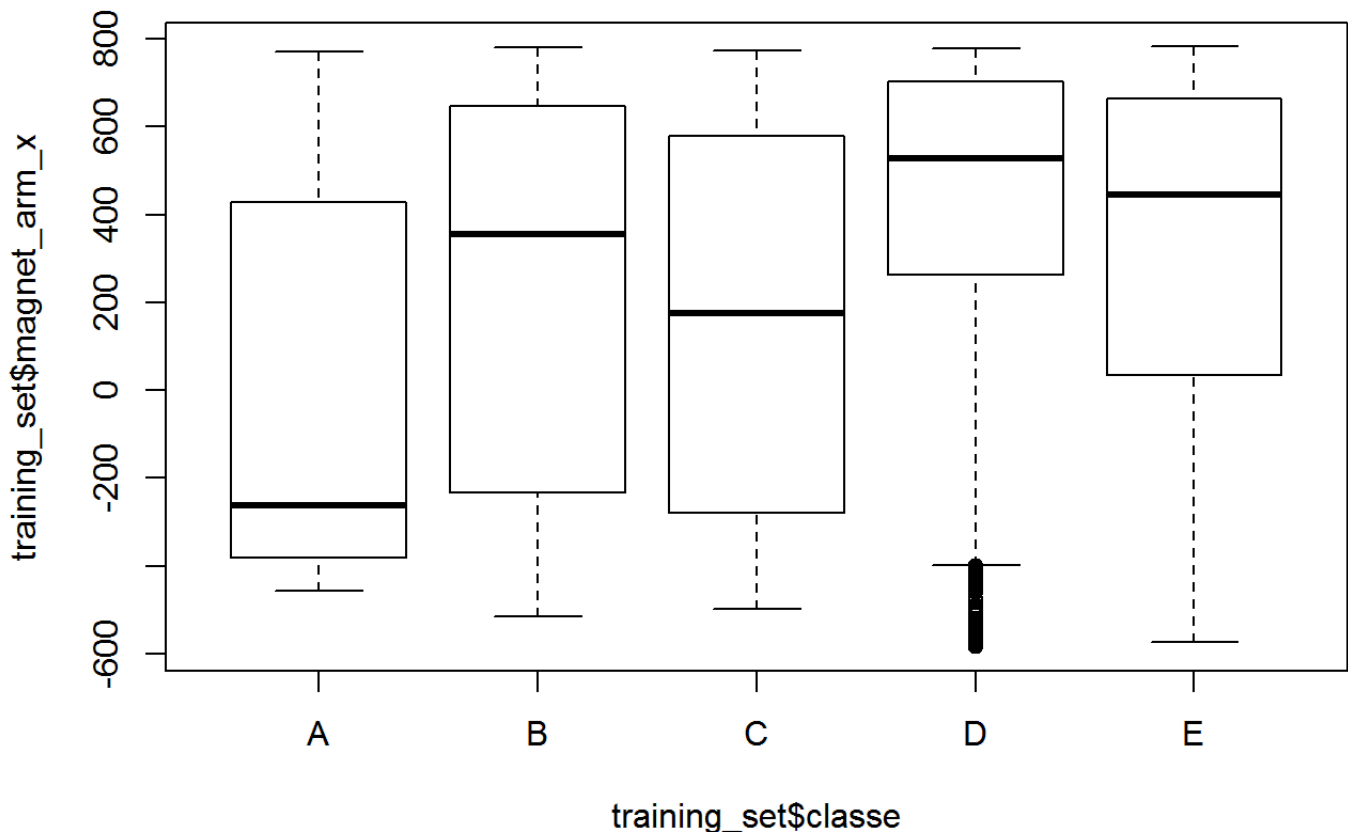
```
training <- read.csv("pml-training.csv",header = TRUE,sep = ",", quote = "\"",na.strings = "NA",dec  
= "." )  
testing <- read.csv("pml-testing.csv",header = TRUE,sep = ",", quote = "\"",na.strings = "NA",dec =  
"." )
```

The first step is to remove the window summary data since it is not defined for the test step. In addition, all the non-critical parameters not correlated to the activity are removed such as time, window numbers. The resulting data set has 52 numerical columns.

```
#data preparation
index <- (training$new_window == "yes")
training_set <- training[-index,]
# integrated parameters
training_ave <- training[index,]
colremove <- c(1:7,11:36,50:59,69:83,87:101,103:112,125:139,141:150)
training_set <- training_set[,-colremove]
testing_set <- testing[,-colremove]
```

For first level analysis, a box plot of the different parameters as a function of the classe outcome can help to detect any issues with the data. As an example, the magnet_arm_x was plotted as a function of the classes.

```
plot(training_set$magnet_arm_x ~ training_set$classe)
```

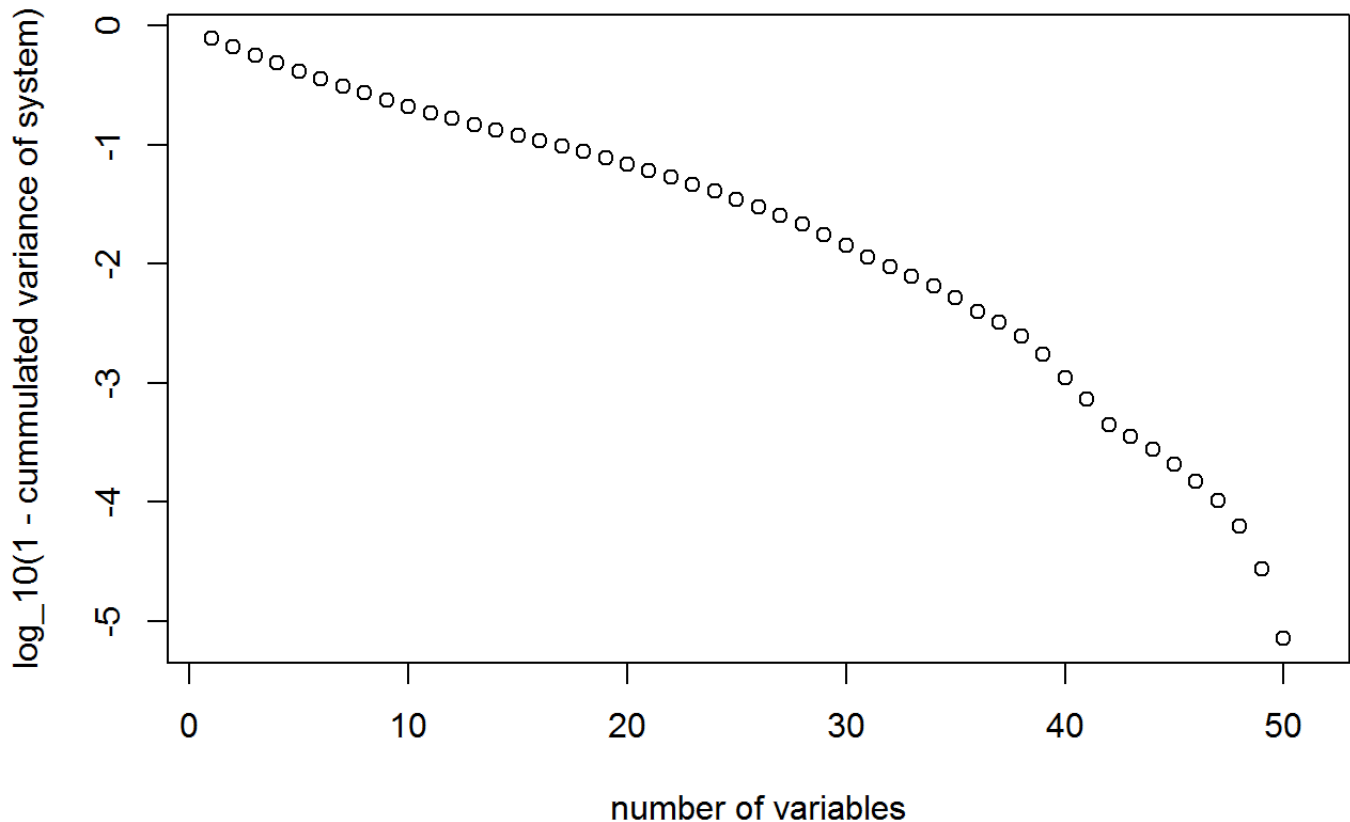


Preprocessing of data

Due to the number of numerical variables, a singular value decomposition was tried to see if we can reduce the number of regressor. the cumulated

The following plot of the residual variance as a function of the number of variables.

```
plot(log10(1.0-cummul_d), ylab = "log_10(1 - cummulated variance of system)",xlab = "number of variables")
```



The decay of the diagonal factor is very slow. To maintain 99% of the variance of the system, we need to include up 30 parameters. this suggests that principal component analysis method might not provide an improvement for the model. Other simple preproceesing methods were used and tested to see if the model accuracy ws improved. Unfortunately, no simple preprocessing was found effective.

Model development.

The training set is split between a training set and a test set in order to evaluate the accuracy of the modeland to pick the best model. 75% of the set is used to train the model.

The outcome is a factor so we need to choice a model for categorization. A standard linear regression model cannot be used. A random forest method was tested and seems to provide the best accuracy.

```
modelFit <- randomForest(classe ~ ., data = training_d)
```

Different method of pre-processing the data was tested (pca, BoxCox) to see if the accuracy of the model could be improved.

verification of the accuracy of the model

The training set split is used to evaluate the accuracy of the model. This training set has seen the same preprocess. The outcome of the prediction using the training set is compared to the actual values.

```
confusionMatrix(training_t$classe,predict(modelFit,training_t))
```

```
## Confusion Matrix and Statistics
##
##           Reference
## Prediction  A   B   C   D   E
##           A 557   0   0   0   0
##           B   0 379   0   0   0
##           C   0   2 340   0   0
##           D   0   0   3 318   0
##           E   0   0   0   1 359
##
## Overall Statistics
##
##           Accuracy : 0.9969
##           95% CI : (0.9933, 0.9989)
##           No Information Rate : 0.2843
##           P-Value [Acc > NIR] : < 2.2e-16
##
##           Kappa : 0.9961
##           McNemar's Test P-Value : NA
##
## Statistics by Class:
##
##           Class: A Class: B Class: C Class: D Class: E
## Sensitivity           1.0000   0.9948   0.9913   0.9969   1.0000
## Specificity           1.0000   1.0000   0.9988   0.9982   0.9994
## Pos Pred Value        1.0000   1.0000   0.9942   0.9907   0.9972
## Neg Pred Value        1.0000   0.9987   0.9981   0.9994   1.0000
## Prevalence            0.2843   0.1945   0.1751   0.1628   0.1833
## Detection Rate        0.2843   0.1935   0.1736   0.1623   0.1833
## Detection Prevalence  0.2843   0.1935   0.1746   0.1639   0.1838
## Balanced Accuracy      1.0000   0.9974   0.9950   0.9975   0.9997
```

```
#confusionMatrix(training_t$classe,predict(modelFit,testPC))
modelFit$finalModel
```

```
## NULL
```

The accuracy of the model is around 99%.

predicting the test set

the final step is to use our current model for predicting the values on a new set of data.

the result of the prediction is given by

```
predict(modelFit,testing_set)
```

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

reference

[1] Ugulino, W.; Cardador, D.; Vega, K.; Velloso, E.; Milidui, R.; Fuks, H. Wearable Computing: Accelerometers' Data Classification of Body Postures and Movements. Proceedings of 21st Brazilian Symposium on Artificial Intelligence. Advances in Artificial Intelligence - SBIA 2012. In: Lecture Notes in Computer Science. , pp. 52-61. Curitiba, PR: Springer Berlin / Heidelberg, 2012. ISBN 978-3-642-34458-9. DOI: 10.1007/978-3-642-34459-6_6