

Machine Learning Analysis with Lagged Predictors

Falco J. Bargagli-Stoffi, Massimo Riccaboni, Armando Rungi

23/2/2020

Introduction

This R Markdown file reproduces the lagged machine learning analysis for the paper "*Machine learning for zombie hunting. Firms' failures, financial constraints, and misallocation*" by Falco J. Bargagli-Stoffi (IMT School for Advanced Studies/KU Leuven), Massimo Riccaboni (IMT School for Advanced Studies) and Armando Rungi (IMT School for Advanced Studies).

R Markdown

This is an R Markdown document. Markdown is a simple formatting syntax for authoring HTML, PDF, and MS Word documents. For more details on using R Markdown see <http://rmarkdown.rstudio.com>.

When you click the **Knit** button a document will be generated that includes both content as well as the output of any embedded R code chunks within the document. You can embed an R code chunks like the following.

Packages Upload

The following packages and functions are the ones used for the analyses performed in the R code. The `functions.R` file contains the functions `F1_score`, `balanced_accuracy`, `model_compare` and `DtD` that were developed to reproduce the following analyses.

```
rm(list=ls()) # to clean the memeory
memory.limit(size=1000000)
```

```
## [1] 1e+06
```

```
options(java.parameters = "-Xmx15000m")
library(rJava)
library(bartMachine)
library(haven)
library(plyr)
library(dplyr)
library(PRRoc)
library(rpart)
library(party)
library(caret)
library(devtools)
library(SuperLearner)
library(Metrics)
library(pROC)
library(Hmisc)
source('functions.R')
```

BART-mia

Run a Bayesian Additive Regression Tree analysis by using the overall data sample (no need to omit the observations with missing values).

```
set.seed(2020)
sample <- sample(seq_len(nrow(data_italy)),
                 size = nrow(omitted),
                 replace=FALSE)
data_italy_bart <- data_italy[sample,]
```

Select the same number of observations as in the previous models for the training and testing samples.

```
set.seed(2020)
train_sample <- sample(seq_len(nrow(data_italy_bart)),
                      size = nrow(data_italy_bart )*0.9, replace=FALSE)
train_bart <- data_italy_bart[train_sample,]
test_bart <- data_italy_bart[-train_sample,]
train_bart$X <- as.data.frame(train_bart[predictors])
test_bart$X <- as.data.frame(test_bart[predictors])
```

Run the analysis.

```
system.time({
bart_machine<-bartMachine(train_bart$X,
                          as.factor(train_bart$failure),
                          use_missing_data=TRUE)
})
```

Depict the performance measures by running the following chunks.

```
fitted.results.bart <- 1- round(predict(bart_machine,
                                       test_bart$X,
                                       type='prob'), 6)
```

```
#Roc
fg.bart<-fitted.results.bart[test_bart$failure==1]
bg.bart<-fitted.results.bart[test_bart$failure==0]
```

```
roc_bart<-roc.curve(scores.class0 = fg.bart,
                    scores.class1 = bg.bart,
                    curve = T)
plot(roc_bart)
```

```
pr_bart<-pr.curve(scores.class0 = fg.bart,
                  scores.class1 = bg.bart,
                  curve = T)
plot(pr_bart)
```

```
#Get Accuracy
fitted.bart <- ifelse(fitted.results.bart> 0.5, 1, 0)
f1_bart <- f1_score(fitted.bart,
                   test_bart$failure,
                   positive.class="1")
```

```
balanced_accuracy_bart <- balanced_accuracy(as.matrix(table(fitted.bart,
                                                            test_bart$failure)))
accuracy_bart <- as.data.frame(rbind(postResample(as.double(fitted.bart),
```

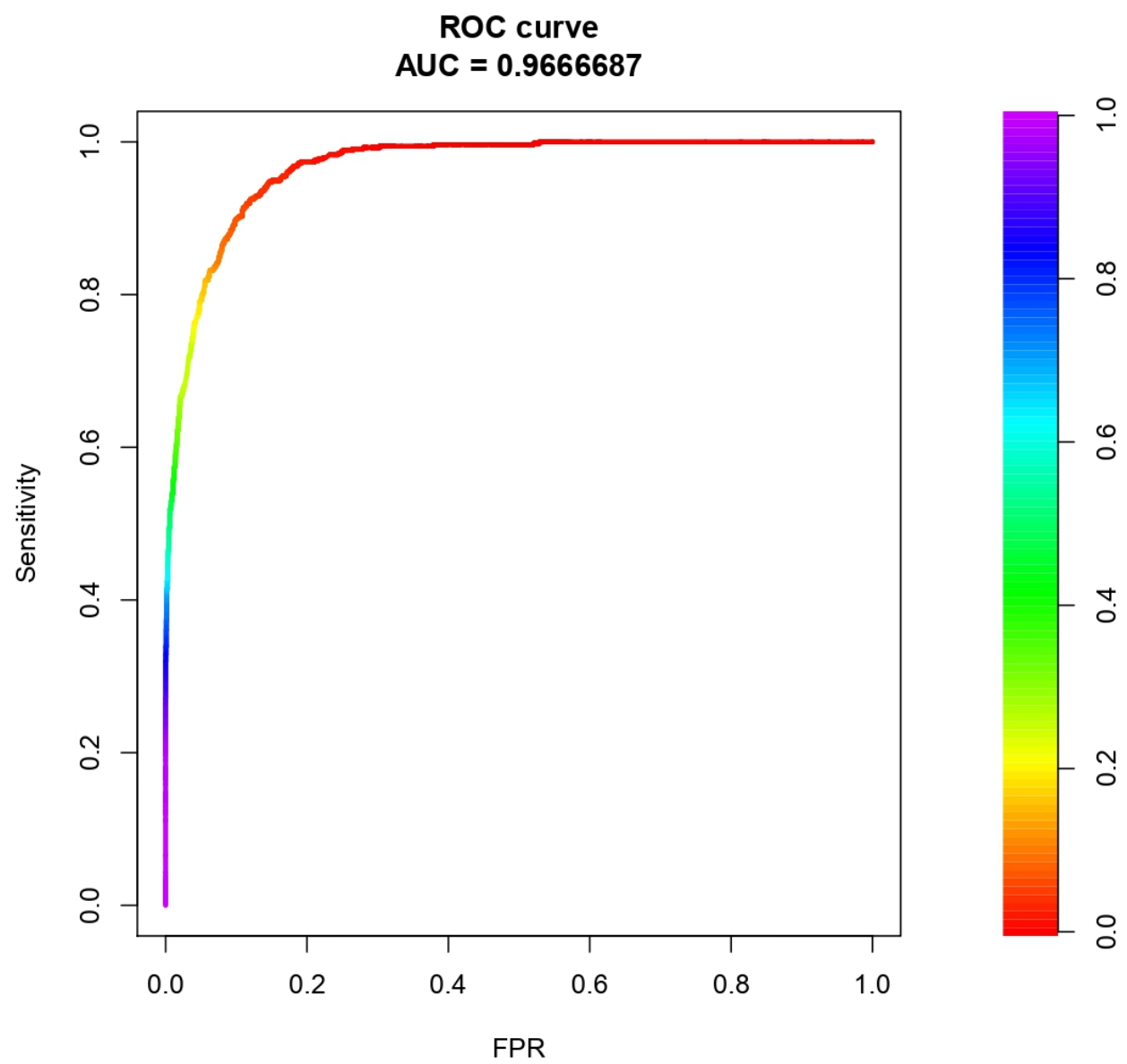


Figure 1: Area under the ROC curve, BART

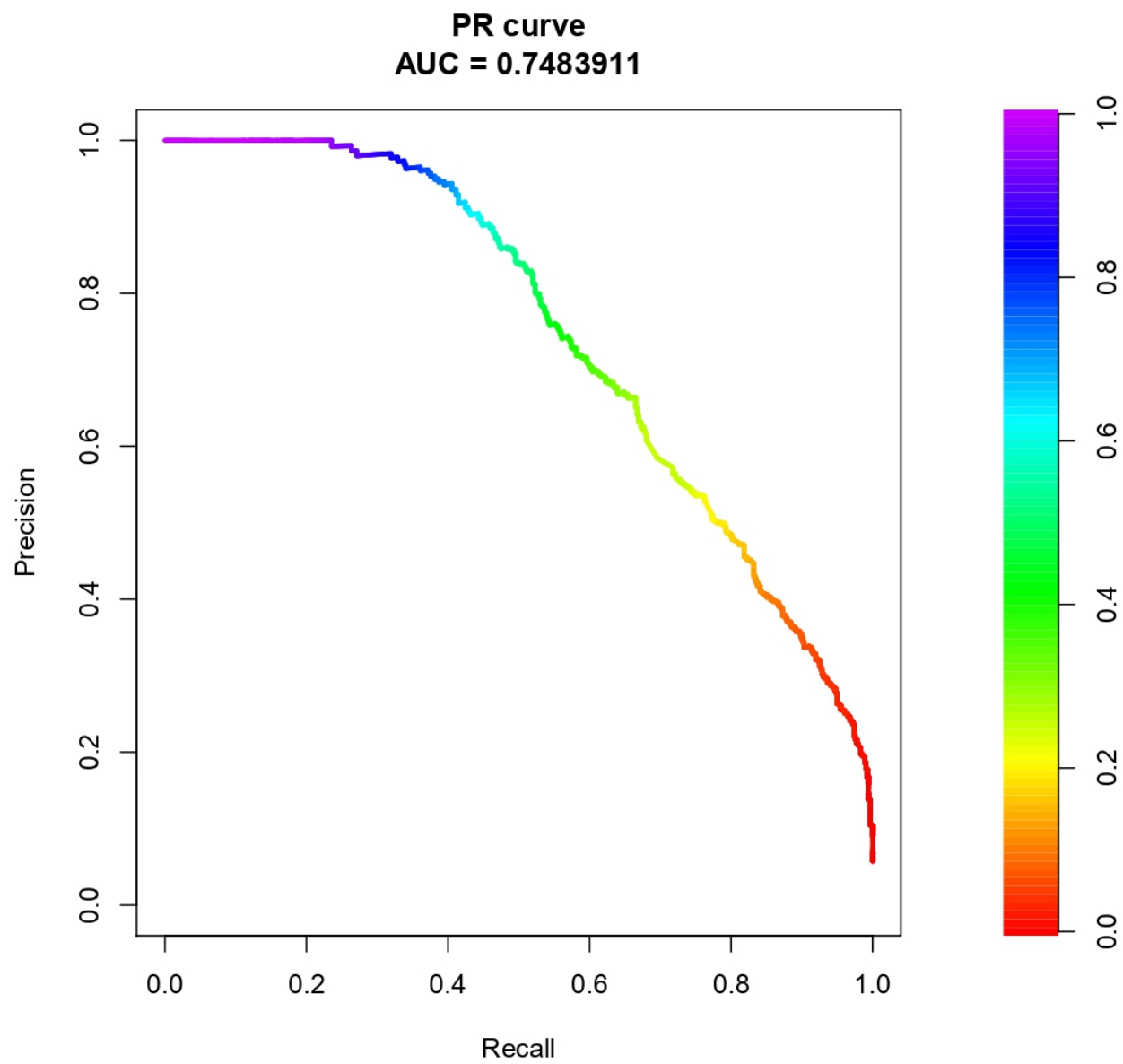


Figure 2: Area under the PR curve, BART

```

test_bart$failure)))

bart_fit <- as.data.frame(cbind(roc_bart$auc,
                              pr_bart$auc.integral,
                              f1_bart,
                              balanced_accuracy_bart,
                              accuracy_bart$Rsquared))
colnames(bart_fit) <- c("AUC", "PR", "f1-score", "BACC", "Rsquared")

```

Save Results

```

model_results <- rbind(logit_fit, ctree_fit, rf_fit, sl_fit, bart_fit)
write.csv(model_results, file = "model_results.csv")

```

Model Comparison

Here, we run an empirical horse race where we define two competitors (benchmark or “usual methods”) of our preferred BART methodology. Natural candidates are “default probability predictors” (credit-ratings type of measures) such as:

1. Altman Z-score;
2. Distance-to-Default.

As these measures do not provide direct predictions for failed firms, we create a series of dummy variables in the following way:

$$Z - dummy_i = \begin{cases} 1 & \text{if } Z - score_i \leq q, \\ 0 & \text{otherwise} \end{cases}$$

and

$$DtD - dummy_i = \begin{cases} 1 & \text{if } DtD_i \leq q, \\ 0 & \text{otherwise} \end{cases}$$

where q is the 1st to 10th percentile distribution of the Z-score and the DtD measures, respectively.

By doing so, we assume to be predicted as “failed”, those observations with values on the left tails of the Z and DtD measures.

We create these variables on the testing set and then we compare their performance, in terms of precision and false discovery rate (FDR), with the one of BART.

Z-score

```

precision_bart <- as.data.frame(t(model_compare(fitted.bart,
                                                test_bart$failure, "1")))
colnames(precision_bart) <- c("Precision BART", "FDR BART")
write.csv(precision_bart, "precision_and_fdr_bart.csv")

seq <- seq(1, 10, 1)
precision_zscore <- matrix(NA, ncol = 2, nrow = length(seq))
for(i in seq) {
  failed_zscore <- ifelse(test$Z_score <= quantile(test$Z_score, i/100),

```

```

precision_zscore[i, c(1:2)] <- t(model_compare(failed_zscore,
                                              1, 0)
                             test$failure, "1"))
}

```

Merton's Distance-to-Default (DtD)

The average equity volatility in the considered time series is 27.9120 (Bank of Italy), while for the average risk free interest rate we use the long term government bond yields" EMU (Eurostat) that has a value of 4.2937.

```
dtd_score <- DtD(mcap = test$shareholders_funds, debt = test$long_term_debt, vol = 27.9120, r = 4.2937)
```

```

seq <- seq(1, 10, 1)
precision_dtd_score <- matrix(NA, ncol = 2, nrow = length(seq))
for(i in seq) {
  failed_dtd_score <- ifelse(dtd_score <= quantile(dtd_score, i/100), 1, 0)
  precision_dtd_score[i, c(1:2)] <- t(model_compare(failed_dtd_score,
                                                    test$failure, "1"))
}

```

```

precision_models <- as.data.frame(cbind(precision_dtd_score, precision_zscore))
colnames(precision_models) <- c("Precision DtD",
                                "FDR DtD",
                                "Precision Z-score",
                                "FDR Z-score")
write.csv(precision_models, "precision_and_fdr_scores.csv")

```