



Predicting Cross-Sell with Artificial Neural Networks

An Empirical Study of ING's Customer Data

Seminar Thesis

submitted to

Hon.-Prof. Dr. Martin Schmidberger
Gabriela Alves Werb

Goethe University Frankfurt am Main
School of Business and Economics
Chair for E-Commerce

by

Lukas Jürgensmeier
(Mat.-Nr.: 6904281)

in partial fulfillment of the requirements
for the degree of

Master of Science in Business Administration

July 31, 2019

Contents

1	Introduction	1
2	Introduction from different file	2
3	Main part	8
3.1	Literature overview and citation	8
3.2	Theory and methods	8
3.3	Data	9
3.4	Empirical Analysis	9
4	Conclusion	9
A	R Code	10
B	Formatting rules	10

List of Figures

1	Mean Bond-Yield-curve (example for a figure)	2
2	Example from the CSCC lecture	3
3	Mean Bond-Yield-curve (example for a figure)	11

List of Tables

1	Cointegration of Bond Yields (example for a table)	11
---	--	----

1 Introduction

The introduction should directly lead to the main topic of the paper. It should not be a historical essay or a deep reaching explanation of the topic, but it should explain concisely what the main questions of the topic are, why they are interesting, and which methods or data will be used. A further goal of the introduction is to define the structure for the paper. This can be achieved by describing the goals, the methods and the main results of the paper. Methods and results do not have to be discussed in detail - this is left to the main part of the paper - but they should be summed up in a short way. The introduction of a paper is often finished by a short “roadmap”. This is not necessary, if the aspects mentioned above have been laid out in a satisfactory way before. (Hastie et al. 2017) This is a reference test from the new .tex file 1

2 Introduction from different file

The introduction should directly lead to the main topic of the paper. It should not be a historical essay or a deep reaching explanation of the topic, but it should explain concisely what the main questions of the topic are, why they are interesting, and which methods or data will be used. A further goal of the introduction is to define the structure for the paper. This can be achieved by describing

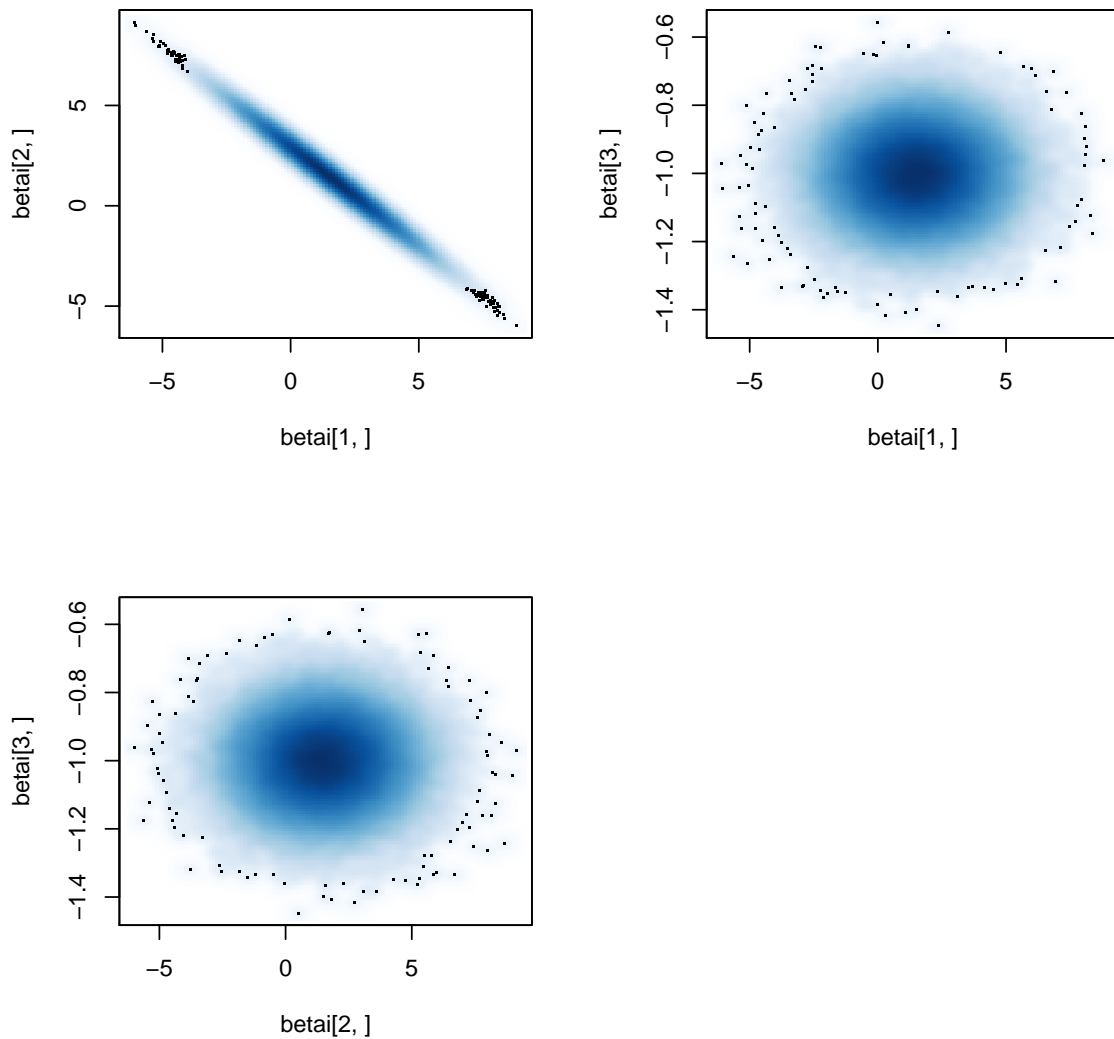
Figure 1: Mean Bond-Yield-curve (example for a figure)



the goals, the methods and the main results of the paper. Methods and results do not have to be discussed in detail - this is left to the main part of the paper - but they should be summed up in a short way. The introduction of a paper is often finished by a short “roadmap”. This is not necessary, if the aspects mentioned above have been laid out in a satisfactory way before. (Hastie et al. 2017)

```
1 # This script follows this blog post
2 # https://blogs.rstudio.com/tensorflow/posts/2018-01-11-keras-customer-churn/
3
4 # clear workspace
5 rm(list = ls())
6
7 #install packages
8 #pkgs <- c("keras", "lime", "tidyquant", "rsample", "recipes", "yardstick", "corrr")
9 #install.packages(pkgs)
10
11 # Load libraries
12 library(keras)
13 library(lime)
14 library(tidyquant)
15 library(rsample)
16 library(recipes)
17 library(yardstick)
18 library(corrr)
19
```

Figure 2: Example from the CSCC lecture



I have no clue what this figure is all about.

```

20
21 # Install Keras if you have not installed before
22 install_keras(method = "conda")
23
24 # read and check data
25 xsell_data_raw <- read.csv("xsell.csv")
26 glimpse(xsell_data_raw)
27
28 # create new variable tenure
29 xsell_data_raw$tenure <- xsell_data_raw$age - xsell_data_raw$entry_age
30
31 # prune data set
32 # Remove unnecessary data
33 xsell_data_tbl <- xsell_data_raw %>%
34   select(-X) %>% #removes ID
35   drop_na() %>% # removes all NA's. Bad Solution! Improve! Removes 70% of
36   # select(xsell, everything())

```

```

37 glimpse(xsell_data_tbl)
38
39 # Split test/training sets
40 set.seed(123)
41 train_test_split <- initial_split(xsell_data_tbl, prop = 0.8)
42 train_test_split
43
44 # Retrieve train and test sets
45 train_tbl <- training(train_test_split)
46 test_tbl <- testing(train_test_split)
47
48 # skipped all feature transformations here
49 # insert if necessary
50
51 # # alternative way for dummy coding all non-numeric variables
52 # non_numeric_var_names <- xsell_data_tbl %>%
53 #   select_if(negate(is.numeric)) %>%
54 #     names
55 #
56 # xsell_data_tbl <- dummy_cols(xsell_data_tbl, non_numeric_var_names)
57 #
58 # # remove non-numeric variables
59 # xsell_data_tbl <- xsell_data_tbl %>%
60 #   select(-non_numeric_var_names)
61 # glimpse(xsell_data_tbl)
62
63 # Create recipe
64 rec_obj <- recipe(xsell ~ ., data = train_tbl) %>%
65   #step_discretize(tenure, options = list(cuts = 6)) %>%
66   #step_log(TotalCharges) %>%
67   step_dummy(all_nominal(), -all_outcomes()) %>%
68   step_center(all_predictors(), -all_outcomes()) %>%
69   step_scale(all_predictors(), -all_outcomes()) %>%
70   prep(data = train_tbl)
71
72 # Apply recipe to predictors (all vars excluding xsell)
73 x_train_tbl <- bake(rec_obj, new_data = train_tbl) %>% select(-xsell)
74 x_test_tbl <- bake(rec_obj, new_data = test_tbl) %>% select(-xsell)
75 glimpse(x_train_tbl)
76
77 # define response variables for training and testing sets
78 y_train_vec <- pull(train_tbl, xsell)
79 y_test_vec <- pull(test_tbl, xsell)
80
81
82 # Building our Artificial Neural Network
83 model_keras <- keras_model_sequential()
84
85 model_keras %>%
86
87   # First hidden layer
88   layer_dense(
89     units = 16,
90     kernel_initializer = "uniform",
91     activation = "relu",
92     input_shape = ncol(x_train_tbl)) %>%
93
94   # Dropout to prevent overfitting
95   layer_dropout(rate = 0.1) %>%
96
97   # Second hidden layer
98   layer_dense(
99     units = 16,
100     kernel_initializer = "uniform",
101     activation = "relu") %>%
102
103   # Dropout to prevent overfitting
104   layer_dropout(rate = 0.1) %>%
105
106   # Output layer

```

```

107 layer_dense(
108   units           = 1,
109   kernel_initializer = "uniform",
110   activation       = "sigmoid") %>%
111
112 # Compile ANN
113 compile(
114   optimizer = 'adam',
115   loss      = 'binary_crossentropy',
116   metrics   = c('accuracy')
117 )
118
119 keras_model
120
121 history <- fit(
122   object      = model_keras,
123   x           = as.matrix(x_train_tbl),
124   y           = y_train_vec,
125   batch_size  = 50,
126   epochs      = 35,
127   validation_split = 0.30
128 )
129
130 # Print a summary of the training history
131 print(history)
132
133 # Plot the training/validation history of our Keras model
134 plot(history)
135
136 # Make predictions
137 # Predicted Class
138 yhat_keras_class_vec <- predict_classes(object = model_keras, x = as.matrix(x_test_tbl)) %>%
139   as.vector()
140
141 # Predicted Class Probability
142 yhat_keras_prob_vec <- predict_proba(object = model_keras, x = as.matrix(x_test_tbl)) %>%
143   as.vector()
144
145 # Evaluate model
146 # Format test data and predictions for yardstick metrics
147 estimates_keras_tbl <- tibble(
148   truth      = as.factor(y_test_vec), # %>% fct_recode(yes = "1", no = "0"),
149   estimate   = as.factor(yhat_keras_class_vec), # %>% fct_recode(yes = "1", no = "0"),
150   class_prob = yhat_keras_prob_vec
151 )
152
153 estimates_keras_tbl
154
155 # change default positive=0 to positive=1
156 options(yardstick.event_first = FALSE)
157
158 # Confusion Table
159 estimates_keras_tbl %>% conf_mat(truth, estimate)
160
161 # Accuracy
162 estimates_keras_tbl %>% metrics(truth, estimate)
163
164 # AUC
165 estimates_keras_tbl %>% roc_auc(truth, class_prob)
166
167 # Precision
168 tibble(
169   precision = estimates_keras_tbl %>% precision(truth, estimate),
170   recall    = estimates_keras_tbl %>% recall(truth, estimate)
171 )
172
173 # F1-Statistic
174 estimates_keras_tbl %>% f_meas(truth, estimate, beta = 1)

```

```

175
176
177 #####
178 ### Evaluate Feature Importance with LIME #####
179 #####
180
181 # Setup
182 class(model_keras)
183
184 #Setup lime::model_type() function for keras
185 model_type.keras.engine.sequential.Sequential <- function(x, ...) {
186   return("classification")
187 }
188
189 # Setup lime::predict_model() function for keras
190 predict_model.keras.engine.sequential.Sequential <- function(x, newdata, type, ...) {
191   pred <- predict_proba(object = x, x = as.matrix(newdata))
192   return(data.frame(Yes = pred, No = 1 - pred))
193 }
194
195
196
197 # Test our predict_model() function
198 predict_model(x = model_keras, newdata = x_test_tbl, type = 'raw') %>%
199   tibble::as_tibble()
200
201 # Run lime() on training set
202 explainer <- lime::lime(
203   x           = x_train_tbl,
204   model       = model_keras,
205   bin_continuous = FALSE
206 )
207
208 # Run explain() on explainer
209 explanation <- lime::explain(
210   x_test_tbl[1:10, ],
211   explainer      = explainer,
212   n_labels       = 1,
213   n_features     = 4,
214   kernel_width   = 0.5
215 )
216
217 # Plot feature importance
218 plot_features(explanation) +
219   labs(title = "LIME Feature Importance Visualization",
220        subtitle = "Hold Out (Test) Set, First 10 Cases Shown")
221
222 plot_explanations(explanation) +
223   labs(title = "LIME Feature Importance Heatmap",
224        subtitle = "Hold Out (Test) Set, First 10 Cases Shown")
225
226 # Feature correlations to xsell
227 corrr_analysis <- x_train_tbl %>%
228   mutate(xsell = y_train_vec) %>%
229   correlate() %>%
230   focus(xsell) %>%
231   rename(feature = rowname) %>%
232   arrange(abs(xsell)) %>%
233   mutate(feature = as_factor(feature))
234
235 corrr_analysis
236
237 # Correlation visualization
238 corrr_analysis %>%
239   ggplot(aes(x = xsell, y = fct_reorder(feature, desc(xsell)))) +
240   geom_point() +
241   # Positive Correlations - Contribute to churn
242   geom_segment(aes(xend = 0, yend = feature),
243               color = palette_light()[[2]],
244               data = corrr_analysis %>% filter(xsell > 0)) +

```



```

245 geom_point(color = palette_light()[[2]],
246             data = corrr_analysis %>% filter(xsell > 0)) +
247 # Negative Correlations - Prevent churn
248 geom_segment(aes(xend = 0, yend = feature),
249              color = palette_light()[[1]],
250              data = corrr_analysis %>% filter(xsell < 0)) +
251 geom_point(color = palette_light()[[1]],
252            data = corrr_analysis %>% filter(xsell < 0)) +
253 # Vertical lines
254 geom_vline(xintercept = 0, color = palette_light()[[5]], size = 1, linetype = 2) +
255 geom_vline(xintercept = -0.25, color = palette_light()[[5]], size = 1, linetype = 2)
256 +
257 geom_vline(xintercept = 0.25, color = palette_light()[[5]], size = 1, linetype = 2)
258 +
259 # Aesthetics
260 theme_tq() +
261 labs(title = "Cross Sell Correlation Analysis",
262       subtitle = paste("Positive Correlations (contribute to xsell)",
263                        "Negative Correlations (prevent xsell)"),
264       y = "Feature Importance")

```

Listing 1: Tensorflow Model

3 Main part

3.1 Literature overview and citation

The literature overview may be kept short in a bachelor thesis. Bachelor theses whose main purpose is to present and discuss the contents of published articles are exceptions. When referring to articles or other literature, it is essential to mark these as references. This is best

$$f(x) = x^2 \tag{1}$$

$$F(x) = \int_b^a \frac{1}{3}x^3 \tag{2}$$

done in the text itself. When referring to papers, the author and the year of publication should be given, e.g. “Imbens (2002) gives an overview for the GMM-estimator and its empirical likelihood”. See for example in (1) for the squared one and (2) for an example of a fancy integral.

If the paper was written by more than two authors, this fact is usually abbreviated as, for example, “Imbens et al. (2002)”. If there was more than one publication in the same year, a small letter should be added to the year such as “Imbens (1997a)”. When referring to a whole chapter, the chapter should be mentioned, e.g. “Wooldridge (2002), ch. 13”.

Direct citations must be enclosed in quotation marks. In this case the year of publication should be added with the author’s name, e.g. “Generalized method of moments (GMM) estimation has become an important unifying framework for inference in econometrics in the last 20 years (Imbens (2002), p. 493)”. The use of direct citations should be kept to a minimum.

3.2 Theory and methods

When writing an empirical bachelor thesis, the theoretical part should be limited to an amount necessary to understand the empirical part. It is better to limit the theory to the special cases rather than striving for a maximum of generality. Of course, when writing a theoretical paper, or a paper on pure methods, the theoretical part will receive more weight.

However, the presentation should always be structured so that it clearly works out the main points, concentrating on the aspects that are really central to the topic. Detailed proofs should be moved to the appendix.

3.3 Data

When writing an empirical paper, it is necessary to give a concise description of the data set that is being used. This description should include information about the data set provider and the variables used. A descriptive analysis of the data is useful, but it may also be moved to the appendix.

3.4 Empirical Analysis

In the empirical section, the main results should be explained first. If this is not possible, because intermediate steps are required to understand the results, then only intermediate results should be explained that are really essential for this purpose. Tables and figures should be used to present the main results. In addition, the tables and figures have to be discussed in the text. Each table and each figure must have their own title and caption.

It is often useful to investigate the robustness of the results with respect to different aspects. If the results were calculated under the homoskedasticity assumption for example, one should discuss what happens if the assumption is violated. More detailed empirical results should be put into the appendix unless there are important reasons not to do so.

4 Conclusion

The conclusion should contain a summary of the main results and its implications. One can also mention directions for future research.

A R Code

```
1 fib <- function(n) {  
2   if (n < 2)  
3     n  
4   else  
5     fib(n - 1) + fib(n - 2)  
6 }  
7 fib(10)
```

Listing 2: Fibonacci Sequence

```
fib <- function(n) {  
  if (n < 2)  
    n  
  else  
    fib(n - 1) + fib(n - 2)  
}  
fib(10)
```

See Test Code 2 for the code example.

B Formatting rules

- Letter size 11 to 12pt, line spacing 1 - 1.5 times letter size, margins left/right 2.5cm, bottom 3cm, top 2.5cm.
- Pre-introduction page numbers should be roman, the ones of the main text arabic.
- The table of contents shows chapters and sections.
- The list of figures and tables list all the figures and tables in the paper.
- Every figure and table should have a short title plus description and should be explained in the text.
- The number of pages of the main text should be between 10-15 pages (appendix excluded).

Figure 3: Mean Bond-Yield-curve (example for a figure)



Table 1: Cointegration of Bond Yields (example for a table)

Rank at least	\mathcal{L}_{trace}	5% crit. value	\mathcal{L}_{max}	5% crit. value
$r_0 = 0$	299.71	76.07	108.08	34.40
$r_0 = 1$	191.64	53.12	91.00	28.14
$r_0 = 2$	100.64	34.91	65.14	22.00
$r_0 = 3$	35.50	19.96	29.29	15.67
$r_0 = 4$	6.21	9.24	6.21	9.24

The Akaike-Information-criterion suggests a maximal lag length of 14

References

Hastie, T., Tibshirani, R. & Friedman, J. (2017), *The Elements of Statistical Learning: Data Mining, Inference, and Prediction*, Springer Series in Statistics, 2 edn, Springer-Verlag, New York.

Statutory Declaration

I herewith declare that I have completed the present thesis independently, without making use of other than the specified literature and aids. Sentences or parts of sentences quoted literally are marked as quotations; identification of other references with regard to the statement and scope of the work is quoted. The thesis in this form or in any other form has not been submitted to an examination body and has not been published. This thesis has not been used, either in whole or part, for another examination achievement.

Frankfurt am Main, July 31, 2019

.....

Lukas Jürgensmeier