

Tutorial 3

Statistical models

École Nationale des Ponts et Chaussées
Département Ingénierie Mathématique et Informatique – Master II

Loïc BRIN • Benoit ROGER

Exercise 1: Term structure of default probability deduced from a Transition Matrix.

This exercise is based on the following S&P transition matrix:

	AAA	AA	A	BBB	BB	B	CCC	D
AAA	90,8%	8,3%	0,7%	0,1%	0,1%	0,0%	0,0%	0,0%
AA	0,1%	91,2%	7,9%	0,6%	0,1%	0,1%	0,0%	0,0%
A	0,9%	2,4%	90,0%	5,4%	0,7%	0,3%	0,1%	0,1%
BBB	0,0%	0,3%	5,9%	86,9%	5,3%	1,2%	0,1%	0,2%
BB	0,0%	0,1%	0,7%	7,7%	80,5%	8,8%	1,0%	1,2%
B	0,0%	0,1%	0,2%	0,5%	6,5%	82,7%	4,1%	5,9%
CCC	0,2%	0,0%	0,2%	1,3%	2,3%	12,9%	60,6%	22,5%
D	0,0%	0,0%	0,0%	0,0%	0,0%	0,0%	0,0%	100,0%

available here http://defaultrisk.free.fr/data/TD2_1.csv.

1. Load the database and define a function that for a given number of year (n) and a given credit rating (l), returns a n -vector with the PD for each n years, deduced from the S&P transition matrix.

```

1 ### loading the transition matrix
2 data<-as.matrix(read.csv2("http://defaultrisk.free.fr/data/TD2_1.csv", header=F))
3
4 ### recursive function to compute the matricial power of a matrix
5 powA <- function(A, n)
6 {
7   if (n==1) return (A)
8   if (n==2) return (A%*%A)
9   if (n>2) return (A%*%powA(A,n-1))
10 }
11
12 ### function that extracts the PD for n years of rating l
13 PDt<-function(n,l)
14 {
15   res<-c()
16   for (i in 1:n)
17   {
18     res<-c(res,powA(data,i)[1,8])
19   }
20   return(res)
21 }

```

PDt is a function that for a given number of years n , and rating l , returns the PD for each n years and stores them in a vector.

2. Plot the PD for all the ratings, for the next 15 years, deduced from the transition matrix.

```

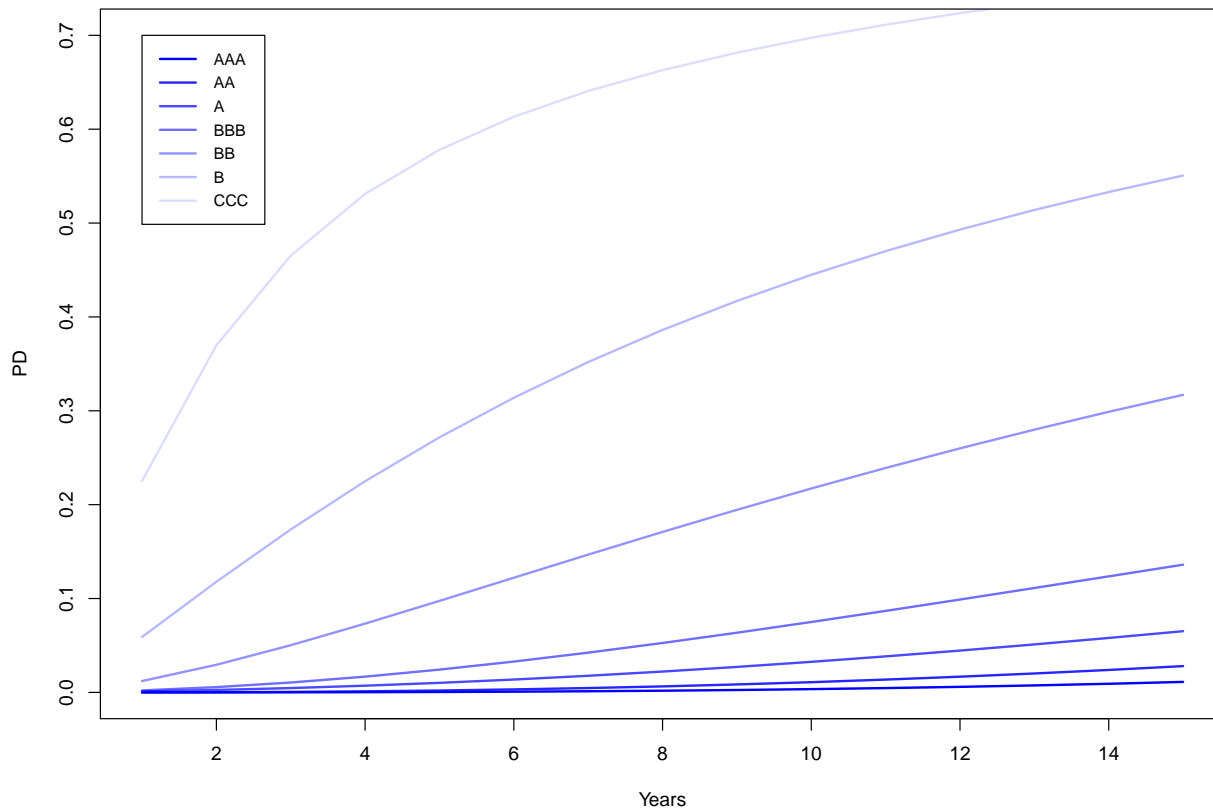
1 ### plotting the term structure of default for each rating
2 library(RColorBrewer)

```

```

3 col<-colorRampPalette(c(rgb(0,0,1,1), rgb(0,0,1,0)), alpha = TRUE)(8)
4 plot(PDt(15,1), type="l", ylim=c(0,0.7), xlab="Years", ylab="PD", lwd=2, col=col[1])
5 for (i in 2:7)
6 {
7   lines(PDt(15,i), lwd=2, col=col[i])
8 }
9 legend(1, 0.7, legend=c("AAA", "AA", "A", "BBB", "BB", "B", "CCC"),
10       col=col, lty=1, lwd = 2, cex=0.8)

```



3. What do you observe?

We observe that in all the cases, the term structure of the probability of default is increasing. Moreover, we can see that investment grade ratings tend to have a convex term structure of probability of default whereas high yield ratings have a concave one.

Exercise 2: Load and explore the "German Credit" database.

1. Load the "German Credit" database available here http://defaultrisk.free.fr/data/TD2_2.csv. Name it `Credit` and display the first ten rows.

```

1 ### Load the German Credit database
2 Credit <- read.csv("http://defaultrisk.free.fr/data/TD2_2.csv")
3
4 ### Display the first ten rows

```

```

5 head(Credit)
6
7 ### You can also display the whole table in a nice way (Rstudio)
8 View(Credit)

```

2. Display the size of the dataframe Credit.

```

1 ### Get the dimension of the dataset Credit
2 dim(Credit)

```

The dataset Credit consists in 1000 observations and 18 variables (17 predictors and 1 target variable).

3. Which columns are numerical? Which are categorical?

```

1 ### Collect the class of each column in Credit
2 col_class <- as.character(lapply(Credit, "class"))
3
4 # or even simpler and safer with the library "purrr"
5 library(purrr)
6 col_class <- map_chr(Credit, class)
7
8 ### Count the columns from each class and display the result
9 table(col_class)
10
11 ### Get the names of the columns from each class and display them
12 writeLines(c("The categorical variables are:",
13             colnames(Credit)[col_class == "factor"],
14             "\n"))
15
16 writeLines(c("The numerical variables are:",
17             colnames(Credit)[col_class == "integer"],
18             "\n"))

```

There are 3 numerical variables (Age, No_Credits and No_dependents) and 15 categorical variables.

4. The variable Default is the variable we want to predict. Is the dataset balanced?

```

1 ### Count the number of observations from each category
2 table(Credit$Default)
3
4 ### Check how this binary variable will be encoded
5 contrasts(Credit$Default)
6
7 ### Relevel Credit$Default to have category "Default" encoded as 1
8 Credit$Default <- relevel(Credit$Default, ref="No Default")
9 contrasts(Credit$Default)

```

The target variable Default is binary ("Default" or "No Default"). The dataset is imbalanced since 70% of the observations belong to the category explicitly named "No Default" and 30% to the category explicitly named "Default". However, it can be considered as relatively balanced for a credit dataset. Indeed, the proportion of "Default" is often much lower. Hence, we will not consider it as an issue for the rest of the exercises.

5. Split the data set into a training set and a test set (70-30%) keeping the same proportion of "Default" and "No Default" in both sets.

```

1 ### Install and load the package named caret
2 install.packages("caret")
3 library(caret)
4
5 ### Set the seed to make results reproducible

```

```

6 set.seed(123)
7
8 ### Get the indices of a training set that maintain the same proportion of "No Default
   "
9 i_train<-createDataPartition(Credit$Default,p=0.7)[[1]]
10
11 ### Split the whole data set into a train set and a test set
12 Train_set <- Credit[i_train,]
13 Test_set <- Credit[-i_train,]
14
15 ### Verify that the stratification worked
16 table(Train_set$Default)
17 table(Test_set$Default)

```

Note that in order to make the results reproducible, we set the seed of the random number generator before randomly splitting the dataset. The same procedure must be followed for non-deterministic algorithms like bagging or random forests.

Exercise 3: Predict Default using a logistic regression.

1. Fit a logistic regression on the training set to predict the binary variable `Default` using the whole set of predictors. Fit another logistic regression using only `Age` and `Status`.

```

1 ### Fit a logistic regression using all the predictors and display the results
2 log_mod <- glm(Default ~ ., data = Train_set, family = binomial("logit"))
3
4 ### Fit a logistic regression using only the predictors "Age" and "Status" (easier to
   interpret)
5 log_mod1<-glm(Default~Age+Status,data=Train_set, family=binomial("logit"))

```

If we fit a model using only `Age` and `Status`, we can see that there is only three out of the four categories from `Status` that actually have a coefficient. The category "Female" is put aside to circumvent colinearity between predictors.

2. What are the significant predictors? As a matter of simplicity use only the predictors `Age` and `Status` ?

```

1 ### Show the summary of log_mod1 (see above)
2 summary(log_mod1)
3
4 ### you can display the results in anexportable way
5 library(broom)
6 tidy(log_mod1)

```

We can see that if we set the threshold at $\alpha = 5\%$ for the coefficient nullity test, only `Status - Male divorced` is considered not significant.

3. Fit a lasso logistic regression and optimize the regularization parameter using cross-validation. AUC will be used to select the regularization/penalty parameter.

```

1 ### Load the packages glmnet to perform the lasso regression
2 library(glmnet)
3
4 ### Create a design matrix of the predictors
5 X_Train<-model.matrix(Default~.,Train_set)
6 X_Test<-model.matrix(Default~.,Test_set)
7
8 ### Using cross-validation, fit lasso regression with different penalties
9 cv_lasso_mod <- cv.glmnet(X_Train, y=Train_set$Default, alpha=1,type.measure="auc",
   family="binomial")
10

```

```

11 ### Plot the results
12 plot(cv_lasso_mod)
13
14 ### Identify the penalty that gives the best AUC
15 (best_lambda <- cv_lasso_mod$lambda.min)
16
17 ### Fit a lasso regression with the best lambda on the whole train set
18 best_lasso_mod<- glmnet(X_Train, y=Train_set$Default, alpha=1, family="binomial", lambda
    =best_lambda)
19
20 ### Show the predictors with non-zero coefficient
21 library(broom)
22 tidy(best_lasso_mod)

```

We can see that out of the 50 variables in the model, only 35 have non-zero coefficients thanks to regularization.

4. Compare the AUC of the two models (before and after regularization) on the test set and conclude.

```

1 ### Load the package "ROCR"
2 library(ROCR)
3
4 ### Make predictions on the test set
5 fit_log <- predict(log_mod, newdata=Test_set, type="response")
6 fit_lasso <- predict(best_lasso_mod, newx=X_Test, type="response")
7
8 ### Create intermediary prediction object to compute the AUC
9 pred_log<-prediction(fit_log, Test_set$Default)
10 pred_lasso<-prediction(fit_lasso, Test_set$Default)
11
12 ### Compute and display the AUC for both models
13 (AUC_log=performance(pred_log, measure="auc")@y.values[[1]])
14 (AUC_lasso=performance(pred_lasso, measure="auc")@y.values[[1]])

```

Exercise 4: Predict Default using tree-based methods.

1. Fit a tree to classify the variable Default.

```

1 ### Load library trees
2 library(tree)
3
4 ### Train a tree on the training set using the Gini index
5 basic_tree <- tree(Default~., data=Train_set, split="gini")

```

2. Display the confusion matrix of the predictions on the test set.

```

1 ### Make prediction on the test set
2 pred_basic_tree <- predict(basic_tree, Test_set, type = "class")
3
4 ### Confusion matrix
5 (confusion_tree <- table(pred_basic_tree, Test_set$Default))

```

3. Try to improve your classifier with bagging (optional: modify the cutoff to improve your detection of Default). Why modifying the cutoff is of particular interest when trying to predict default?

```

1 ### Load the package randomForest
2 library(randomForest)
3
4 ### Set the seed to make results reproducible
5 set.seed(123)
6

```

```

7 ### Train a bagging algorithm on the training set (cutoff is set to improve "Default"
  detection)
8 bagging_tree <-
9   randomForest(
10    Default ~ .,
11    data = Train_set,
12    mtry = ncol(Train_set) - 1, #number of predictors used at each split
13    importance = TRUE,
14    cutoff = c(0.2, 0.8), # see ?randomForest
15    ntree = 500
16   )
17
18 ### Make predictions on the test set and display the confusion matrix
19 pred_bagging_tree <- predict(bagging_tree, Test_set, type = "class")
20 (confusion_bagging <- table(pred_bagging_tree, Test_set$Default))

```

4. Use a random forest algorithm to improve your classifier. What are the most important variables?

```

1 ### Set the seed to make results reproducible
2 set.seed(123)
3
4 ### Train a random forest algorithm on the training set (cutoff is set to improve "
  Default" detection)
5 RF <-
6   randomForest(
7     Default ~ .,
8     data = Train_set,
9     mtry = round(sqrt(ncol(Train_set) - 1), 0),
10    importance = TRUE,
11    cutoff = c(0.2, 0.8),
12    ntree = 500
13   )
14
15 ### Make predictions on the test set and display the confusion matrix
16 pred_RF <- predict(RF, Test_set, type = "class")
17 (confusion_RF <- table(pred_RF, Test_set$Default))
18
19 ### Display the variable importance Plot
20 importance(RF)
21 varImpPlot(RF)

```