# Machine Learning: Assignment 2

Karsten Standnes - STNKAR012

October 2017

## 1   Introduction

In the world today there are many methods that fall under the category of "Machine Learning", some being very simular and some very different. All of them share in common that they use data to produce a model that can be used on unseen data. When sucessful this is very appeling in todays society when we got lots of data on situations where there is likely to be a underlying pattern, another requirement for learning. There is broad agreement that machine learning is a good way to make prediction and classification models, and often the only computational feasible way. This makes it a task to decide which method in machine learning to chose for a given problem. The answer is not always the same and several methods can be good, but for different reasons. In this task I will show several different machine learning algorithms and how they perform on classifying pictures of handwritten numbers.
COMMENT on data snooping

## 2   Methods

### 2.1   Tree based methods

#### 2.1.1   Classification tree

Classification tree is maybe the method in Machine Learning which is easiest to interpret due to it's intuitive construction and clear visualization. The method uses a greedy approach using recursive binary splitting to structure a tree that can classify input based on it's variables. The goal of the classification tree is to classify a set of data as best as possible while have a low complexity to avoid overfitting. Below we see that one classification tree is not enough to make a great model for the problem, but combining many of them gives us a "Random Forest" which is discussed in subsection 2.1.2. Classification trees also gives a nice visual image of the classification.

#### 2.1.2   Random Forest

Random Forest is as mentioned (subsection 2.1.1) constructed of many classification tree. After such a construction by using a set of training data, new data can be ran through the "forest". The new data is classified with the label the majority of trees labelled it. The trees in other words "vote" for the best classification for the data. As in real life, voting makes little difference if all the votes are the same. To avoid making a forest out of $n$ ($n \in \mathbb{Z}_{>0}$) identical trees two things are changed in the construction of the tree from the classification tree. One is that a the training data for each tree is a sample of size $N$ from the whole set of $N$ data sampled with replacement. The other is that $m$ randomly selected variables are used for each tree, where $m << M$ and $M$ is the whole set of variables in the problem. This gives a rich variety of trees where with enough trees a predictive model can be created.

## 2.2 Neural Networks

For Neural Networks I have decided to use a convulutional neural network because of it great performance in problems that can be represented as an "image". This because it convolutional neural networks(CNN's) works by recognizing paterns in images or represented as a matrix or tensor(multidimension matrix) in the computer. It does so by using convolutional, pooling and voting layers. Multiple of these can be stacked to create a very a precise recongising images.
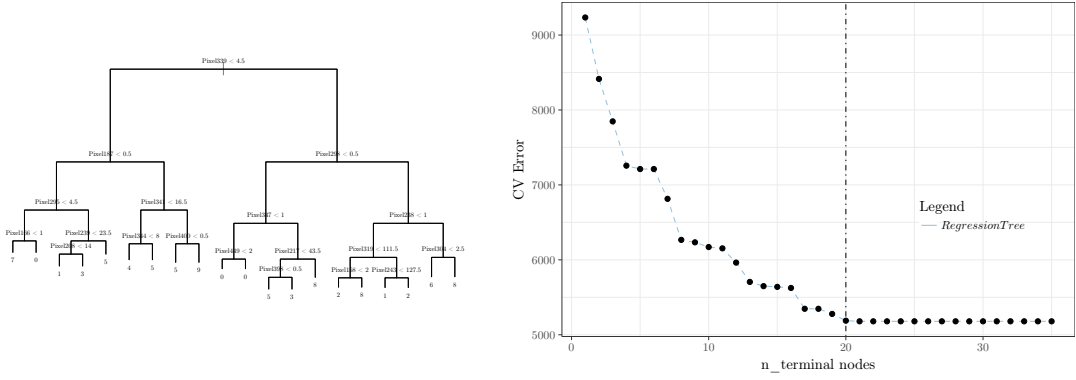
Figure 1: Error measures $E_{out}$ and $E_{val}$ are plotted for different sizes of training and validation sets.

## 2.3 Support vector machines

# 3 Results

## 3.1 Tree based methods

### 3.1.1 Classification tree



(a) $E_{CV}$ for different amounts of leaf nodes.

(b) Regression tree after pruning

Figure 2: A figure with two subfigures

|       | 0  | 1  | 2  | 3  | 4  | 5  | 6  | 7  | 8  | 9  | Error | Rate    |
|-------|----|----|----|----|----|----|----|----|----|----|-------|---------|
| 0     | 31 | 0  | 5  | 1  | 1  | 3  | 4  | 1  | 0  | 1  | 0.340 | 16/47   |
| 1     | 0  | 44 | 1  | 2  | 0  | 2  | 3  | 0  | 1  | 0  | 0.170 | 9/53    |
| 2     | 0  | 3  | 32 | 2  | 1  | 0  | 0  | 0  | 2  | 0  | 0.200 | 8/40    |
| 3     | 0  | 1  | 1  | 26 | 1  | 0  | 1  | 1  | 0  | 3  | 0.235 | 8/34    |
| 4     | 0  | 0  | 0  | 0  | 26 | 0  | 10 | 0  | 1  | 1  | 0.316 | 12/38   |
| 5     | 4  | 2  | 2  | 12 | 9  | 31 | 6  | 10 | 3  | 7  | 0.640 | 55/86   |
| 6     | 0  | 1  | 3  | 1  | 2  | 1  | 23 | 0  | 0  | 2  | 0.303 | 10/33   |
| 7     | 1  | 0  | 1  | 2  | 0  | 0  | 2  | 34 | 0  | 5  | 0.244 | 11/45   |
| 8     | 0  | 2  | 5  | 3  | 6  | 5  | 6  | 0  | 39 | 3  | 0.435 | 30/69   |
| 9     | 3  | 0  | 2  | 2  | 7  | 3  | 2  | 1  | 6  | 27 | 0.491 | 26/53   |
| Total | 0  | 0  | 0  | 0  | 0  | 0  | 0  | 0  | 0  | 0  | 0.371 | 185/498 |

### 3.1.2 Random Forest

|       | 0  | 1  | 2  | 3  | 4  | 5  | 6  | 7  | 8  | 9  | Error | Rate   |
|-------|----|----|----|----|----|----|----|----|----|----|-------|--------|
| 0     | 46 | 0  | 0  | 1  | 0  | 0  | 0  | 0  | 0  | 0  | 0.021 | 1/47   |
| 1     | 0  | 42 | 0  | 0  | 0  | 1  | 0  | 1  | 0  | 0  | 0.045 | 2/44   |
| 2     | 0  | 2  | 48 | 0  | 0  | 0  | 0  | 1  | 0  | 0  | 0.059 | 3/51   |
| 3     | 0  | 0  | 0  | 41 | 0  | 1  | 0  | 0  | 3  | 0  | 0.089 | 4/45   |
| 4     | 0  | 1  | 0  | 0  | 52 | 1  | 1  | 4  | 0  | 1  | 0.133 | 8/60   |
| 5     | 0  | 0  | 0  | 4  | 0  | 42 | 1  | 0  | 1  | 0  | 0.125 | 6/48   |
| 6     | 0  | 0  | 0  | 0  | 1  | 0  | 47 | 0  | 2  | 0  | 0.060 | 3/50   |
| 7     | 0  | 1  | 0  | 0  | 0  | 0  | 0  | 48 | 0  | 2  | 0.059 | 3/51   |
| 8     | 1  | 2  | 0  | 1  | 0  | 0  | 0  | 0  | 42 | 1  | 0.106 | 5/47   |
| 9     | 0  | 0  | 0  | 1  | 2  | 1  | 0  | 0  | 1  | 48 | 0.094 | 5/53   |
| Total | 0  | 0  | 0  | 0  | 0  | 0  | 0  | 0  | 0  | 0  | 0.081 | 40/496 |

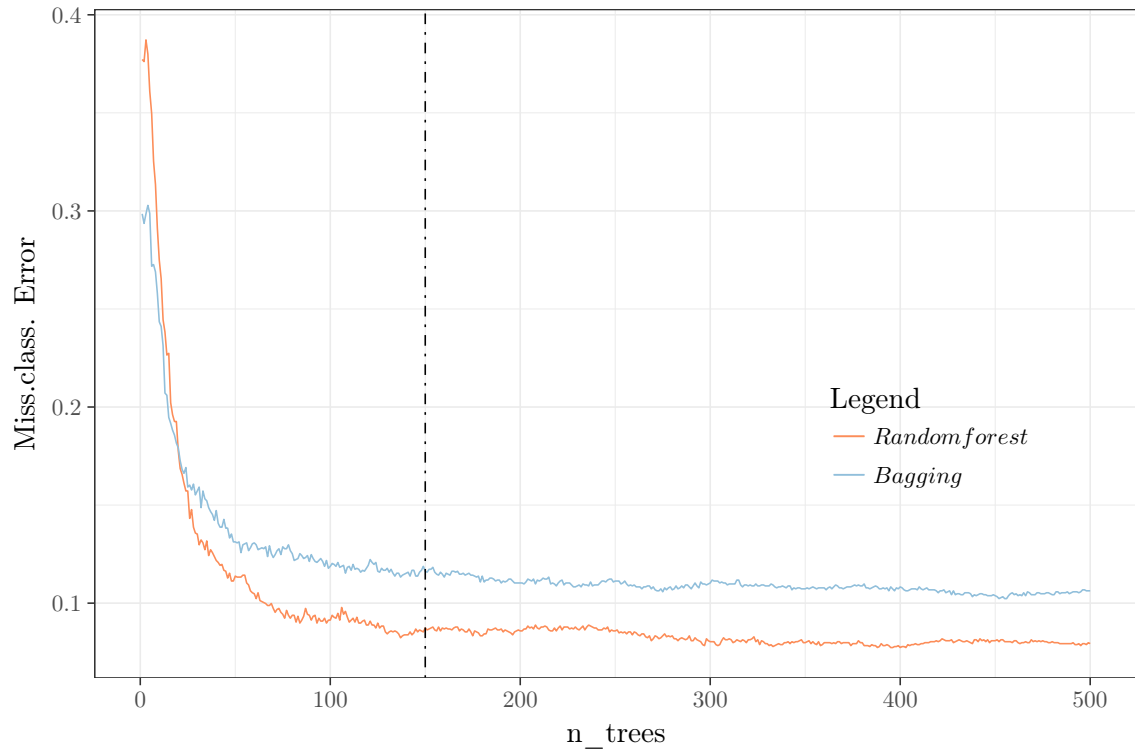|       | 0  | 1  | 2  | 3  | 4  | 5  | 6  | 7  | 8  | 9  | Error | Rate   |
|-------|----|----|----|----|----|----|----|----|----|----|-------|--------|
| 0     | 47 | 0  | 0  | 2  | 0  | 3  | 0  | 0  | 0  | 0  | 0.096 | 5/52   |
| 1     | 0  | 42 | 0  | 0  | 0  | 0  | 0  | 0  | 0  | 0  | 0.000 | 0/42   |
| 2     | 0  | 2  | 46 | 0  | 0  | 0  | 0  | 1  | 0  | 0  | 0.061 | 3/49   |
| 3     | 0  | 1  | 0  | 41 | 0  | 2  | 0  | 1  | 2  | 0  | 0.128 | 6/47   |
| 4     | 0  | 1  | 0  | 0  | 49 | 1  | 1  | 2  | 0  | 2  | 0.125 | 7/56   |
| 5     | 0  | 0  | 0  | 3  | 0  | 38 | 1  | 0  | 1  | 0  | 0.116 | 5/43   |
| 6     | 0  | 0  | 1  | 0  | 1  | 1  | 47 | 0  | 2  | 1  | 0.113 | 6/53   |
| 7     | 0  | 0  | 0  | 0  | 0  | 0  | 0  | 47 | 0  | 3  | 0.060 | 3/50   |
| 8     | 0  | 2  | 1  | 1  | 1  | 0  | 0  | 2  | 43 | 1  | 0.157 | 8/51   |
| 9     | 0  | 0  | 0  | 1  | 4  | 1  | 0  | 1  | 1  | 45 | 0.151 | 8/53   |
| Total | 0  | 0  | 0  | 0  | 0  | 0  | 0  | 0  | 0  | 0  | 0.103 | 51/496 |

Figure 3: Error measures $E_{out}$ and $E_{val}$ are plotted for different sizes of training and validation sets.

## 3.2 Neural Networks

## 3.3 Support vector machines

# 4 Discussion

# 5 Acknowledgments

# 6 Appendices

## 6.1 R-Code: Regression Tree

```
1  ## Libraries and seed
2  rm(list = ls())
3  library(caret)
4  library(readr)
5  library(tree)
6  library(randomForest)
7  library(gbm)
8  library(tikzDevice)      # library to export plots to .tex files
9
10 options(tikzMetricPackages = c("\\usepackage[utf8]{inputenc}", "\\usepackage
     [T1]{fontenc}",
11                                "\\usetikzlibrary{calc}", "\\usepackage{
                                    amssymb}"))
12
13 set.seed(420)
```

```r
14
15 if(!exists("create_confusion_matrix", mode = "function")){
16     source("Help_Scripts/to_latex_functions.R")
17 }
18
19 #-------------------#
20
21 ## Data
22 path_data <- paste0(getwd(), "/data")
23 path_to_here <- paste0(getwd(), "/Tree_Based_Methods")   # getwd give path
       to project
24 # which is one folder over
25
26 train_data <- read.csv(paste0(path_data, "/Train_Digits_20171108.csv"),
       header = TRUE)
27 unclassified_data <- read.csv(paste0(path_data, "/Test_Digits_20171108.csv")
       , header = TRUE)
28
29 # Remove unnessesary varibles which have a low variance
30 split_train_test <- createDataPartition(train_data$Digit, p = 0.8, list =
       FALSE)
31 test_data <- train_data[-split_train_test, ]
32 train_data <- train_data[split_train_test, ]
33
34 # Remove variable with low variance which are near zero. Doing it after
35 # splitting in train/test set to avoid contaminating the data.
36 near_zero_variables <- nearZeroVar(train_data[,-1], saveMetrics = T, freqCut
        = 10000/1,uniqueCut = 1/7)
37 cut_variables <- rownames(near_zero_variables[near_zero_variables$nzv ==
       TRUE,])
38 variables <- setdiff(names(train_data), cut_variables)
39 train_data <- train_data[, variables]
40 test_data <- test_data[, variables]
41
42 train_data[, 1] <- as.factor(train_data[, 1])
43 test_data[, 1] <- as.factor(test_data[, 1])
44
45 #train_data[,1] <- as.factor(train_data[,1])
46
47 unclassified_data[,1] <- as.factor(unclassified_data[,1])
48
49 sum(near_zero_variables$nzv)
50
51 #-------------------#
52
53 ## REGRESSION - tree
54
55 regression <- function(
56     minimum_development,
57     train_data,
58     test_data
59     ){
60     # Chanage name of pixel columns to work with tikz library
61
62     print(getClass(class(train_data)))
63     colnames(train_data)[ 2:length(train_data[1,])] <- c(paste0("Pixel", 1:(
           length(train_data[1,]) - 1)))
64     colnames(test_data)[ 2:length(train_data[1, ])] <- c(paste0("Pixel", 1:(
           length(train_data[1,]) - 1)))
```

```r
65      minimum_development <- 0.005
66      tree_model <- tree(Digit ~ ., data = train_data, mindev = minimum_
            development)
67      plot(tree_model)
68      text(tree_model, cex = .5)
69      print(summary(tree_model))
70
71      cross_validation <- cv.tree(tree_model, K = 10)
72      cross_validation$k[1] <- 0
73      alpha <- round <- round(cross_validation$k)
74
75      plot(cross_validation$size, cross_validation$dev, type = "b",
76          xlab = "Number of terminal nodes", ylab = "CV error")
77
78      ggplot_df <- data.frame(size = cross_validation$size, dev = cross_
            validation$dev)
79
80      destination_path <- paste0(path_to_here, "/Results_TBM/Regression_Tree")
81
82      ggplot1 <- ggplot(data = ggplot_df, aes(x = size, y = dev)) +
83                  geom_line(aes(colour = "$RegressionTree$"), linetype = "
                        dashed") +
84                  geom_point() +
85                  geom_vline(xintercept = 20, color = "black", linetype = "
                        dotdash") +
86                  xlab("n\\_{terminal nodes}") +
87                  ylab("CV Error") +
88                  scale_colour_manual("Legend",
89                                      breaks = c("$RegressionTree$"),
90                                      values = c("#91bfdb"),
91                                      guide = guide_legend(override.aes = list(
92                                          linetype = c("solid"),
93                                          shape = c(16)
94                                      ))) +
95                  theme_bw() +
96                  theme(legend.position = c(0.8, 0.355),
97                      legend.background = element_rect(fill=alpha('white', 0)))
98      ggsave(paste0(destination_path, ".png"))
99
100     ggplot_to_latex(ggplot1, destination_path, width = 5, height = 5)
101     tree_prune <- prune.tree(tree_model, best = 20)
102     summary(tree_prune)
103
104     tikz(file = paste0(destination_path, "_Tree.tex"), width = 6, height =
            4)
105     plot(tree_prune)
106     text(tree_prune, cex = .5)
107     dev.off()
108
109     predicted <- predict(tree_prune, test_data, type = "class")
110     create_confusion_matrix(predicted, test_data[,1], destination_path)
111     }
112
113 regression(0.05, train_data, test_data)
114 #------------------#
115
116 ## RANDOM FORREST -randomForest
117
118
```

```
119
120 #------------------#
121
122 ## BOOSTING - gbm
123
124 #------------------#
```

../R_scripts/Tree_Based_Methods/Regression_Tree.R

## 6.2 R-Code: Random Forest

```
1  ## Libraries and seed
2  rm(list = ls())
3  library(randomForest)    # library giving a easy-to-use random forest method
4  library(caret)           # useful library to split up data set
5  library(tikzDevice)      # library to export plots to .tex files
6  library(xtable)          # library to export data frames to tables in .tex
       files
7  set.seed(420)            # seed to replicate results and get consistent test
       and training set
8
9  # Load help script with functions to export the results to latex
10 # These functions gathered to avoid duplicate code
11 if(!exists("create_confusion_matrix", mode = "function")){
12     source("Help_Scripts/to_latex_functions.R")
13 }
14
15 #-------------------#
16
17 ## Data
18 path_data <- paste0(getwd(), "/data")
19 path_to_here <- paste0(getwd(), "/Tree_Based_Methods")   # getwd give path
       to project
20                                                    # which is one
                                                         folder over
21
22 train_data <- read.csv(paste0(path_data, "/Train_Digits_20171108.csv"),
       header = TRUE)
23 unclassified_data <- read.csv(paste0(path_data, "/Test_Digits_20171108.csv")
       , header = TRUE)
24
25 train_data[,1] <- as.factor(train_data[, 1])
26
27 # split training set into training and test set
28
29 split_train_test <- createDataPartition(train_data$Digit, p = 0.8, list =
       FALSE)
30 test_data <- train_data[-split_train_test, ]
31 train_data <- train_data[split_train_test, ]
32
33 #-------------------#
34
35 ## Random forest
36 # Train forest
37 train_random_forest <- function(
38     data,
39     n_trees,
40     minimum_development = 0.01
41     ){
42         random_forest <- randomForest(Digit ~ .,
43                                       data = data,
44                                       ntree = n_trees,
45                                       #mindev = minimum_development,
46                                       importance = TRUE,
47                                       na.action = na.exclude)
48         return(random_forest)
49     }
50
```

```r
51  # Plot error as the number of trees increase
52
53  plot_error_development <- function(
54      random_forest_data,
55      destination_path
56      ){
57          error_data <- data.frame(n_trees = 1:nrow(random_forest_data$err.
              rate),
58                                   error <- random_forest_data$err.rate[,"OOB"
                                       ])
59
60          write.csv(error_data, file = paste0(destination_path, ".csv"))
61
62          ggplot1 <- ggplot(data = error_data, aes(x = n_trees)) +
63              geom_line(aes(y = error, colour = "$Random forest")) +
64              xlab("$n\\_{trees}$") +
65              ylab("Miss. class. Error") +
66              scale_colour_manual("Legend",
67                                  breaks = c("$Random forest$"),
68                                  values = c("black"),
69                                  guide = guide_legend(override.aes = list(
70                                      linetype = c("solid"),
71                                      shape = c( 16)
72                                  ))) +
73              theme(legend.position = c(0.9, 0.2))
74          ggsave(paste0(destination_path, ".png"))
75
76          ggplot_to_latex(ggplot1, destination_path, width = 6, height = 4)
77  }
78
79  main <- function(){
80      n_trees = 50
81      random_forest <- train_random_forest(train_data, n_trees)
82      plot_error_development(random_forest, paste0(path_to_here, "/Results_TBM
          /Random_Forest_",
83                                                   n_trees, "trees_Error_plot"
                                                       ))
84
85      prediction <- predict(random_forest, newdata = test_data)
86      create_confusion_matrix(predicted_value = prediction, true_value = test_
          data$Digit,
87                              paste0(path_to_here, "/Results_TBM/Random_Forest
                                  _",
88                                                           n_trees, "
                                                               trees"))
89  }
90
91  main()
```

../R_scripts/Tree_Based_Methods/Random_Forest.R

## 6.3 R-Code: Bagging

```
1  ## Libraries and seed
2  rm(list = ls())
3  library(randomForest)     # library giving a easy-to-use random forest method
4  library(caret)            # useful library to split up data set
5  library(tikzDevice)       # library to export plots to .tex files
6  library(xtable)           # library to export data frames to tables in .tex
       files
7  set.seed(420)             # seed to replicate results and get consistent test
       and training set
8
9  # Load help script with functions to export the results to latex
10 # These functions gathered to avoid duplicate code
11 if(!exists("create_confusion_matrix", mode = "function")){
12     source("Help_Scripts/to_latex_functions.R")
13 }
14
15 #-------------------#
16
17 ## Data
18
19 path_data <- paste0(getwd(), "/data")
20 path_to_here <- paste0(getwd(), "/Tree_Based_Methods")   # getwd give path
       to project
21 # which is one folder over
22
23 train_data <- read.csv(paste0(path_data, "/Train_Digits_20171108.csv"),
       header = TRUE)
24 unclassified_data <- read.csv(paste0(path_data, "/Test_Digits_20171108.csv")
       , header = TRUE)
25
26 train_data[,1] <- as.factor(train_data[, 1])
27
28 # split training set into training and test set
29
30 split_train_test <- createDataPartition(train_data$Digit, p = 0.8, list =
       FALSE)
31 test_data <- train_data[-split_train_test, ]
32 train_data <- train_data[split_train_test, ]
33
34 #-------------------#
35
36 ## Random forest
37 # Train forest
38 train_bagging <- function(
39     data,
40     n_trees,
41     minimum_development = 0.01
42 ){
43     n_features <- ncol(data) - 1
44     bagging <- randomForest(Digit ~ .,
45                             data = data,
46                             ntree = n_trees,
47                             #mindev = minimum_development,
48                             mtry = n_features,
49                             importance = TRUE,
50                             na.action = na.exclude)
51     return(bagging)
```

```r
52 }
53
54 # Plot error as the number of trees increase
55
56 plot_error_development <- function(
57     random_forest_data,
58     destination_path
59 ){
60     error_data <- data.frame(n_trees = 1:nrow(random_forest_data$err.rate),
61                              error <- random_forest_data$err.rate[,"OOB"])
62
63     write.csv(error_data, file = paste0(destination_path ,".csv"))
64     ggplot1 <- ggplot(data = error_data, aes(x = n_trees)) +
65         geom_line(aes(y = error, colour = "$Bagging$")) +
66         xlab("n\\_{trees}") +
67         ylab("Miss.class. Error") +
68         scale_colour_manual("Legend",
69                             breaks = c("$Bagging$"),
70                             values = c("black"),
71                             guide = guide_legend(override.aes = list(
72                                 linetype = c("solid"),
73                                 shape = c( 16)
74                             ))) +
75         theme(legend.position = c(0.9, 0.2))
76     ggsave(paste0(destination_path, ".png"))
77
78     ggplot_to_latex(ggplot1, destination_path, width = 6, height = 4)
79 }
80
81 main <- function(){
82     n_trees <- 50
83     bagging <- train_bagging(train_data, n_trees)
84     plot_error_development(bagging, paste0(path_to_here, "/Results_TBM/
85         Bagging_",
86                                                         n_trees, "trees_Error_plot"))
86
87     prediction <- predict(bagging, newdata = test_data)
88
89     create_confusion_matrix(predicted_value = prediction, true_value = test_
90         data$Digit,
90                             paste0(path_to_here, "/Results_TBM/Bagging_",
91                                     n_trees, "trees"))
92 }
93
94 main()
```

../R_scripts/Tree_Based_Methods/Bagging.R

11

## 6.4 R-Code: Boosting

```
1  ## Libraries and seed
2  rm(list = ls())
3  library(caret)            # useful library to split up data set
4  library(tikzDevice)       # library to export plots to .tex files
5  library(gbm)              # library with powerful boosting method
6  library(xtable)           # library to export data frames to tables in .tex
       files
7  set.seed(420)             # seed to replicate results and get consistent test
       and training set
8
9  # Load help script with functions to export the results to latex
10 # These functions gathered to avoid duplicate code
11 if(!exists("create_confusion_matrix", mode = "function")){
12     source("Help_Scripts/to_latex_functions.R")
13 }
14
15 #-------------------#
16
17 ## Data
18
19 path_data <- paste0(getwd(), "/data")
20 path_to_here <- paste0(getwd(), "/Tree_Based_Methods")   # getwd give path
       to project
21
22 train_data <- read.csv(paste0(path_data, "/Train_Digits_20171108.csv"),
       header = TRUE)
23 unclassified_data <- read.csv(paste0(path_data, "/Test_Digits_20171108.csv")
       , header = TRUE)
24
25 train_data[,1] <- as.factor(train_data[, 1])
26
27 # split training set into training and test set
28
29 split_train_test <- createDataPartition(train_data$Digit, p = 0.8, list =
       FALSE)
30 test_data <- train_data[-split_train_test, ]
31 train_data <- train_data[split_train_test, ]
32
33 #-------------------#
34
35 ## Boosting
36 # Train booster
37 boosting <- function(
38     data,
39     n_trees,
40     minimum_development = 0.01,
41     interaction_depth = 2,
42     shrinkage = 0.001
43 ){
44     # boosting <- gbm(Digit ~ .,
45     #                 data = data,
46     #                 distribution = "multinomial",
47     #                 n.trees = n_trees,
48     #                 interaction.depth = interaction_depth,
49     #                 shrinkage = shrinkage,
50     #
51     #                 bag.fraction = 1,
```

```r
52      #                     cv.folds = 10,
53      #                     n.cores = 4)
54
55
56      tune_control <- trainControl(method = "cv",
57                                   number = 5,
58                                   repeats = 1)
59      training_grid <- expand.grid(n.trees = c(n_trees),
60                                   interaction.depth = c(interaction_depth),
61                                   shrinkage = c(shrinkage),
62                                   n.minobsinnode = c(10))
63      print(training_grid)
64      boosting <- train(Digit ~ ., data = data, method = "gbm",
65                        trControl = tune_control,
66                        tuneGrid = training_grid)
67      return(boosting)
68 }
69
70
71 # Plot error as the number of trees increase
72
73 plot_error_development <- function(
74     boosting_data,
75     destination_path
76 ){
77     error_data <- data.frame(n_trees = 1:length(boosting_data$cv.error),
78                              error <- boosting_data$cv.error)
79     write.csv(error_data, file = paste0(destination_path ,".csv"))
80
81     ggplot1 <- ggplot(data = error_data, aes(x = n_trees)) +
82         geom_line(aes(y = error, colour = "$Boosting$")) +
83         xlab("$n_{trees}$") +
84         ylab("Miss.class. Error") +
85         scale_colour_manual("Legend",
86                             breaks = c("$Boosting$"),
87                             values = c("black"),
88                             guide = guide_legend(override.aes = list(
89                                 linetype = c("solid"),
90                                 shape = c( 16)
91                             ))) +
92         theme(legend.position = c(0.9, 0.2)) +
93         theme_bw() +
94         theme(legend.position = c(0.8, 0.355),
95               legend.background = element_rect(fill=alpha('white', 0))) +
96     ggsave(paste0(destination_path, ".png"))
97
98     ggplot_to_latex(ggplot1, destination_path, width = 6, height = 4)
99 }
100
101 main <- function(){
102     n_trees = 10
103     boosting_train <- boosting(train_data,n_trees)
104     plot_error_development(boosting_train, paste0(path_to_here,
105                                                   "/Results_TBM/Boosting_",
106                                                   n_trees,
107                                                   "trees_Error_plot"))
108     #predicted <- predict(boosting_train, test_data)
109     #create_confusion_matrix(predicted, test_data$Digit, paste0(path_to_here
           ,
```

```
110        #                                                    "/Results_
             TBM/Boosting_",
111        #                                                    n_trees))
112 }
113
114 main()
```

## 6.5   R-Code: Neural Network

```
1  ## Libraries and seed
2  library(h2o)
3  library(caret)
4  library(reshape2)
5
6  set.seed(420)
7
8  #-------------------#
9
10 ## Data
11
12 path_to_here <- getwd()
13
14 train_data <- read.csv(paste0(path_to_here, "/data/Train_Digits_20171108.csv
       "))
15 unclassified_data <- read.csv(paste0(path_to_here, "/data/Test_Digits_
       20171108.csv"))
16
17 local.h2o <- h2o.init(ip = "localhost", port = 54321, startH2O = TRUE,
       nthreads = -1)
18
19 train_data[,1] <- as.factor(train_data[, 1])
20 split_train_test <- createDataPartition(train_data$Digit, p = 0.8, list =
       FALSE)
21 test_data <- train_data[-split_train_test, ]
22 train_data <- train_data[split_train_test, ]
23
24 train_data <- as.h2o(train_data)
25 unclassified_data <- as.h2o(unclassified_data)
26 test_data <- as.h2o(test_data)
27
28 #-------------------#
29
30 ## Getting useful data from grid run of neural networkss
31
32 get_data_in_df <- function(
33     data
34 )
35 {
36     n <- length(data@model_ids)
37     mse_errors <- rep(0,n)
38     mean_per_class_errors <- rep(0,n)
39     hidden <- rep("", n)
40     str(hidden)
41     rate <- rep(0,n)
42     l1 <- rep(0,n)
43     epochs <- rep(0,n)
44     model_numbers <- rep(0,n)
45     train_error <- rep(0,n)
46     train_mse <- rep(0,n)
47     test_error <- rep(0,n)
48     test_mse <- rep(0,n)
49     activation <- rep("",n)
50     input_dropout_ratio <- rep(0,n)
51     nesterov_accelerated_gradient <- rep("", n)
52
53     model_df <- data.frame(model_numbers = mse_errors, hidden, rate, l1,
```

```r
            epochs,
                            train_error, test_error, train_mse, test_mse,
                                activation, input_dropout_ratio,
                                stringsAsFactors = FALSE)
    str(model_df)

    for(i in 1:n){
        model <- h2o.getModel(data@model_ids[[i]])
        model_df$mse_errors[i] <- h2o.mse(model)
        #model_df$mean_per_class_error[i] <- model@model$cross_validation_
            metrics@metrics$mean_per_class_error
        model_df$mean_per_class_error[i] <- h2o.performance(model, xval = T)
            @metrics$mean_per_class_error

        model_paramaters <- model@allparameters
        model_name <- model@model_id
        model_number <- sub(".*model_(.*)$", "\\1", model_name)
        model_df$model_numbers[i] <- as.integer(model_number)
        model_df$hidden[i] <- paste(as.character(model_paramaters$hidden),
            sep = " ", collapse = ", ")
        model_df$rate[i] <- model_paramaters$rate
        model_df$l1[i] <- model_paramaters$l1
        model_df$epochs[i] <- model_paramaters$epochs
        model_df$activation[i] <- model_paramaters$activation
        model_df$input_dropout_ratio[i] <- model_paramaters$input_dropout_
            ratio
        model_df$nesterov_accelerated_gradient[i] <- model_paramaters$
            nesterov_accelerated_gradient

        #print(model)
        train_performance <- h2o.performance(model, train_data)@metrics
        train_performance_error <- train_performance$mean_per_class_error
        train_performance_mse <- train_performance$MSE

        model_df$train_error[i] <- train_performance_error
        model_df$train_mse[i] <- train_performance_mse

        test_performance <- h2o.performance(model, test_data)@metrics
        test_predictions <- h2o.predict(model, test_data)
        test_accuracy <- test_predictions$predict == test_data$Digit
        test_performance_error <- 1 - mean(test_accuracy)
        #test_performance_error <- test_performance$mean_per_class_error
        test_performance_mse <- test_performance$MSE

        model_df$test_error[i] <- test_performance_error
        model_df$test_mse[i] <- test_performance_mse

    }
    model_df <- model_df[with(model_df, order(model_numbers)),]
    model_df
}

activation <- list("Rectifier", "RectifierWithDropOut")# "Tanh")
hidden <- list(c(100,100), c(150, 150), c(100, 100, 100)) #c(100, 100, 100))
    #, c(150, 150, 150))
input_dropout_ratio <- list(0, 0.2)
nesterov_accelerated_gradient <- list( TRUE, FALSE)
epochs <- list(20)#, 20)
l1 = list(1.4e-5)
```

```
104 hyper_params <- list(activation = activation, hidden = hidden, input_dropout
        _ratio = input_dropout_ratio, nesterov_accelerated_gradient = nesterov_
        accelerated_gradient, epochs = epochs, l1 = l1)
105
106 grid_deep_learning <- h2o.grid(algorithm = "deeplearning",
107                                 x = 2:785,
108                                 y = 1,
109                                 training_frame = train_data,
110                                 nfolds = 10,
111                                 stopping_metric = "MSE",
112                                 stopping_tolerance = 0.0025,
113                                 hyper_params = hyper_params)
114     save_results <- function(results){
115     write.csv(results, file = paste0(path_to_here, "/Neural_Networks/results
            _NN/grid_run_evenodd2.csv"))
116 }
117
118 df <- get_data_in_df(grid_deep_learning)
119 save_results(df)
120
121 results_df <- df
122
123 results_df <- results_df[with(results_df, order(mean_per_class_error)),]
124 results_df$row_names <- 1:length(results_df[,1])
125
126 melt_datas <- melt(results_df[c("test_error","mean_per_class_error", "row_
        names",
127                                 "model_numbers")], id = c("row_names", "
                                    model_numbers"))
128
129 # Plot classification error
130 plot_list[[1]] <- ggplot(data=melt_datas,
131                         aes(x=row_names, y=value)) +
132     geom_point(aes(colour = as.factor(model_numbers), group = as.factor(
            model_numbers)), size = 3) +
133     geom_line(aes(group = variable)) +
134     labs(y = "Missclassification error in range 0 to 1",
135         x = "Models",
136         title = "Missclassification error for training model and test set",
137         caption = "Top - Training model, Bottom - Test set",
138         colour = "Model id") +
139     scale_y_continuous(limits = c(0, 0.2))
140 ggsave(paste0(path_to_here,"/Neural_Networks/results_NN/per_class_error3.png
        "))
141
142                                 deep_learning_predicting <- h2o.predict(object
                                        = deep_learning_results, newdata = test_
                                        data)
143 deep_learning_performance <- h2o.performance(model = deep_learning_results3,
        newdata = test_data)
144 deep_learning_performance
145 deep_learning_predicting_data_frame <- as.data.frame((deep_learning_
        predicting))
146
147
148 deep_learning_results2 <- h2o.deeplearning(x = 2:785,
149                                             y = 1,
150                                             training_frame = train_data,
151                                             activation = "Tanh",
```

```
152                                             hidden = c(160, 160, 160, 160,
                                                    160),
153                                             nfolds = 10,
154                                             keep_cross_validation_predictions
                                                    = TRUE,
155                                             epochs = 40)
156
157 deep_learning_results3 <- h2o.deeplearning(x = 2:785,
158                                             y = 1,
159                                             training_frame = train_data,
160                                             #activation = "
                                                    RectifierWithDropout",
161                                             activation = "Rectifier",
162                                             input_dropout_ratio = 0.2,
163                                             #hidden_dropout_ratios = c(0.2,
                                                    0.2, 0.2),
164                                             nfolds = 10,
165                                             balance_classes = TRUE,
166                                             hidden = c(150, 150, 150),
167                                             momentum_stable = 0.99,
168                                             nesterov_accelerated_gradient =
                                                    TRUE,
169                                             epochs = 15)
170
171 h2o.performance(deep_learning_results3, test_data)
```
../R_scripts/Neural_Networks/uneural_network.R

## 6.6   R-Code: Convolutional Neural Network

```
1  ## Libraries and seed
2  rm(list = ls())
3  library(mxnet)
4  library(caret)
5  set.seed(420)
6
7  #-------------------#
8
9  ## Data
10
11 path_to_here <- getwd()
12
13 train_data <- read.csv(paste0(path_to_here, "/data/Train_Digits_20171108.csv
       "), header = TRUE)
14 unclassified_data <- read.csv(paste0(path_to_here, "/data/Test_Digits_
       20171108.csv"), header = TRUE)
15
16 train_data[,1] <- as.factor(train_data[, 1])
17
18 # split training set into training and test set
19
20 split_train_test <- createDataPartition(train_data$Digit, p = 0.8, list =
       FALSE)
21 test_data <- train_data[-split_train_test, ]
22 train_data <- train_data[split_train_test, ]
23
24 # convert to matrix, required by "mxnet"
25
26 train <- data.matrix(train_data)
27 test <- data.matrix(test_data)
28
29 train_x <- t(train[, -1]/255)
30 train_y <- train[, 1]
31
32 train_array <- train_x
33 dim(train_array) <- c(28, 28, 1, ncol(train_x))
34
35 #train_x <- t(train_x)#/255)
36
37 test_x <- test[, -1]
38 test_y <- test[, 1]
39
40 test_x <- t(test_x/255)
41
42 # transpose and normalize to more
43
44 #-------------------#
45
46 ## Setting up Convolutional Neural Network(CNN)
47
48 data <- mx.symbol.Variable("data")
49 fc1 <- mx.symbol.FullyConnected(data, name="fc1", num_hidden=128)
50 act1 <- mx.symbol.Activation(fc1, name="relu1", act_type="relu")
51 fc2 <- mx.symbol.FullyConnected(act1, name="fc2", num_hidden=64)
52 act2 <- mx.symbol.Activation(fc2, name="relu2", act_type="relu")
53 fc3 <- mx.symbol.FullyConnected(act2, name="fc3", num_hidden=10)
54 softmax <- mx.symbol.SoftmaxOutput(fc3, name="sm")
```

```
55
56 devices <- mx.cpu()
57
58 mx.set.seed(0)
59
60 model <- mx.model.FeedForward.create(softmax, X=train_x, y=train_y,
61                                      ctx=devices, num.round=10, array.batch.
                                          size=100,
62                                      learning.rate=0.07, momentum=0.9,  eval
                                          .metric=mx.metric.accuracy,
63                                      initializer=mx.init.uniform(0.07),
64                                      epoch.end.callback=mx.callback.log.
                                          train.metric(100))
65
66 preds <- predict(model, test_x)
```
../R_scripts/Neural_Networks/convolutional_neural_network.R

## 6.7 R-Code: K-nearest Neighbours

```
1  ## Libraries and seed
2  rm(list = ls())
3
4  library(caret)          # useful library to split up data set
5  library(tikzDevice)     # library to export plots to .tex files
6  library(xtable)         # library to export data frames to tables in .tex
       files
7  library(kknn)
8  set.seed(420)           # seed to replicate results and get consistent test
       and training set
9
10 options(tikzMetricPackages = c("\\usepackage[utf8]{inputenc}", "\\usepackage
       [T1]{fontenc}",
11                                 "\\usetikzlibrary{calc}", "\\usepackage{
                                      amssymb}"))
12 #------------------#
13
14 ## Data
15 path_data <- paste0(getwd(), "/data")
16 path_to_here <- paste0(getwd(), "/Tree_Based_Methods")   # getwd give path
       to project
17 # which is one folder over
18
19
20 train_data <- read.csv(paste0(path_data, "/Train_Digits_20171108.csv"),
       header = TRUE)
21 unclassified_data <- read.csv(paste0(path_data, "/Test_Digits_20171108.csv")
       , header = TRUE)
22
23 train_data[,1] <- as.factor(train_data[, 1] )
24
25 # split training set into training and test set
26
27 split_train_test <- createDataPartition(train_data$Digit, p = 0.8, list =
       FALSE)
28 test_data <- train_data[-split_train_test, ]
29 train_data <- train_data[split_train_test, ]
30
31 knn_pred <- kknn(Digit ~ ., train = train_data[], test = test_data[-1], k
       =1)
32 k_pred <- fitted.values(knn_pred)
33 confusionMatrix(k_pred, test_data[,1])
```

<div align="center">../R_scripts/K_Nearest_Neighbors/kknn.R</div>

## 6.8 R-Code: Support Vector Machines

```
1  ## Libraries and seed
2  rm(list = ls())
3  library(e1071)
4  library(caret)          # useful library to split up data set
5  #library(tikzDevice)     # library to export plots to .tex files
6  #library(xtable)         # library to export data frames to tables in .tex
       files
7  library(readr)
8
9  set.seed(420)            # seed to replicate results and get consistent test
       and training set
10
11 #options(tikzMetricPackages = c("\\usepackage[utf8]{inputenc}", "\\
       usepackage[T1]{fontenc}",
12 #                               "\\usetikzlibrary{calc}", "\\usepackage{
       amssymb}"))
13
14 #-------------------#
15
16 ## Data
17 path_data <- paste0(getwd(), "/data")
18 path_to_here <- paste0(getwd(), "/Support_Vector_Machines")  # getwd give
       path to project
19 # which is one folder over
20
21 train_data <- read.csv(paste0(path_data, "/Train_Digits_20171108.csv"),
       header = TRUE)
22 unclassified_data <- read.csv(paste0(path_data, "/Test_Digits_20171108.csv")
       , header = TRUE)
23
24 train_data[,1] <- as.factor(train_data[, 1])
25
26 nzr <- nearZeroVar(train_data[,-1], saveMetrics = TRUE, freqCut = 10000/1,
       uniqueCut = 1/7)
27 sum(nzr$zeroVar)
28
29 sum(nzr$nzv)
30
31 cut_variables <- rownames(nzr[nzr$nzv == TRUE, ])
32 variable <- setdiff(names(train_data), cut_variables)
33 train_data <- train_data[, variable]
34
35 split_train_test <- createDataPartition(train_data$Digit, p = 0.8, list =
       FALSE)
36 test_data <- train_data[-split_train_test, ]
37 train_data <- train_data[split_train_test, ]
38
39 label <- train_data[1]
40 train_data$Digit <- NULL
41 train_data <- train_data/255
42 cov_train <- cov(train_data)
43
44 train_pc <- prcomp(cov_train)
45 varex <- train_pc$sdev^2/sum(train_pc$sdev^2)
46 varcum <- cumsum(varex)
47 result <- data.frame(num = 1:length(train_pc$sdev),
48                      ex = varex,
```

```r
49                              cum = varcum)
50
51 plot(result$num, result$cum, type = "b", xlim = c(0,100))
52 abline(v=25, lty=2)
53
54
55 train_score <- as.matrix(train_data) %*% train_pc$rotation[,1:25]
56 train_data <- cbind(label, as.data.frame(train_score))
57 colors <- rainbow(length(unique(train_data$Digit)))
58 names(colors) <- unique(train_data$label)
59
60 plot(train_data$PC1, train_data$PC2, type = "n", main = "First two Principal
       Components")
61 text(train_data$PC1, train_data$PC2, label = train_data$Digit, col = colors[
      train_data$Digit])
62
63 svm_model <- svm(Digit ~ ., data = train_data, cost = 8, kernel = "radial")
64
65 test_data2 <- test_data[-1]/255
66 test_data2 <- as.matrix(test_data2) %*% train_pc$rotation[,1:25]
67 test_data2 <- as.data.frame(test_data2)
68
69 predicted <- predict(svm_model, test_data2)
70
71 confusion_matrix <- confusionMatrix(predicted, test_data$Digit)
```
../R_scripts/Support_Vector_Machines/support_with_pca.R