

Context / Keyword (roof)

IDA: Project 6 Caravan Insurance – Bridging the Gap between Machine Learning and Business Strategy

Keyword (backbone)

Keyword (backbone)

Chenpo HU

Reviewer: Prof. Dr. Tobias Scheffer

Tutor: Dr. Lena Jäger

Dr. Paul Prasse

Silvia Makowski

Potsdam, September 2018

Agenda



Block 1: Introduction

09:30 – 09:40

1 Tasks and datasets

2 Traps that I run into

3 Trained machine learning process



Block 2: Implementation and Result

09:40 – 10:10

4 Introduction of applied machine learning methods

5 Implementation and result in Jupyter

This is the final Version of the project, in which the trained machine learning process is applied.

6 Recommendation of good performed models



Block 3: Discussion

10:10 – 10:20

7 In-depth discussion

In-depth discussion focusing on key use-cases and customer journeys leveraging Blockchain based solutions as well as alternatives

Aggregation of discussion results and outline of possible courses of action, including the evaluation based on the fit to business vision and analysis of strategic implications, risks and chances



Block 4: Q & A

10:20 – 10:30

8 Exchange of thoughts

Do you have further improvement advices or innovative ideas?

8 Literature

Literatures cited in the PPT.

1 Tasks and datasets

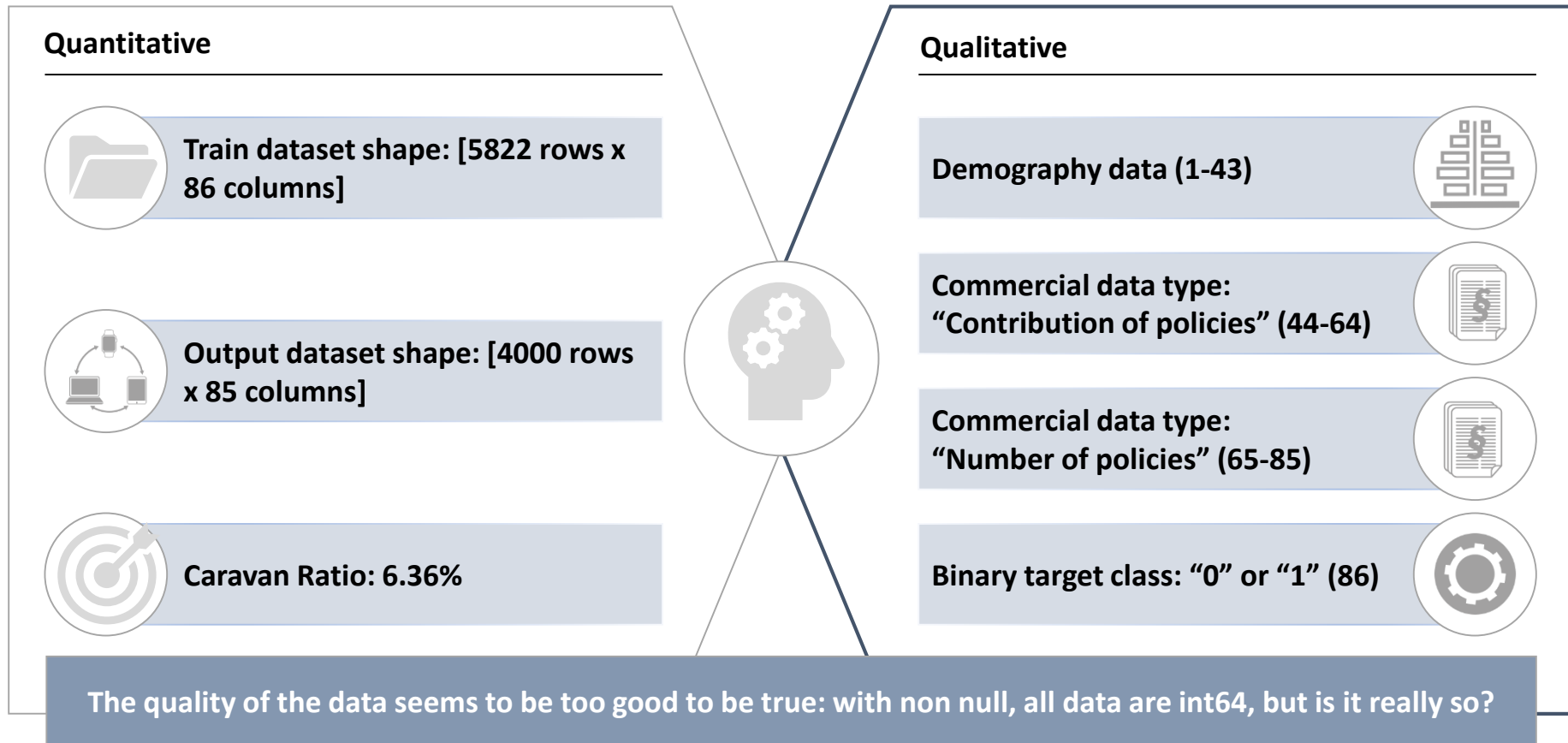


Task

Develop **Models** to determination of **target customers** of for Caravan insurance, translate into data understanding: predict as many **class „1“** correctly as possible



Dataset

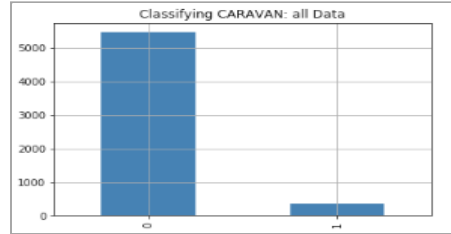


2 Traps that I run into

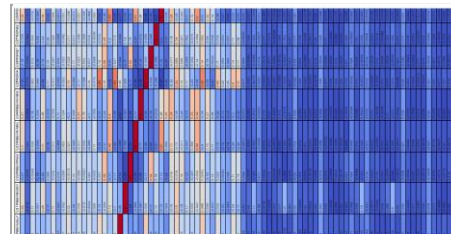


Detours

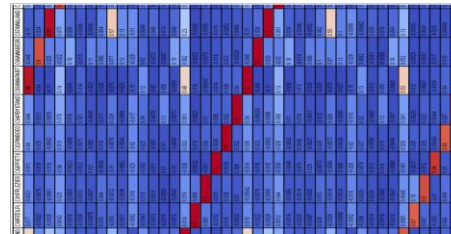
Unbalanced data



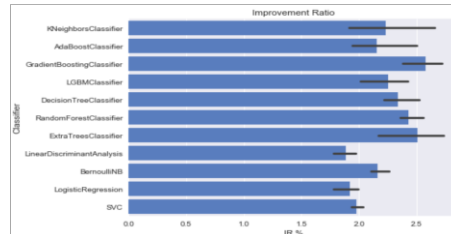
Into Interval coded sensitive nominal data



Highly correlated features



Useless standard scores



Description

It is clear that our dataset is highly unbalanced with only **6.36% of observations** actually buying the insurance.

The **nominal data** of dataset 1 – 42 are coded into **interval** in the raw data. Some datasets in 1 - 42 should be „-1“ correlated and binary, but they are coded into scala.

Feature 44-64 and 65-85 are **strongly (at least 0,85) correlated**, since there are three red line run through the correlation matrix.

The **high accuracy score tells us nothing**, since we can say that every proband will not buy the caravan, which has the accuracy of 93.64%.

Solution

- Up-sample the min
- Down-sample the majority class
- Change your performance metric
- Penalize algorithms (cost-sensitive)
- Use tree-based algorithms

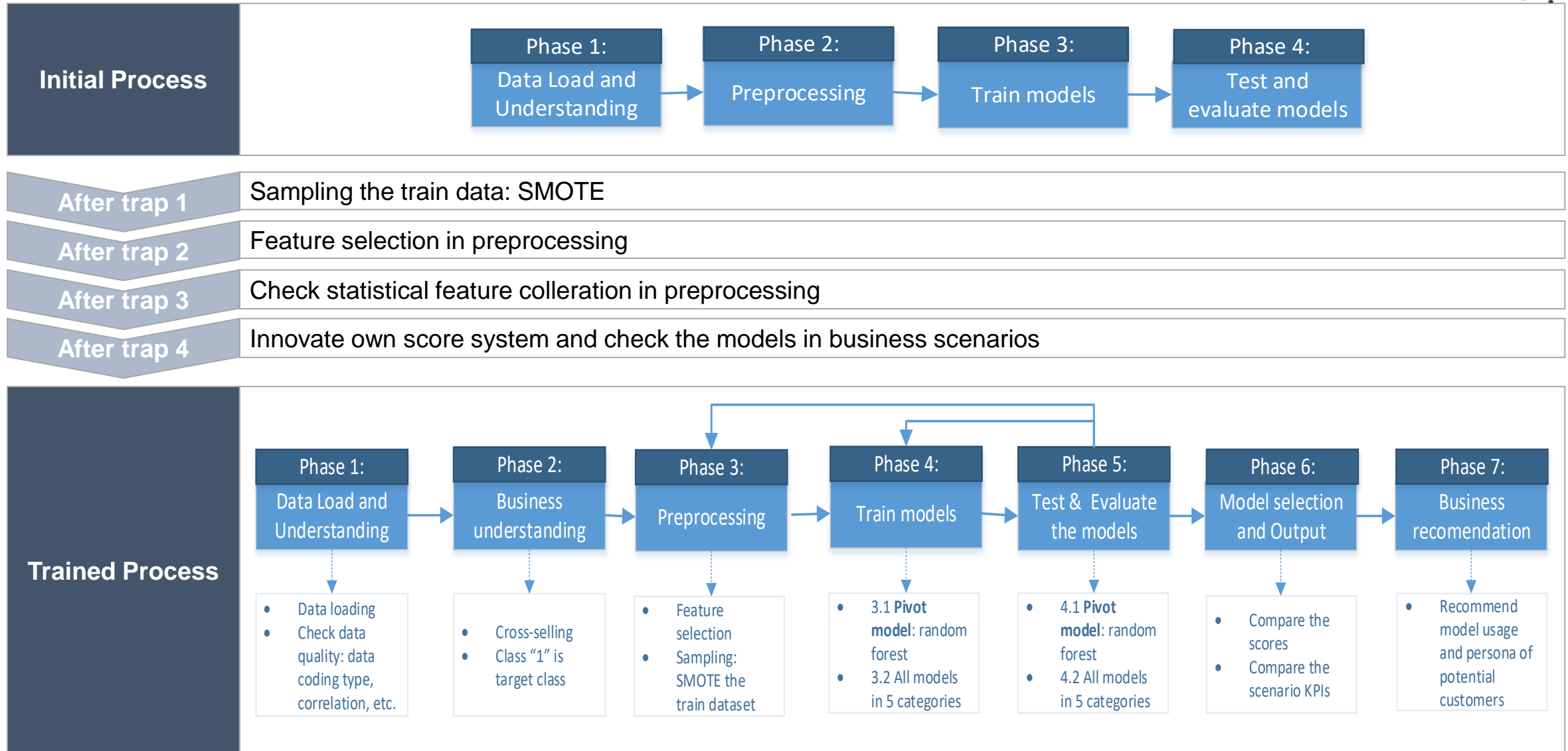
Since those the **Demography data** are encrypted, we focus on the commercial data of 43-85.

I used three kinds of features:

- Feature group 1: 43-85
- Feature group 2: 43-64
- Feature group 3: Top 8 most important features

I use mainly the no. of benefit Items (TP) and the improvement ratio based on confusion matrix to evaluate the performance of the model.

3 Trained machine learning process

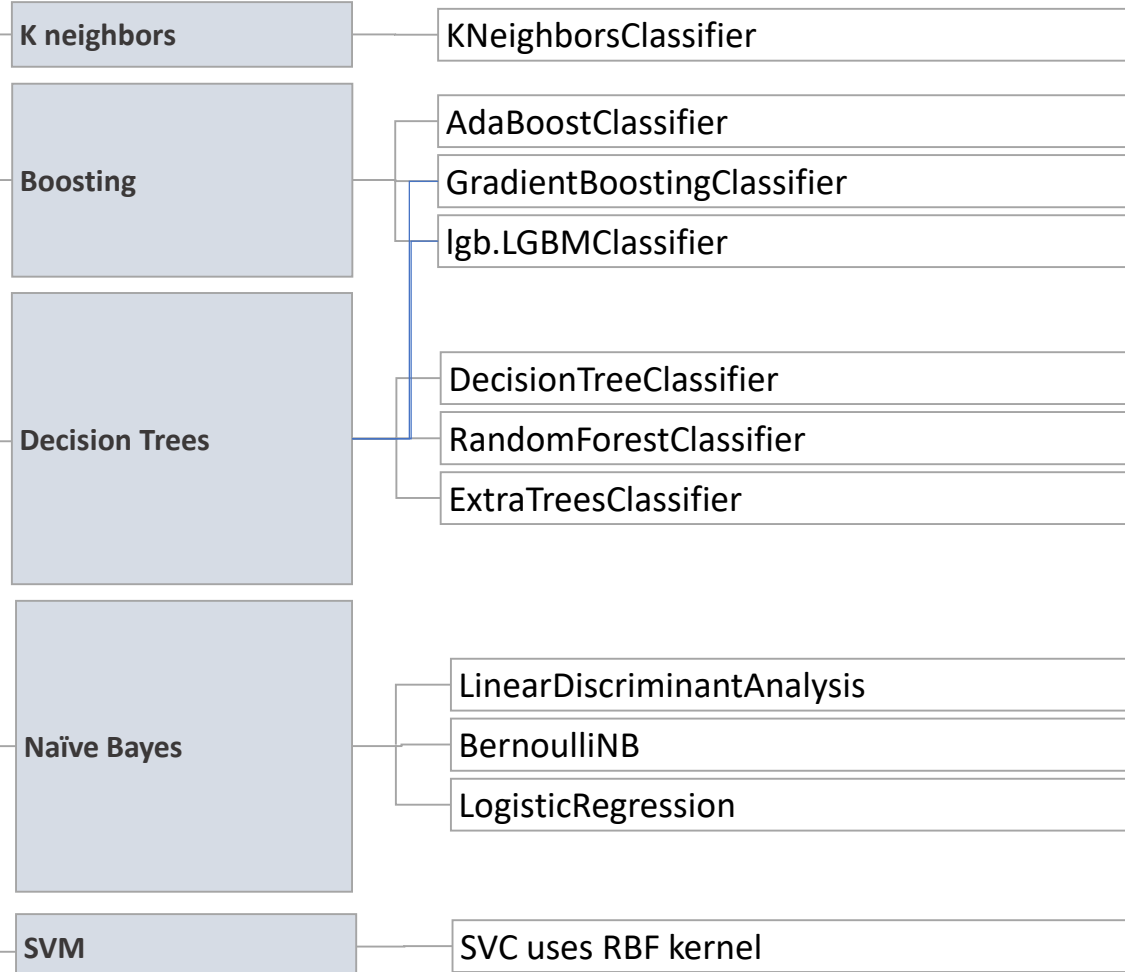


4 Introduction of applied machine learning methods

Overview of categorized classifiers in this project

5 Categories of classifiers

To select the best performed models and learn during the project more technologies, I decide to generate several categories of models instead of generating only 2-3 best performed models



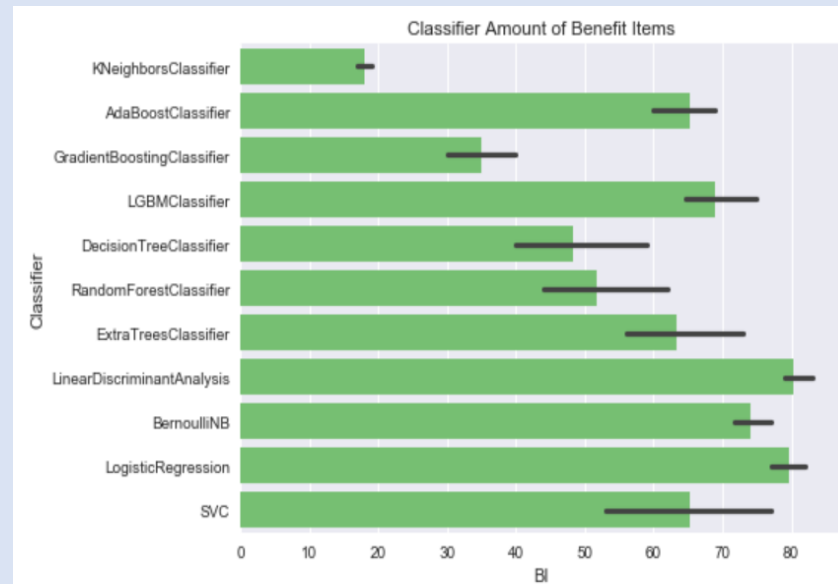
Summary

- The used classifiers can be classified into **5 categories**:
 - **K neighbors** can be used
 - **Boosting** has a family of machine learning algorithms that convert weak learners to strong ones.
 - **Decision trees** often perform well on imbalanced datasets because their hierarchical structure allows them to learn signals from both classes.
 - **Naïve Bayes** is a frequently used in imbalanced dataset problems.
 - In handling a binary classification task, **Support Vector Machine (SVM)** is one of the methods reported to give a high accuracy in predictive modeling compared to the other techniques such as Discriminant Analysis (García et al., 2015).

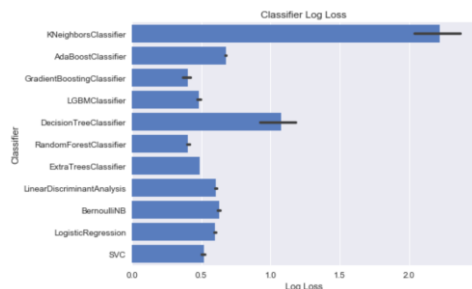
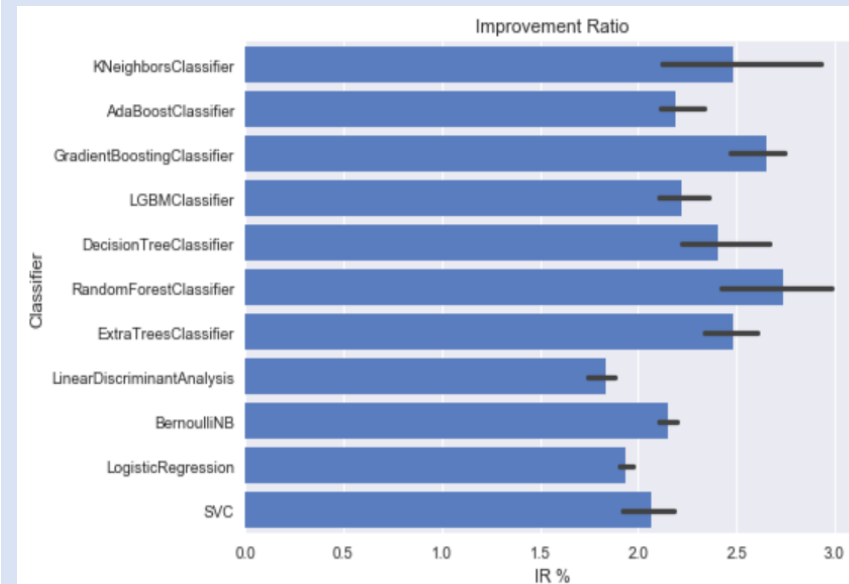
6 Recommendation of good performed models

Select models

Amount of Benefit Items (correctly predicted Caravan)



Improvement Ratio



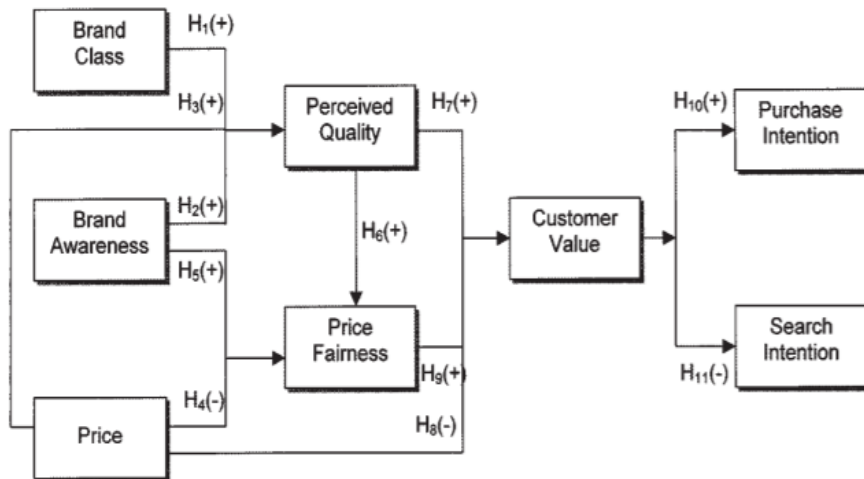
- **Improvement Ratio** = predicted true caravan rate / true caravan rate in raw data = $(TP / (TP + FP)) / ((TP + FN) / (TP + TN + FP + FN))$
- Pros: By using the models, companies can focus more on the customers, who may like to buy caravan with less marketing cost. Companies can select the models based on their company KPIs and strategies, as realized in the scenario in the Jupyter notebook both in case to keep balance or reach the goal.
- Cons: the prediction decreased the data size, 38% - 40% (Naïve Bayes) caravan customers may not be targeted in the marketing.

Conclusion

- **Naïve Bayes** can predict the largest amount of benefit items (**60% - 72%**) correctly, with the improvement ratio between 1.74 to 2.20. More importantly, these models have almost no overfitting problem, which may occur in other models.
- Besides the **Naïve Bayes**, the performance of AdaBoosting, LGBMClassifier, RandomForest (but has slightly overfitting), ExtraTrees and SVC (but has slightly overfitting) are also good.
- All the models have an improvement Ratio of at least 1.74.
- The performance of RandomForestClassifier, ExtraTreesClassifier, LogisticRegression and SVC in the models clusters using top 8 features are good and with slightly or no overfitting problem.
- The **persona of the caravan customers** are: (1) **Contribution to Car policy** with 6 (2) **high purchase power** (3) **high Contribution private third party insurance** (4) **has contribution to boat policies**.

7 In - depth Discussion

Figure 1
A Conceptual Model of Brand and Price Effects on Perceived Quality, Price Fairness, Perceived Value, and Purchase and Search Intention



Source: Oh, H. (2000)

- 1 Maybe the false positive will become true positive in the future.

Are the false positive (customers that are predicted to be caravan, but in in turth not caravam) really false positive?

Maybe they are the potential customers that are not digged out jet? Since this raw data is only valid in the date that it was extracted. What will happen after a few years?

- 2 The customers that are frequently predicted to be „1“ by different models, are most likely to be the caravan customers.

We can analyse the output excel file further and identify these customers, that are predicted to be „1“ frequently by different models.

Since then the possibility that the customer is not caravan would be the multiplication of (1-Benefitratio) of those models, which will be relativ small.

Think outside of the Box!



8 Q & A: Exchange of thoughts



Source: <https://blitzmetrics.com/tag/big-data/>

9 Literatures



Oh, H. (2000). The effect of brand class, brand awareness, and price on customer value and behavioral intentions. *Journal of Hospitality & Tourism Research*, 24(2), 136-162.

García, S., Luengo, J., & Herrera, F. (2015). *Data preprocessing in data mining* (pp. 195-243). Switzerland: Springer International Publishing.