



Un-Black-Boxing Artificial Neural Networks
Predicting and Explaining Bank Customer's Cross-Sell Likelihood

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

in partial fulfillment of the requirements
for the degree of

Master of Science in Business Administration

July 15, 2019

Contents

1	Introduction	1
2	A Short Primer on Cross-Sell	1
3	Methodology	2
3.1	Feature Engineering and Data Set Preparation	2
3.2	Artificial Neural Network with Keras	3
3.3	Cross-validation and Hyperparameter Tuning	4
3.4	Local Interpretable Model-agnostic Explanation (LIME)	6
4	Results and Model Comparison	7
4.1	Choosing the Neural Network Architecture	8
4.2	Comparing the Predictions of the ANN with the Benchmark Logit Model	9
4.3	Beyond the Black Box: Feature Importance with LIME	10
5	Managerial and Research Implications	13
6	Conclusion	15
A	Appendix	16
B	R Code	17

List of Figures

1	Architecture of the Artificial Neural Network implemented in this study	4
2	Loss and accuracy for training and test data over training history of the tuned model	5
3	ROC curve comparison between benchmark logit (black) and hyperparameter- tuned neural network predictions (red)	11
4	LIME feature importance of the ten most important features for each of the first four customers in the test data set	12
5	Individual variable assessment—violin plot of xsell vs. age	13
6	Correlation of feature engineered variables with xsell	16

List of Tables

1	Resulting statistics of 324 hyperparameter tuning model runs	8
2	Model performance comparison—statistics and confusion matrix	10

List of Abbreviations

ANN	Artificial Neural Network
API	Application Programming Interface
AUC	Area Under Curve
LIME	Local Interpretable Model-agnostic Explanation
ReLU	Rectified Linear Unit
ROC	Receiver Operating Characteristics
SP	Submodular Pick

1 Introduction

The rise of machine learning techniques has come with a great increase of predictive accuracy in many settings (Jordan & Mitchell 2015). The hype around those techniques has led to a widespread use of them in many fields in science and practice—from medicine to business or self-driving cars (Chen & Asch 2017). However, some of those methods are increasingly computationally expensive and might not in every case be superior to traditional regression-based methods such as simple multivariate linear regression or logit models.

This paper uses a data set from a large German bank with numerous customer characteristics to predict whether an existing customer buys another product—also referred to as cross-selling. The goal of this research is first to predict whether a customer will cross-buy and open a checking account and second to identify certain characteristics of typical customers that increase the likelihood of opening that account. Therefore, this paper introduces a complex Artificial Neural Network (ANN) with optimized predictive power and compares the predictions to a simple logit model. One downside of ANN models in particular the lack of interpretability—they are widely regarded as black-boxes (Benítez et al. 1997, Dayhoff & DeLeo 2001). This paper addresses this shortcoming by implementing LIME—a method to explain individual predictions of neural networks (Ribeiro et al. 2016). The paper is structured as follows: The first chapter gives a short introduction to cross-selling theory and is followed by the explanation of the implemented statistical methods. The subsequent chapter presents the empirical results and especially compares the ANN to the simple logit model. The paper concludes with implications for practitioners and outlines further research avenues.

2 A Short Primer on Cross-Sell

Cross-selling is an important action field for established companies with an existing customer base (Li et al. 2011). Selling additional products to present customers taps into a readily accessible target group for additional sales potential. Felvey (1982) argues that a businesses’ client base might already be more inclined to buy another product from the company that they already have a relationship with, compared to non-

customers. Also, the cost of selling additional products to already existing customers is cheaper than targeting new ones (Reichheld & Sasser Jr 1990). Benefits of cross-selling from the company point-of-view include increased switching costs for customers, while customers benefit from an easy integration of different products and a higher certainty about the quality of the new product (Kamakura et al. 1991, 2003). However, as in most business settings, there is no free lunch and hence cross-selling efforts are also associated with costs and risks. Apart from trivial costs for cross-selling efforts such as mailing or telephone campaigns, there is a significant churn risk involved. If not accurately targeted, customers can be annoyed by irrelevant and continuous advertisement and might eventually churn as a result (Keaveney 1995). Due to the significant sales potential on the one hand and costs and risks involved on the other hand, it is therefore paramount to accurately identify the target group that is most likely to cross-buy another product. This paper addresses this issue and presents an Artificial Neural Network to accomplish that task in the best possible way.

3 Methodology

This section provides an overview over the implemented methodology by first outlining the feature engineering process that transforms the original data set to suitable format for neural networks. The second part then describes the ANN and its initial architecture, while the third part describes the hyperparameter tuning process that leads to the *best* model¹. Lastly, this section addresses a common criticism of machine learning technologies in general and neural networks in particular—the uninterpretable black box—by introducing the Local Interpretable Model-agnostic Explanation (LIME) as a method to "look under the hood" of a neural network and determine which features contributed by how much to a prediction.

3.1 Feature Engineering and Data Set Preparation

In order to feed seamlessly into a neural network, extensive feature engineering and data preparation is required (Hastie et al. 2017). Before training the model, I replace

¹This model is definitely not the best model *existing*, but the best model *found* through the trial-and-error tuning process described in detail in section 3.3.

missing values in `pref_device` and `occupation` by `None_or_missing`, transform all `char` variables to a `factor` type, and remove the misleading feature `ID`. However, there are still numerous features including missing values. Since they are not missing at random, they need to be imputed in order to not bias the prediction. Deleting all observations with missing values would do exactly that. Hence, I impute missing values with a Random Forest from the `randomForest::rfImpute()` function. This algorithm first replaces missing numeric (factor) variables with the median (mode) and then uses proximity measures from `randomForest` as weights to replace NA's with the weighted average of the non-missing values and repeats this for a pre-defined number of iterations² (Liaw & Wiener 2002). I furthermore bin the continuous features `age`, `entry_age` and `last_acc`³ in four distinct groups. Since a neural network cannot handle multi-categorical features, the routine dummies-codes those. Finally—to not distort the input weights for the first hidden layer and thus not negatively affect the prediction quality (Hastie et al. 2017, pp. 398)—I standardize all input features to exhibit $\bar{x} = 0$ and $s = 1$. The feature engineered data set consists of $n_x = 74$ input variables that will be used to predict cross-sell.

3.2 Artificial Neural Network with Keras

While the previous section discussed the steps required to transform the data to a suitable format for an ANN, this part describes the basic architecture of this classification method. This study uses the R interface `keras` (Chollet & Allaire 2017) to the popular same name Python library. As a high-level API, it enables an easier way to access the features of Google's `TensorFlow` (Abadi et al. 2015).

The neural network in this study is a common feedforward network which uses the backpropagation algorithm to adjust its weights (Werbos 1990). Based on the feature

²After four iterations the out-of-bag error did not decrease notably anymore, hence four iterations were chosen.

³Those three features exhibit a highly non-linear relationship. According to the universal approximation theorem, a neural network with enough units in a hidden layer can approximate any continuous function (Hornik 1991). The best neural network will be found by trying out different hyperparameters, including some combinations which would not have *enough* hidden layer units for a solid approximation. Hence, I bin those three variables to make sure that this non-linear relationship can be taken into account by the model in any case. Also, the benchmark logit model could not capture this effect without further specification.

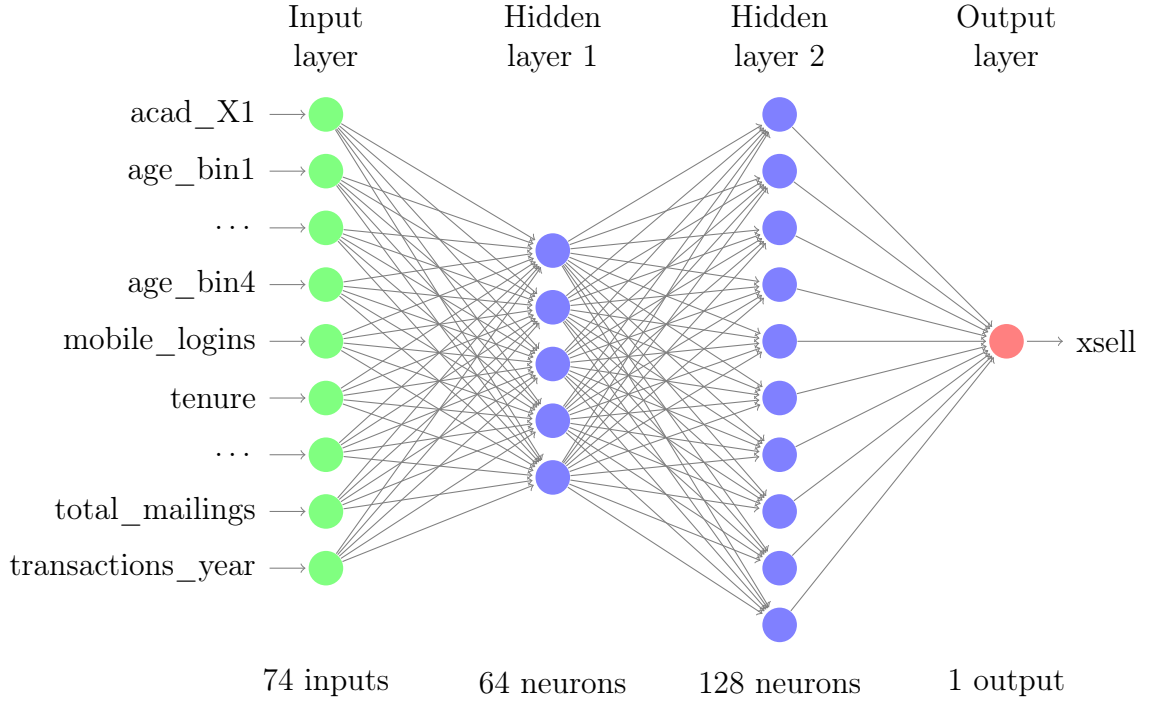


Figure 1: Architecture of the Artificial Neural Network implemented in this study

engineering process described in the previous section, the ANN has $n_x = 74$ inputs X_i , two hidden layers with neurons H_i^l referring to the i^{th} neuron in the l^{th} hidden layer. The number of neurons in the hidden layers is subject to hyperparameter tuning in the following section. As a final layer, the ANN predicts the single binary output \hat{Y}_i , in our case `xsell`. Figure 3.2 graphically visualizes the architecture. This network uses a random uniform distribution $U(-0.05, 0.05)$ to initialize the weights for each neuron H_i^l . Also, a bias initialized at zero feeds into the two hidden layers, for which each H_i^l is activated by a rectifier function (ReLU). Glorot et al. (2011) show that rectifying units in an ANN yield superior outcomes compared to others such as the hyperbolic tangent activation function. Since the output \hat{Y}_i is a binary feature, I use a sigmoid activation function for the output layer. To avoid overfitting, the model implements dropout for the hidden layers. This process randomly drops neurons and their connections from the ANN and thus leads to a more generalizable model (Srivastava et al. 2014). Mitigating the risk of overfitting to the training data is furthermore a key objective of the next section.

3.3 Cross-validation and Hyperparameter Tuning

After setting up the basic architecture of the neural network, this section describes the process that leads to the *best* model⁴. An important feature of the *best* model is that it does not memorize the exact relationships in the training data set—overfitting—but rather learns the generalizable underlying mechanisms and thereby performs well on the holdout validation and subsequently on the test or even an out-of-time data set. For this purpose I first describe the optimization of the ANN by cross-validating predictions of several training epochs with a holdout sample and then introduce hyperparameter tuning—a process that leads to the most desirable⁵ configuration of the ANN by running many combinations of different model specifications.

To compile the ANN, I use the adaptive moment estimation algorithm **Adam** (Kingma & Ba 2014). Ruder (2016) recommend this gradient descent optimizer as computationally inexpensive and performing best in most applications. This optimizer minimizes the cross-entropy loss (Zhang & Sabuncu 2018) for each training epoch and accordingly updates the weights for the next model. Figure 2 shows the training history of the chosen ANN over ten epochs. The cross-entropy loss is plotted in the top pane, and the accuracy displayed in the bottom part, each for the training and validation data set.

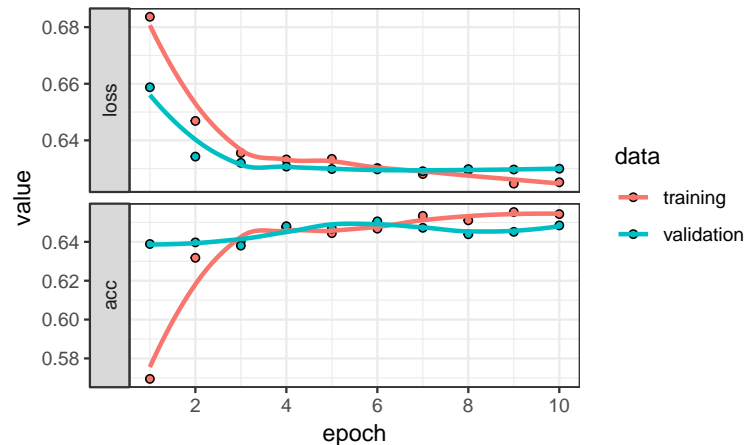


Figure 2: Loss and accuracy for training and test data over training history of the tuned model

⁴as mentioned earlier, this is not the best model *existing*, but *found*.

⁵"Desirable" is defined as most accurate in this study. However, other metrics such as loss are also valid optimization criteria.

The validation accuracy in comparison with the training accuracy in this case indicates first how accurate the ANN predicts true positives and true negatives in both data sets, but more importantly, gives an indication whether overfitting is an issue. If the training accuracy kept steadily increasing, while the validation accuracy remained relatively constant, the model would overfit—which means that the model is not generalizable beyond the training data set (Hansen & Salamon 1990). As exhibited in figure 2, both accuracies are relatively constant and rather similar during the last training epochs. This is a strong indication that the model does not overfit and will perform similarly well on the test or out-of-time data.

However, one question remains unanswered with cross-validation: Which model architecture leads to the most accurate predictions? Hyperparameter tuning provides a remedy to this issue (Bergstra & Bengio 2012). In short, this process runs numerous models with different combinations of hyperparameters⁶. The model with the highest validation accuracy `val_acc` is then chosen to be the *best* model for predicting `xsell`. This study uses the "naive" grid method, which runs all possible combinations of pre-defined values for hyperparameters⁷. However, this method is computationally expensive. It requires $N = a^b$ different models, where a is the number of discrete values for each parameter and b is the number of hyperparameters to be tuned. Running two separate sets of hyperparameter combinations for 324 different models in total took roughly five hours on a standard consumer laptop. It can be assumed that this limited set of combinations did by far not yield the best model. However, this naive grid method runs an exponentially increasing number of different models and thus gets difficult to handle quite quickly. A possible solution to this problem is random search, which is even more efficient than grid search (Bergstra & Bengio 2012).

By cross-validating the predictions with a loss function and by tuning the ANN's hyperparameter, this section described the approach to finding the most suitable model for the prediction task at hand. The next part turns to prediction explanation—a key

⁶Hyperparameters refer to the specifications of an ANN, such as number of units per hidden layer, dropout rate and weights, but also different optimization algorithms and activation functions (Bengio 2000)

⁷I use each three different values of the hyperparameters `val_acc`, `acc`, `units1`, `dropout1`, `units2`, `dropout2`, `epochs` and `batch`

discussion point and critique of neural networks.

3.4 Local Interpretable Model-agnostic Explanation (LIME)

In contrast to traditional regression-based methods that yield beta coefficients and enable the construction of a functional form of the underlying mechanisms in the data, a neural network does not provide an easily interpretable function that links Y_i to X_i . A common criticism of ANN hence is the "black-box" argument (Benítez et al. 1997, Dayhoff & DeLeo 2001). Though the model might yield more accurate predictions, there is no clear answer to the question why the model predicted those particular values or which features contributed by how much.

With the development of new methods that address this problem, an ANN is less a black box. Consistent with the research question to uncover factors that drive cross-sell, this section focuses on the Local Interpretable Model-agnostic Explanation (LIME) introduced by Ribeiro et al. (2016) as a methodology to assess feature importance on a single customer level⁸. LIME uses a trick to come up with an interpretable explanation: While the highly non-linear function of the ANN is unknown over the total space, it can be approximated locally by an easier, readily interpretable model such as a linear regression or a decision tree.

While those models globally do not match the functional form of the neural network at all, they can yield accurate *local* approximations (Ribeiro et al. 2016, see figure 3 for a visual example). Turning to the task at hand, this paper uses a multiple linear regression with the ten most important features for every customer. LIME permutes every observation and calculates the **gower**-distance between the observation and the permutation (Pedersen & Benesty 2018). It then fits a forward-selection procedure on the basis of a ridge regression (Pedersen 2018) to identify the ten most important features, weighted by the distance calculated earlier. Following Ribeiro et al. (2016), the beta coefficients ("weights") are then reported as a locally accurate approximation of feature importance. This method hence enables to go beyond the typical "black-

⁸Other tools such as neural interpretation diagrams, Garson's algorithm, sensitivity analysis or randomized approaches (Olden & Jackson 2002) offer similar explanations on an aggregate level. However, covering those is beyond the scope of this paper.

box" description of a neural network. While this chapter described the implemented methods theoretically, the remainder of this paper reports and interprets the results of those machine learning techniques applied on the cross-sell data set.

4 Results and Model Comparison

This section proceeds as follows: First, it reports the results of the hyperparameter tuning and cross-validation as well as the resulting optimal ANN architecture. The second part then compares the predictions of the neural network to a simple benchmark logit model. It concludes by reporting the feature importance for individual customers derived by LIME.

4.1 Choosing the Neural Network Architecture

The task of choosing the *best* ANN architecture is not trivial. Using the approach described in section 3.3, I ran a total of 324 different models, each with different combinations of hyperparameters. Ranked by validation accuracy, table 1 shows the top five and bottom two models with its respective model specification. It is noteworthy

Table 1: Resulting statistics of 324 hyperparameter tuning model runs

rank	val_acc	acc	units1	dropout1	units2	dropout2	epochs	batch
1	0.656	0.654	64	0.600	128	0.600	10	100
2	0.654	0.692	128	0.600	64	0.600	30	150
3	0.653	0.674	64	0.600	64	0.600	30	100
4	0.651	0.658	64	0.600	64	0.600	10	100
5	0.651	0.673	64	0.600	128	0.600	30	50
...
323	0.604	0.731	64	0.200	32	0.400	30	128
324	0.602	0.735	64	0.200	32	0.200	30	128

that this evaluation metric only varies by five percentage points from the best (65.6%) to the worst model (60.2%). Also, the dropout rate for all five best performing models is 0.6 for both hidden layers, while the number of units in those varies between 64 and 128. In contrast to that the least preferable models exhibit only 32 neurons in the second hidden layer. In this case, a higher number of neurons in those layers appears to

better capture the effects of the input features on `xsell`. The winning model used ten training epochs, that means the model adjusted weights by backpropagation ten times before stopping. This value is rather low⁹ and suggests an advantage of stopping early to avoid overfitting. Finally, the size of the batches—that is the number of samples sent back and forth through the network at once (Chollet & Allaire 2017)—appears to not have an obvious impact on model performance.

Even though this routine identified the best out of 324 models, it only tested three different values for eight hyperparameters in two separate runs. Due to the computational complexity, important parameters such as the *learning rate* were not included. Furthermore, I only tested three more or less arbitrarily chosen discrete values for the hyperparameters—it is therefore obvious that better models do exist. It is however difficult to identify them with the grid method and limited computational resources.

4.2 Comparing the Predictions of the ANN with the Benchmark Logit Model

The ANN from the previous section with the highest validation accuracy is furthermore used to predict `xsell` in the remaining unseen test data set¹⁰. To assess the quality of the model, this paper now compares the resulting model statistics and confusion matrices of the discussed ANN with those of a simple logit model based on the same 74 input features used for the neural network¹¹. Table 2 displays all relevant model performance metrics from accuracy to kappa and each confusion matrix. To a surprise, there is no apparent difference in model accuracy between the hyperparameter-tuned ANN and the simple benchmark logit model. The reference model even predicts `xsell` 0.4 percentage points more accurate¹². Since 50% of all customers in the data set made a purchase, this is the benchmark accuracy for any predictive model. Compared to a

⁹I tested epoch values of 5, 10, and 30 for the tuning routine.

¹⁰20% of the 10,000 observations in the data set were used for testing, while the remaining 8,000 observations form the training set. Out of those, I reserved 2,400 rows (30%) as a holdout validation sample for cross-validation during the fitting process.

¹¹In fact, this logit model has a relatively simple functional form, since it throws in every feature at hand without any selection mechanism or concern for multicollinearity. However, those features have undergone the same extensive feature engineering process as described for the ANN. Hence, a lot of concerns are dealt with, e.g. though not functionally mapped, the logit model incorporates non-linear effects of `age`, `entry_age`, and `last_acc` through binning.

¹²Due to the small test set sample size of 2,000 this effect might not be statistically significant.

coin toss, both the ANN and logit model yield a roughly 29% increase in accuracy. Nevertheless, both models produce structurally different predictions. This becomes

Table 2: Model performance comparison—statistics and confusion matrix

<i>Benchmark Logit</i>			<i>Hyperparameter-tuned ANN</i>		
Accuracy	0.6463		Accuracy	0.6423	
Sensitivity	0.6653		Sensitivity	0.7331	
Specificity	0.6271		Specificity	0.5508	
Pos Pred Value	0.6429		Pos Pred Value	0.6221	
Neg Pred Value	0.6500		Neg Pred Value	0.6716	
Kappa	0.2925		Kappa	0.2841	
Prevalence	0.5023		Prevalence	0.5023	

<i>Prediction</i>	<i>Actual</i>		<i>Prediction</i>	<i>Actual</i>	
	0	1		0	1
0	624	336	0	548	268
1	371	668	1	447	736

clear when comparing sensitivity (true positive rate) and specificity (true negative rate). The ANN is relatively better at predicting customers that actually buy another product ($\mathbf{x_{sell}} = 1$), while the logit model more accurately predicts negative outcomes. This has far-reaching implications for business applications, since certain wrong predictions might be more costly than others¹³.

Another widely used visual tool and metric for model comparison in machine learning is the receiver operating characteristics (ROC) curve and the associated area under curve (AUC) (Fawcett 2006). Figure 3 plots both ROC curves and the AUC in a single plot—the curves are visually inspected almost identical. Quantified by the AUC, the ANN performs 0.9 percentage points better than the logit model at 69.3%.

4.3 Beyond the Black Box: Feature Importance with LIME

The previous section introduced the implemented ANN model specification and reported the results in terms of predictive performance. However, prediction is only one part of the research question. The other equally important question to address is explainability. Out-of-the-box, machine learning algorithms such as ANN are indeed a

¹³The scope of this paper does unfortunately not allow a thorough discussion of different predictive error costs. See Domingos (1999) for a detailed discussion on this topic and a proposed solution that incorporates costs into any classifier.

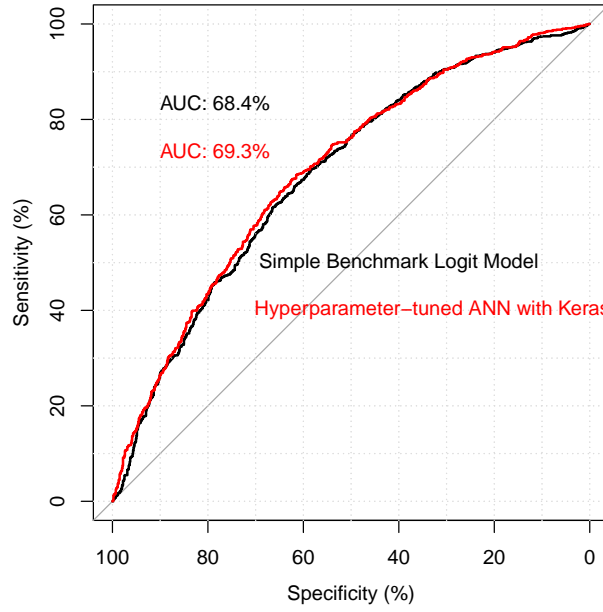


Figure 3: ROC curve comparison between benchmark logit (black) and hyperparameter-tuned neural network predictions (red)

"black-box" and thus it is impossible to understand why the model predicted which outcome. LIME, as introduced theoretically in section 3.4, provides insights into the reasoning of why the model came to a certain conclusion by locally fitting an explainable model and reporting the resulting feature weights as a measure for its contribution to the prediction (Ribeiro et al. 2016). In this paper, I focus on explanations for individual observations, i.e. individual customers. Explaining feature importance on a global level would be the next step. However, the limited scope of this paper does not permit a thorough discussion of both. In a business setting, it is crucial to understand why the model predicted a certain outcome for cross-selling of an individual customer, especially for building human trust in the machine learning algorithm's prediction (Ribeiro et al. 2016). The remainder of this section therefore focuses on this explanation only.

Figure 4 visualizes the results of `lime::explain()` for the first four customers in the test data set. For each of those customers, the ten most important features are reported in descending order of their influence on the model's prediction. *Label* is the model's prediction whether a customer buys another product, *probability* gives a degree of the

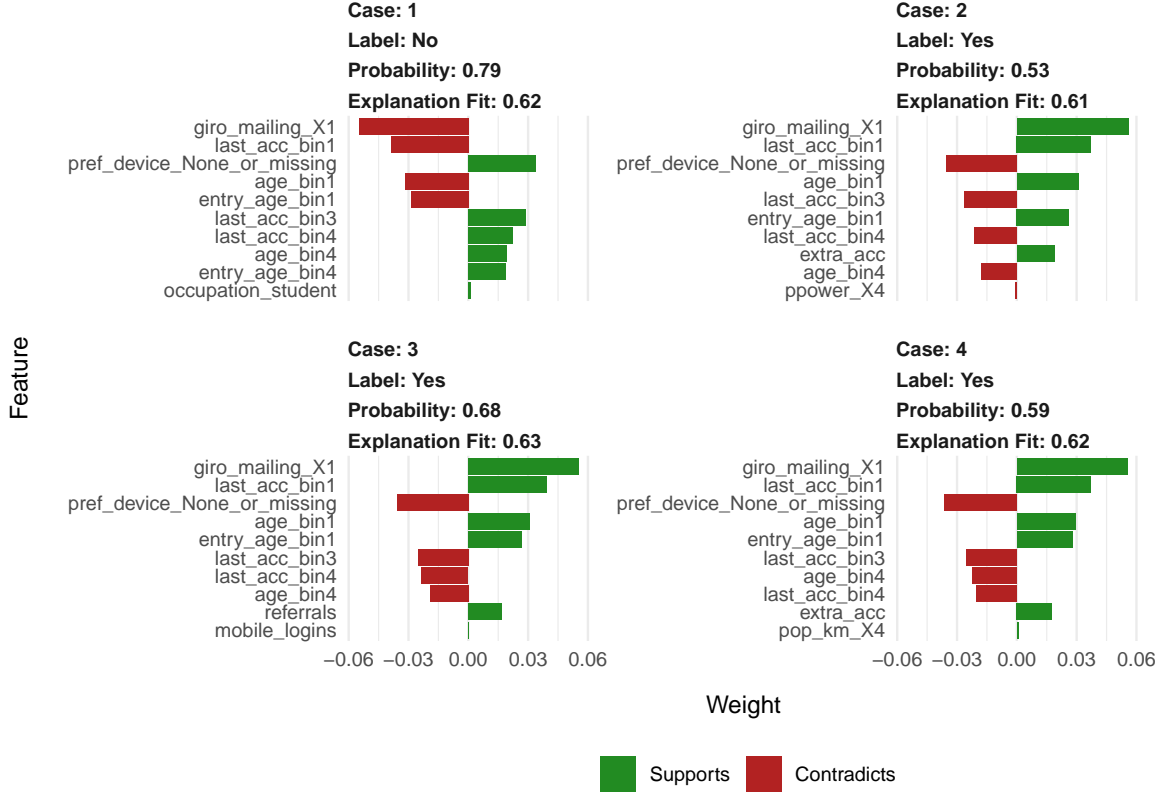


Figure 4: LIME feature importance of the ten most important features for each of the first four customers in the test data set

model’s confidence in the prediction, and *explanation fit* reports the locally approximated model’s R^2 with the ten most important features utilized. It is apparent that the features exhibiting the highest correlation with `xsell` also play the biggest role for those four customers (see figure 6 in the appendix for a graphical representation of correlation in the data set).

A red bar indicates that the particular value of this feature contradicts the model’s conclusion, while a green bar signals supporting evidence to the model’s prediction. Whether the customer received an advertisement mail about opening a checking account (`giro_mailing_X1`) shows the largest weight from the regression model in every one of the four cases. Also the group with the lowest days passed since the customer last opened an account (`last_acc_bin1`) shows a very large contribution to the prediction in those four cases. Additionally, the customers’ age seems to have a big impact on cross-sell predictions—its binned version is consistently included as a top reason for

the predictions. However, note that no generalization on the entire data set is possible from this interpretation, since it only explains individual predictions.

A first step towards a general statement about the customer base’s characteristics is looking at the distribution of features split by whether the customer opened a checking account or not. For **age**, the violin plot in figure 5 shows a striking difference in the age distribution. Younger customers tend to more often cross-buy. Looking at the correlation of all features with **xsell** in figure 6 in the appendix, the youngest age group shows the highest positive correlation after **giro_mailing_X1**. Receiving a mail with an offer to buy another product in every of the four individual explanations is the feature that has the highest influence on the model’s prediction, while it also displays the highest correlation with cross-sell. On the contradicting side, a missing preferred device displays the highest negative correlation, while it also comes up as the third most important feature in the individual LIME explanations. Though a mere correlation does not imply causality, it can provide a first clue about general underlying mechanism and can validate the LIME results. Judging from the first four customers only, the locally fitted models’ weights are in line with the overall correlations.

Combining the LIME explanations with the actual values for those customers would now enable to act on those results. The following section further elaborate on this issue and discusses the managerial implications that those methods bring with.

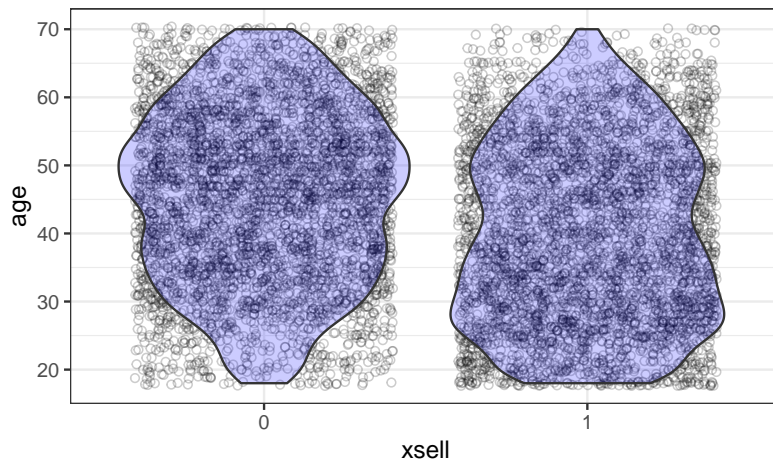


Figure 5: Individual variable assessment—violin plot of xsell vs. age

5 Managerial and Research Implications

This section first highlights managerial implications that can be derived from this study and then outlines future research areas in this field. Looking at the ANN’s accuracy compared to a simple benchmark logit model one recommendation for practice is obvious: Do not use neural networks in this particular business setting. They are computationally expensive and provide no overall predictive advantage to way simpler models such as logit. Possible reasons for the disappointing ANN performance could be the lack of enough data—neural networks require larger data sets compared to well-specified simple regression-based models for accurately learning the underlying relationships (Chan et al. 1999). Also, neural networks are better suited for other tasks such as image recognition (Esteva et al. 2017) and speech (Hinton et al. 2012) or text recognition (Jaderberg et al. 2014). In those cases, a simple multivariate regression would not be suitable.

Nevertheless, LIME provides an important way to understand why the black-box ANN predicted cross-selling. This is important to first get an understanding which features are contributing by how much and in which direction. Out-of-the-box, this would not be possible with ANNs. More importantly however, LIME builds trust in the prediction of an otherwise mysterious algorithm. By fitting a locally reasonable and interpretable model, this new approach enables decision makers to assess the validity of the neural network’s predictions beyond typical performance measures such as accuracy and AUC (Ribeiro et al. 2016).

Still, this paper lacks a generalizeable explanation beyond simple descriptive metrics such as correlation (see figure 6) or the distribution of observations by `xsell` as shown in figure 5. Though individual predictions can be explained by LIME, one cannot derive a truly general statement about the sample from it. In their research on LIME, Ribeiro et al. (2016) suggest submodular picking of individual observation-level explanations in order to derive generalizable insights into the black-box model. Simplified, they propose an algorithm that optimally picks the most representative individual observations for the entire sample. Unfortunately, the LIME-SP routine is so far only implemented in the original `Python` package and currently under development for the corresponding

keras R interface (Kavicky 2017).

This empirical analysis suggests which customer groups are likely to purchase another product. Most importantly, young customers exhibit a higher probability to purchase than old ones. The more recently a customer already cross-bought, the more likely it is that the bank can sell an additional product to this client again. Targeting those customer groups in particular could be a valid business strategy. However, more research would be necessary to establish a clear causal relationship between those variables. A neural network is definitely not the method of choice for this task—even though new developments such as LIME can greatly expand the trust and explainability of black-box models.

6 Conclusion

The first goal of this paper was to predict whether a bank customer opens another checking account based on their individual characteristics. A sophisticated hyperparameter-tuned Artificial Neural Network was therefore introduced and its predictions compared to a simple logit model. Surprisingly, the ANN performed slightly worse than the benchmark logit model in terms of accuracy and only marginally outperformed it in terms of AUC. Due to the enormous computational resources required for finding a satisfactory ANN, this paper argues against the implementation of neural networks for predicting cross-sell in this particular business setting and for simpler methods such as a logit model.

The second research question—identifying underlying characteristics that drive cross-selling—was initially impossible to answer with a black-box ANN. However, this paper implemented a local approximation of the neural network with LIME by an easily interpretable regression model. This enabled detailed insights into which customer features support or prevent cross-selling on an individual customer level. Most importantly, sending mail about opening a new account is the most contributing feature for the four analyzed customers, followed directly by the customer’s age. Younger clients are more likely to cross-buy, which is also supported on the global level by looking at correlations and individual scatter plots. Even though those two methods mostly answer the

research question, there is still more research to be done to establish causality. On a global level, the methodology implemented in this study only reports which features are associated with cross-selling. However, no reasonable statement can be made about which features *cause* an increased cross-selling probability across all customers.

A Appendix

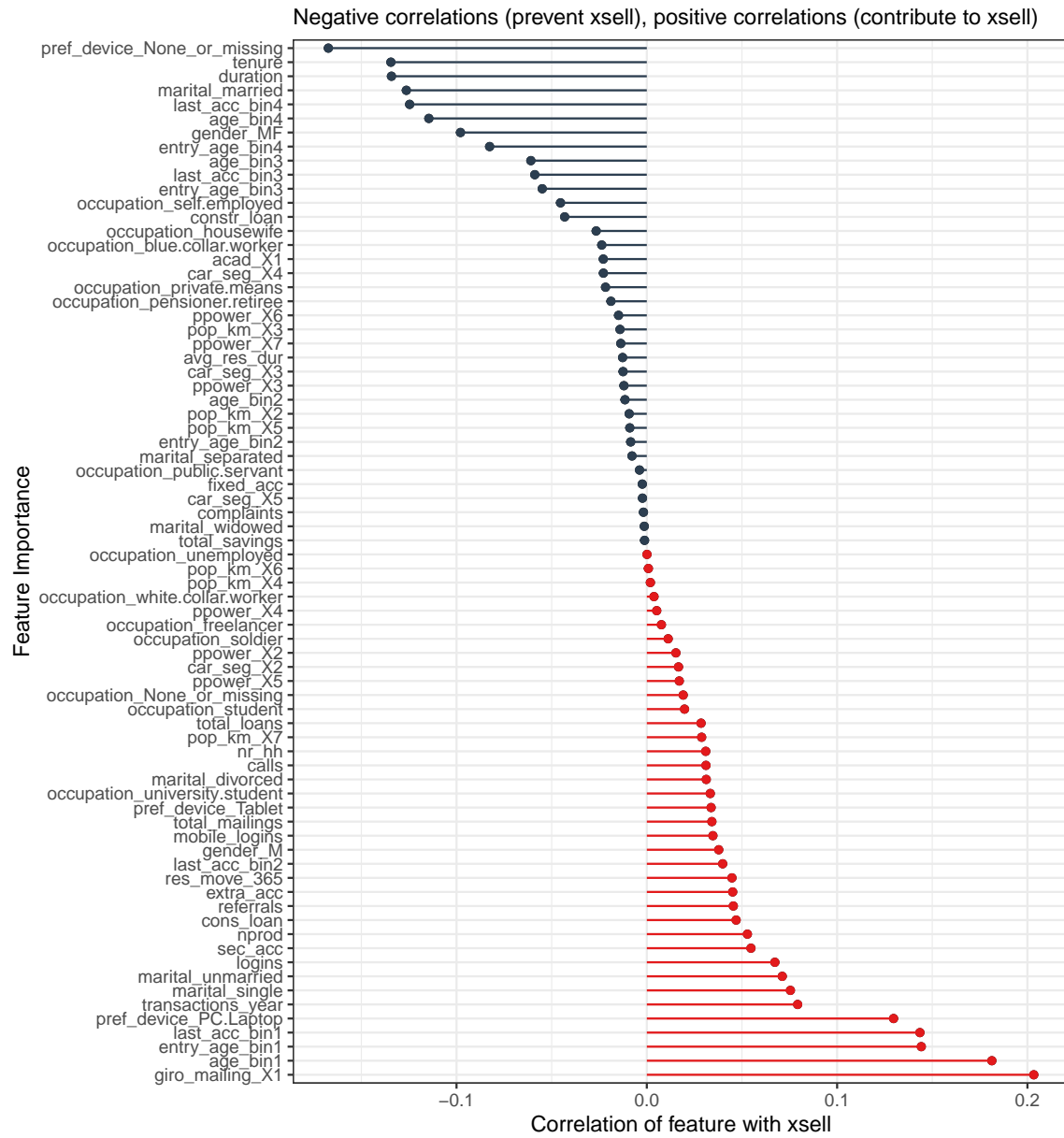


Figure 6: Correlation of feature engineered variables with xsell

B R Code

```
1 #####
2 ##### Libraries #####
3 #####
4
5 # clear workspace
6 rm(list = ls())
7
8 #install packages as required
9
10 # Exception: Installing Keras is a little tricky.
11 # You additionally need a python installation on your machine
12 #install_keras(method = "conda")
13
14 # Load libraries
15 library(keras)
16 library(lime)
17 library(tfruns)
18 library(tidyquant)
19 library(rsample)
20 library(recipes)
21 library(yardstick)
22 library(corr)
23 library(randomForest)
24 library(missForest)
25 library(neuralnet)
26 library(caret)
27 library(dplyr)
28 library(corrplot)
29 library(pROC)
30 library(processx)
31 library(caret)
32 library(e1071)
33 library(stargazer)
34
35 # set same seed for R, Python, NumPy and Tensorflow
36 use_session_with_seed(42)
37
38 #####
39 ## Data cleaning and Feature Engineering ##
40 #####
41
42 # read and check data
43 xsell_data_raw <- read.csv("xsell.csv", na.strings=c("", "NA"), stringsAsFactors =
44 FALSE)
45 glimpse(xsell_data_raw)
46 summary(xsell_data_raw)
47
48 # shuffle rows
49 xsell_data_raw <- xsell_data_raw[sample(nrow(xsell_data_raw)),]
50
51 # create new variable tenure
52 xsell_data_raw$tenure <- xsell_data_raw$age - xsell_data_raw$entry_age
53
54 # data cleaning, replace NAs in char-variables with "None_or_missing"
55 xsell_data_raw <- xsell_data_raw %>%
56   replace_na(list(pref_device = "None_or_missing"))
57 xsell_data_raw <- xsell_data_raw %>%
58   replace_na(list(occupation = "None_or_missing"))
59
60 # all character columns to factor:
61 xsell_data_raw <- mutate_if(xsell_data_raw, is.character, as.factor)
62 #additional numeric variables that should rather be treated as factors
63 xsell_data_raw$car_seg <- as.factor(xsell_data_raw$car_seg)
64 xsell_data_raw$acad <- as.factor(xsell_data_raw$acad) # remove if strange results
65 xsell_data_raw$giro_mailing <- as.factor(xsell_data_raw$giro_mailing) # remove if
66   strange results
67 xsell_data_raw$pop_km <- as.factor(xsell_data_raw$pop_km)
```

```

66 xsell_data_raw$ppower <- as.factor(xsell_data_raw$ppower)
67
68 # Remove unnecessary data and clean data set
69 xsell_data_tbl <- xsell_data_raw %>%
70   select(-X) %>% #removes ID
71   # if you don't want to run a random forest for NA imputation, you can do apply of
72   # the two easier fixes to NA's:
73   #drop_na() %>% # removes all NA's. Bad Solution! Improve! Removes 70% of
74   # observations
75   #na.roughfix(xsell_data_raw) #replaces NA's: Numeric with median, factor with mode
76
77 # Impute NAs with a Random Forest
78 xsell_data_tbl$xsell <- as.factor(xsell_data_tbl$xsell) # transform to factor for
79   random forest imputation
80 xsell_data_tbl <- rfImpute(xsell ~ ., xsell_data_tbl, iter = 4, ntree=100)
81 xsell_data_tbl$xsell <- as.integer(ifelse(xsell_data_tbl$xsell == "0", 0, 1)) #
82   transform back to numeric
83
84 glimpse(xsell_data_tbl)
85
86 # Split test/training sets
87 train_test_split <- initial_split(xsell_data_tbl, prop = 0.8)
88 train_test_split
89
90 # Retrieve train and test sets
91 train_tbl <- training(train_test_split)
92 test_tbl <- testing(train_test_split)
93
94 # define features for binning
95 to_bin <- c("age", "entry_age", "last_acc")
96
97 # Create recipe
98 rec_obj <- recipe(xsell ~ ., data = train_tbl) %>%
99   step_discretize(to_bin, options = list(cuts = 4)) %>%
100   step_dummy(all_nominal(), -all_outcomes()) %>%
101   step_center(all_predictors(), -all_outcomes()) %>%
102   step_scale(all_predictors(), -all_outcomes()) %>%
103   prep(data = train_tbl)
104
105 # Apply recipe to predictors (all vars excluding xsell)
106 x_train_tbl <- bake(rec_obj, new_data = train_tbl) %>% select(-xsell)
107 x_test_tbl <- bake(rec_obj, new_data = test_tbl) %>% select(-xsell)
108 glimpse(x_train_tbl)
109
110 # define response variables for training and testing sets
111 y_train_vec <- pull(train_tbl, xsell)
112 y_test_vec <- pull(test_tbl, xsell)
113
114 # visually check transformed data set with a histogram for each feature
115 x_train_tbl %>%
116   gather(colnames, xsell) %>%
117   ggplot(aes(x = xsell)) +
118     geom_histogram() +
119     facet_wrap(~colnames,
120               scales = 'free',
121               ncol = 9)+
122   theme_bw()
123
124 # correlation matrix
125 correl_matrix <- cor(x_train_tbl, use="pairwise.complete.obs")
126 # correlation plot
127 corrplot(correl_matrix)
128
129 #####
130 # Building the Artificial Neural Network #####
131 #####
132
133 # I tuned the hyperparameters by trying out 324 different models
134 # See the scripts hyperpar_tuning_tf_runs.R and keras_nnet_architecture
135 # I then use the "best" hyperparameters below according to the results of those runs

```

```

132
133 # Setting up the ANN with Keras
134 model_keras <- keras_model_sequential()
135
136 # the hyperparameters inserted here were tuned with the hyperpar_tuning_tf_runs.R
    script
137 model_keras %>%
138
139 # First hidden layer
140 layer_dense(
141   units           = 64,
142   kernel_initializer = "uniform",
143   activation       = "relu",
144   use_bias         = TRUE,
145   bias_initializer  = 'zeros',
146   input_shape      = ncol(x_train_tbl)) %>%
147
148 # Dropout to prevent overfitting
149 layer_dropout(rate = 0.6) %>%
150
151 # Second hidden layer
152 layer_dense(
153   units           = 128,
154   kernel_initializer = "uniform",
155   use_bias         = TRUE,
156   bias_initializer  = 'zeros',
157   activation       = "relu") %>%
158
159 # Dropout to prevent overfitting
160 layer_dropout(rate = 0.6) %>%
161
162 # Output layer
163 layer_dense(
164   units           = 1,
165   kernel_initializer = "uniform",
166   activation       = "sigmoid") %>%
167
168 # Compile ANN
169 compile(
170   optimizer = 'adam',
171   loss       = 'binary_crossentropy',
172   metrics    = c('accuracy')
173 )
174
175 # print model architecture
176 model_keras
177
178 # Tensorboard can be useful to follow along the training of complex models as it
    trains
179 # Also, you can easily compare different models in one graph
180 # in this case, this is not necessary, hence commented
181 # launch TensorBoard (data won't show up until after the first epoch)
182
183 # tensorboard("logs/run_1")
184
185 history <- fit(
186   object      = model_keras,
187   x           = as.matrix(x_train_tbl),
188   y           = y_train_vec,
189   batch_size  = 100,
190   epochs      = 10,
191   validation_split = 0.3,
192   # include callback below if you want to use tensorflow
193   #callbacks = callback_tensorboard("logs/run_1", write_images = TRUE),
194   verbose     = 1
195 )
196
197 ## compare runs
198 ## In this example, I used three different callbacks for three different model
    specifications

```



```

199 # tensorboard("logs")
200
201 # Plot the training/validation history of the Keras model
202 plot(history) +
203   theme_bw()
204
205 #####
206 ## Extract Predictions from the Keras ANN ###
207 #####
208
209 # Predicted Class
210 yhat_keras_class_vec <- predict_classes(object = model_keras, x = as.matrix(x_test_tbl
211   )) %>%
212   as.vector()
213
214 # Predicted Class Probability
215 yhat_keras_prob_vec <- predict_proba(object = model_keras, x = as.matrix(x_test_tbl))
216   %>%
217   as.vector()
218
219 #####
220 ### Model Comparison ANN/Logit #####
221 #####
222 # run script that sets up a simple benchmark Logit model
223 source("benchmark_logit_model.R", echo = TRUE)
224
225 # compare ANN Model ROC curve to benchmark logit model
226 roc_comp <- roc(estimates_keras_tbl$truth, estimates_keras_tbl$class_prob, percent=TRUE
227   , plot=TRUE, print.auc=TRUE,
228     print.auc.x = 90, print.auc.y = 85, grid=TRUE)
229 roc_comp <- roc(xsell_valid$xsell, xsell_valid$pred_logit, percent=TRUE, plot=TRUE,
230   print.auc=TRUE, grid=TRUE, col = "red", print.auc.x = 90, print.auc.y
231     = 75, add = TRUE)
232
233 text(40,50, "Simple Benchmark Logit Model")
234 text(33,40, "Hyperparameter-tuned ANN with Keras", col = "red")
235
236 # compare model statistics (Confusion Matrix, Accuracy, Sensitivity, Recall, etc.)
237 conf_matrix_keras <- confusionMatrix(as.factor(yhat_keras_class_vec), as.factor(
238   estimates_keras_tbl$truth),
239   positive="1", dnn = c("Prediction", "Actual"))
240 conf_matrix_logit <- confusionMatrix(as.factor(xsell_valid$predict), as.factor(xsell_
241   valid$xsell),
242   positive="1", dnn = c("Prediction", "Actual"))
243
244 conf_matrix_keras
245 conf_matrix_logit
246
247 ConfMatK <- as.data.frame.matrix(conf_matrix_keras$table)
248 stargazer(ConfMatK, head = FALSE, title = "Table", summary = FALSE)
249
250 #####
251 ### Evaluate Feature Importance with LIME #####
252 #####
253
254 # Setup
255 class(model_keras)
256
257 #Setup lime::model_type() function for keras
258 #This specifies that the task at hand is a classification task
259 model_type.keras.engine.sequential.Sequential <- function(x, ...) {
260   return("classification")
261 }
262
263 # Setup lime::predict_model() function for keras
264 predict_model.keras.engine.sequential.Sequential <- function(x, newdata, type, ...) {
265   pred <- predict_proba(object = x, x = as.matrix(newdata))
266   return(data.frame(Yes = pred, No = 1 - pred))
267 }

```

```

263 }
264
265
266 # Test the predict_model() function
267 predict_model(x = model_keras, newdata = x_test_tbl, type = 'raw') %>%
268   tibble::as_tibble()
269
270 # Run lime() on training set
271 explainer <- lime::lime(
272   x = x_train_tbl,
273   model = model_keras,
274   bin_continuous = FALSE
275 )
276
277 # Run explain() on explainer
278 explanation <- lime::explain(
279   x_test_tbl[1:4, ],
280   explainer = explainer,
281   n_labels = 1,
282   n_features = 10,
283   kernel_width = 0.5,
284   feature_select = "forward_selection"
285 )
286
287 # Plot feature importance
288 plot_features(explanation, ncol = 2) #+
289   labs(title = "LIME Feature Importance Visualization",
290        caption = "Test data set, first four customers")
291
292 plot_explanations(explanation) +
293   labs(title = "LIME Feature Importance Heatmap",
294        subtitle = "Hold Out (Test) Set, First Four Cases Shown") +
295   theme_bw()
296
297 # Feature correlations to xsell
298 corrr_analysis <- x_train_tbl %>%
299   mutate(xsell = y_train_vec) %>%
300   correlate() %>%
301   focus(xsell) %>%
302   rename(feature = rowname) %>%
303   arrange(abs(xsell)) %>%
304   mutate(feature = as_factor(feature))
305
306 # Correlation visualization
307 corrr_analysis %>%
308   ggplot(aes(x = xsell, y = fct_reorder(feature, desc(xsell)))) +
309   geom_point() +
310   # Positive Correlations - Contribute to xsell
311   geom_segment(aes(xend = 0, yend = feature),
312               color = palette_light()[[2]],
313               data = corrr_analysis %>% filter(xsell > 0)) +
314   geom_point(color = palette_light()[[2]],
315             data = corrr_analysis %>% filter(xsell > 0)) +
316   # Negative Correlations - Prevent xsell
317   geom_segment(aes(xend = 0, yend = feature),
318               color = palette_light()[[1]],
319               data = corrr_analysis %>% filter(xsell < 0)) +
320   geom_point(color = palette_light()[[1]],
321             data = corrr_analysis %>% filter(xsell < 0)) +
322   # Vertical lines
323   # geom_vline(xintercept = 0, color = palette_light()[[5]], size = 1, linetype = 2) +
324   # geom_vline(xintercept = -0.05, color = palette_light()[[5]], size = 1, linetype =
325     2) +
326   # geom_vline(xintercept = 0.05, color = palette_light()[[5]], size = 1, linetype =
327     2) +
328   # Aesthetics
329   theme_bw() + # theme_tq replaced because in conflict with randomForest package
330   labs(title = "Cross Sell Correlation Analysis",
331        subtitle = paste("Negative correlations (prevent xsell)",
332                          "positive correlations (contribute to xsell)"),

```

```

331     x = "Correlation of feature with xsell",
332     y = "Feature Importance")
333
334 # individual variable assessment: age
335 train_tbl %>%
336   ggplot(aes(x=as.factor(xsell), y=age))+
337   geom_jitter(shape = 1, alpha = 0.2) +
338   geom_violin(fill="blue", alpha = 0.2) +
339   labs(x = "xsell")+
340   theme_bw()
341
342 # individual variable assessment: giro mailing
343 table(as.factor(train_tbl$xsell), train_tbl$giro_mailing)
344
345 # individual variable assessment: logins
346 train_tbl %>%
347   ggplot(aes(x=as.factor(xsell), y=logins))+
348   geom_jitter(shape = 1, alpha = 0.5) +
349   geom_violin(fill="blue", alpha = 0.2) +
350   ylim(0,10) +
351   theme_bw()
352
353 table(train_tbl$xsell, train_tbl$gender)

```

Listing 1: Main analysis

```

1 #####
2 ### Hyperparameter tuning #####
3 #####
4
5 # run the external script with defined flags as default
6 training_run("keras_nnet_architecture.R")
7
8 # !!CAUTION!! Combines every single combination (3^(#tuned parameters)), thus long
   runtime
9 # run various combinations of dropout1 and dropout2
10
11 # Only Dropout tuning
12 runs <- tuning_run("keras_nnet_architecture.R", runs_dir = "dropout_tuning", flags =
   list(
13   dropout1 = c(0.2, 0.4, 0.6),
14   dropout2 = c(0.2, 0.4, 0.6)
15 ))
16
17 # run combinations of dense units, epochs and batch size
18 runs <- tuning_run("keras_nnet_architecture.R", runs_dir = "runs/20190608_neurons_
   epochs_batch_tuning", flags = list(
19   #dropout1 = c(0.2, 0.4, 0.6),
20   #dropout2 = c(0.2, 0.4, 0.6),
21   dense_units1 = c(8, 64, 128),
22   dense_units2 = c(8, 64, 128),
23   epochs = c(5, 10, 30),
24   batch_size = c(50, 100, 150)
25 ))
26
27 # only number of neurons per layer tuning
28 runs <- tuning_run("keras_nnet_architecture.R", runs_dir = "number_neurons_tuning",
   flags = list(
29   dense_units1 = c(8, 32, 64),
30   dense_units2 = c(8, 32, 64)
31 ))
32
33
34 #####
35 ### After-tunig processing #####
36 #####
37
38 # check latest run
39 latest_run()

```

```

40
41 # compare runs in interactive RStudio viewer
42 compare_runs()
43
44 # View the run with the highest val_acc
45 view_run("runs/20190608_neurons_epochs_batch_tuning/2019-06-08T08-56-30Z")
46
47 # compare runs in data frame
48 runs1 <- ls_runs(runs_dir = "runs/20190607_dropout_neurons_epochs_tuning", order =
  metric_val_acc)
49 runs2 <- ls_runs(runs_dir = "runs/20190608_neurons_epochs_batch_tuning", order =
  metric_val_acc)
50
51 # combine both runs in one df
52 runs <- rbind(runs1, runs2)
53
54 # extract important information only for printing with stargazer
55 runs_print <- runs %>%
56   arrange(desc(metric_val_acc)) %>%
57   select(metric_val_acc, metric_acc, flag_dense_units1, flag_dropout1,
58     flag_dense_units2, flag_dropout2, flag_epochs, flag_batch_size)
59
60 stargazer(runs_print[1:5,], summary = FALSE)
61 stargazer(runs_print[323:324,], summary = FALSE)
62
63 # move all runs into run/archive
64 clean_runs()
65
66 # extract and save certain model
67 copy_run_files("runs/2019-06-08T00-38-43Z", to = "20190607-best-model")

```

Listing 2: Hyperparameter tuning routine

```

1 #####
2 ### keras neural network architecture ###
3 #####
4
5 # to be called by hyperpar_tuning_tf_runs.R
6
7 # setting up flags for hyperparameter tuning
8 FLAGS <- flags(
9   flag_integer("dense_units1", 8),
10  flag_numeric("dropout1", 0.6),
11  flag_integer("dense_units2", 64),
12  flag_numeric("dropout2", 0.6),
13  flag_integer("epochs", 30),
14  flag_integer("batch_size", 128)
15 )
16
17
18 # Setting up the ANN with Keras
19 model_keras <- keras_model_sequential()
20
21 model_keras %>%
22
23   # First hidden layer
24   layer_dense(
25     units           = FLAGS$dense_units1,
26     kernel_initializer = "uniform",
27     activation       = "relu",
28     use_bias         = TRUE,
29     bias_initializer  = 'zeros',
30     input_shape      = ncol(x_train_tbl)) %>%
31
32   # Dropout to prevent overfitting
33   layer_dropout(rate = FLAGS$dropout1) %>%
34
35   # Second hidden layer
36   layer_dense(

```

```

37     units                = FLAGS$dense_units2,
38     kernel_initializer   = "uniform",
39     use_bias              = TRUE,
40     bias_initializer      = 'zeros',
41     activation            = "relu") %>%
42
43   # Dropout to prevent overfitting
44   layer_dropout(rate = FLAGS$dropout2) %>%
45
46   # Output layer
47   layer_dense(
48     units                = 1,
49     kernel_initializer   = "uniform",
50     activation            = "sigmoid") %>%
51
52   # Compile ANN
53   compile(
54     optimizer = 'adam',
55     loss      = 'binary_crossentropy',
56     metrics   = c('accuracy')
57   )
58
59
60 history <- fit(
61   object      = model_keras,
62   x            = as.matrix(x_train_tbl),
63   y            = y_train_vec,
64   batch_size  = FLAGS$batch_size,
65   epochs      = FLAGS$epochs,
66   validation_split = 0.3,
67   verbose     = 0
68 )

```

Listing 3: Flagged neural net architecture for hyperparameter tuning

```

1 #####
2 ###   Apply model to out-of-time sample   ###
3 #####
4
5 # import oot data set
6 xsell_oot <- read.csv("xsell_oot.csv", na.strings=c("", "NA"), stringsAsFactors = FALSE)
7
8 # create new variable tenure
9 xsell_oot$tenure <- xsell_oot$age - xsell_oot$entry_age
10
11 # data cleaning, replace NAs in char-variables with "None_or_missing"
12 xsell_oot <- xsell_oot %>%
13   replace_na(list(pref_device = "None_or_missing"))
14 xsell_oot <- xsell_oot %>%
15   replace_na(list(occupation = "None_or_missing"))
16
17 # all character columns to factor:
18 xsell_oot <- mutate_if(xsell_oot, is.character, as.factor)
19 #additional numeric variables that should rather be treated as factors
20 xsell_oot$car_seg <- as.factor(xsell_oot$car_seg)
21 xsell_oot$acad <- as.factor(xsell_oot$acad) # remove if strange results
22 xsell_oot$giro_mailing <- as.factor(xsell_oot$giro_mailing) # remove if strange
   results
23 xsell_oot$pop_km <- as.factor(xsell_oot$pop_km)
24 xsell_oot$ppower <- as.factor(xsell_oot$ppower)
25
26 # extract ID (X)
27 X <- xsell_oot$X
28
29 # Remove unnecessary data and clean data set
30 xsell_oot_data_tbl <- xsell_oot %>%
31   select(-X) %>% #removes ID

```

```

32 # if you don't want to run a random forest for NA imputation, you can do apply of
    the two easier fixes to NA's:
33 #drop_na() #>% # removes all NA's. Bad Solution! Improve! Removes 70% of
    observations
34 na.roughfix() #replaces NA's: Numeric with median, factor with mode
35
36 # Impute NAs with a Random Forest
37 # can't do here, because there is no xsell
38 #xsell_oout_data_tbl <- rfImpute(xsell ~ . ,xsell_oout_data_tbl, iter = 4, ntree=100)
39
40 # apply recipe to normalize, etc.
41 x_valid_tbl <- bake(rec_obj, new_data = xsell_oout_data_tbl)
42
43 # Predicted Class
44 yhat_keras_class_vec_oout <- predict_classes(object = model_keras, x = as.matrix(x_
    valid_tbl)) %>%
45   as.vector()
46
47 xsell_oout_incl_pred <- cbind(X, xsell_pred = yhat_keras_class_vec_oout)
48 write.csv(xsell_oout_incl_pred, "xsell_oout_predicted.csv", row.names = FALSE)

```

Listing 4: Using the ANN to predict the previously unseen out-of-time data set

References

- Abadi, M., Agarwal, A., Barham, P., Brevdo, E., Chen, Z., Citro, C., Corrado, G. S., Davis, A., Dean, J., Devin, M., Ghemawat, S., Goodfellow, I., Harp, A., Irving, G., Isard, M., Jia, Y., Jozefowicz, R., Kaiser, L., Kudlur, M., Levenberg, J., Mane, D., Monga, R., Moore, S., Murray, D., Olah, C., Schuster, M., Shlens, J., Steiner, B., Sutskever, I., Talwar, K., Tucker, P., Vanhoucke, V., Vasudevan, V., Viegas, F., Vinyals, O., Warden, P., Wattenberg, M., Wicke, M., Yu, Y. & Zheng, X. (2015), ‘TensorFlow: Large-Scale Machine Learning on Heterogeneous Distributed Systems’.
- Bengio, Y. (2000), ‘Gradient-based optimization of hyperparameters’, *Neural computation* **12**(8), 1889–1900.
- Benítez, J. M., Castro, J. L. & Requena, I. (1997), ‘Are artificial neural networks black boxes?’, *IEEE Transactions on neural networks* **8**(5), 1156–1164.
- Bergstra, J. & Bengio, Y. (2012), ‘Random search for hyper-parameter optimization’, *Journal of Machine Learning Research* **13**(Feb), 281–305.
- Chan, H.-P., Sahiner, B., Wagner, R. F. & Petrick, N. (1999), ‘Classifier design for computer-aided diagnosis: Effects of finite sample size on the mean performance of classical and neural network classifiers’, *Medical physics* **26**(12), 2654–2668.
- Chen, J. H. & Asch, S. M. (2017), ‘Machine Learning and Prediction in Medicine — Beyond the Peak of Inflated Expectations’, *The New England journal of medicine* **376**(26), 2507–2509.
- Chollet, F. & Allaire, J. (2017), ‘R Interface to Keras’.
- Dayhoff, J. E. & DeLeo, J. M. (2001), ‘Artificial neural networks: Opening the black box’, *Cancer: Interdisciplinary International Journal of the American Cancer Society* **91**(S8), 1615–1635.
- Domingos, P. M. (1999), MetaCost: A General Method for Making Classifiers Cost-Sensitive, in ‘Fifth ACM SIGKDD International Conference on Knowledge Discovery and Data Mining’, San Diego, CA, pp. 155–164.
- Esteva, A., Kuprel, B., Novoa, R. A., Ko, J., Swetter, S. M., Blau, H. M. & Thrun, S.

- (2017), ‘Dermatologist-level classification of skin cancer with deep neural networks’, *Nature* **542**(7639), 115–118.
- Fawcett, T. (2006), ‘An introduction to ROC analysis’, *Pattern Recognition Letters* **27**(8), 861–874.
- Felvey, J. (1982), ‘Cross-selling by computer’, *Bank Marketing* **7**, 25–27.
- Glorot, X., Bordes, A. & Bengio, Y. (2011), ‘Deep Sparse Rectifier Neural Networks’, *Proceedings of the 14th International Conference on Artificial Intelligence and Statistics* **15**, 9.
- Hansen, L. K. & Salamon, P. (1990), ‘Neural network ensembles’, *IEEE Transactions on Pattern Analysis & Machine Intelligence* (10), 993–1001.
- 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.
- Hinton, G., Deng, L., Yu, D., Dahl, G., Mohamed, A.-r., Jaitly, N., Senior, A., Vanhoucke, V., Nguyen, P., Kingsbury, B. & Sainath, T. (2012), ‘Deep Neural Networks for Acoustic Modeling in Speech Recognition’, *IEEE Signal Processing Magazine* **29**.
- Hornik, K. (1991), ‘Approximation capabilities of multilayer feedforward networks’, *Neural Networks* **4**(2), 251–257.
- Jaderberg, M., Simonyan, K., Vedaldi, A. & Zisserman, A. (2014), Synthetic Data and Artificial Neural Networks for Natural Scene Text Recognition, *in* ‘2011 International Conference On Computer Vision’.
- Jordan, M. I. & Mitchell, T. M. (2015), ‘Machine learning: Trends, perspectives, and prospects’, *Science* **349**(6245), 255–260.
- Kamakura, W. A., Ramaswami, S. N. & Srivastava, R. K. (1991), ‘Applying latent trait analysis in the evaluation of prospects for cross-selling of financial services’, *International Journal of Research in Marketing* **8**(4), 329–349.
- Kamakura, W. A., Wedel, M., de Rosa, F. & Mazzon, J. A. (2003), ‘Cross-selling

- through database marketing: A mixed data factor analyzer for data augmentation and prediction’, *International Journal of Research in Marketing* **20**(1), 45–65.
- Kavicky, R. (2017), ‘Local Interpretable Model-Agnostic Explanations (R port of original Python package)’.
- Keaveney, S. M. (1995), ‘Customer Switching Behavior in Service Industries: An Exploratory Study’, *Journal of Marketing* **59**(2), 71–82.
- Kingma, D. P. & Ba, J. (2014), Adam: A Method for Stochastic Optimization, in ‘International Conference on Learning Representations 2015’, San Diego, CA.
- Li, S., Sun, B. & L. MONTGOMERY, A. (2011), ‘Cross-Selling the Right Product to the Right Customer at the Right Time’, *Journal of Marketing Research* **48**, 683–700.
- Liaw, A. & Wiener, M. (2002), ‘Classification and Regression by randomForest’, *R News* **2**(3), 18–22.
- Olden, J. D. & Jackson, D. A. (2002), ‘Illuminating the “black box”: A randomization approach for understanding variable contributions in artificial neural networks’, *Ecological modelling* **154**(1-2), 135–150.
- Pedersen, T. L. (2018), ‘Package ‘lime’: Local Interpretable Model-Agnostic Explanations’.
- Pedersen, T. L. & Benesty, M. (2018), ‘Understanding lime’, https://cran.r-project.org/web/packages/lime/vignettes/Understanding_lime.html.
- Reichheld, F. F. & Sasser Jr, W. E. (1990), ‘Zero defections: Quality comes to services’, *Harvard Business Review* **68**(5), 105–111.
- Ribeiro, M. T., Singh, S. & Guestrin, C. (2016), "Why Should I Trust You?": Explaining the Predictions of Any Classifier, in ‘22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining’, San Francisco, CA, USA, pp. 1135–1144.
- Ruder, S. (2016), ‘An overview of gradient descent optimization algorithms’, <http://sebastianruder.com/optimizing-gradient-descent/index.html>.

- Srivastava, N., Hinton, G., Krizhevsky, A., Sutskever, I. & Salakhutdinov, R. (2014), ‘Dropout: A Simple Way to Prevent Neural Networks from Overfitting’, *Journal of Machine Learning Research* **15**, 1929–1958.
- Werbos, P. J. (1990), ‘Backpropagation through time: What it does and how to do it’, *Proceedings of the IEEE* **78**(10), 1550–1560.
- Zhang, Z. & Sabuncu, M. (2018), ‘Generalized Cross Entropy Loss for Training Deep Neural Networks with Noisy Labels’, *Advances in Neural Information Processing Systems 31* pp. 8778–8788.

Statutory Declaration

Ich versichere hiermit, dass ich die vorliegende Arbeit selbständig und ohne Benutzung anderer als der angegebenen Quellen und Hilfsmittel verfasst habe. Wörtlich übernommene Sätze oder Satzteile sind als Zitat belegt, andere Anlehnungen, hinsichtlich Aussage und Umfang, unter Quellenangabe kenntlich gemacht. Die Arbeit hat in gleicher oder ähnlicher Form noch keiner Prüfungsbehörde vorgelegen und ist nicht veröffentlicht. Sie wurde nicht, auch nicht auszugsweise, für eine andere Prüfungs- oder Studienleistung verwendet.

I herewith declare that I have completed the present term paper 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 15, 2019



Lukas Jürgensmeier