

# Churn Prediction

Beat the logistic regression model of 'Die ZEIT'

Author: Sebastian Fuhrer

# TABLE OF CONTENTS

01

THE COMPILATION  
OF THE DATA SET

02

EXPLORATORY  
DATA ANALYSIS

03

FEATURE  
IMPORTANCE

04

MACHINE LEARNING  
MODELS

05

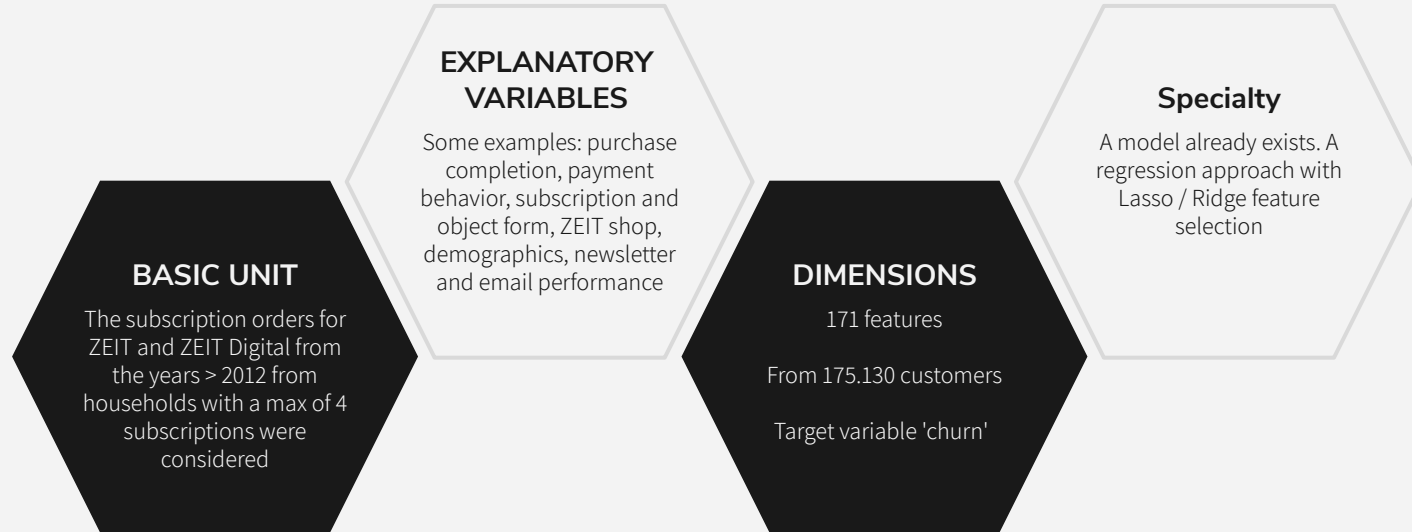
FUTURE WORK



**01**

THE COMPILATION OF THE  
DATA SET

# THE COMPILATION OF THE DATA SET







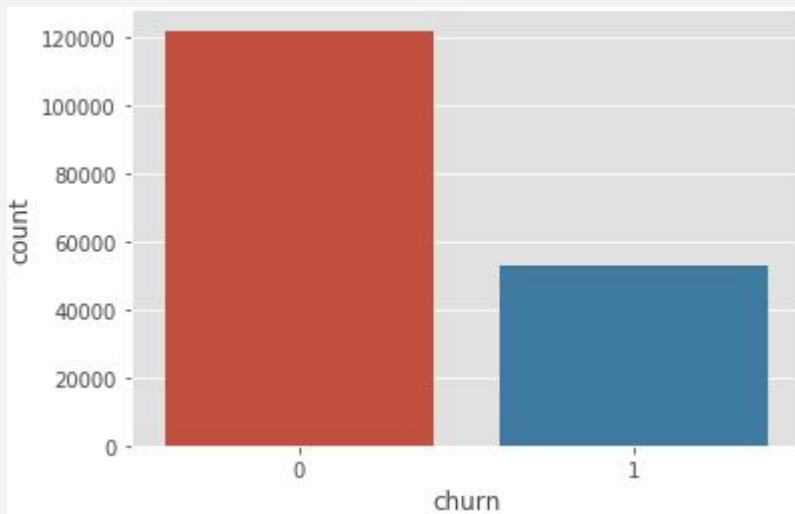
02

## EXPLORATORY DATA ANALYSIS



## CUSTOMER BEHAVIOR

The distribution of the target variable 'churn' is unbalanced. About 52.905 customers churned and 122.140 renew or still have a subscription from 'Die ZEIT' in the period considered.



**CHURN**

## CHURN AND RETENTION RATE

**69,78%**

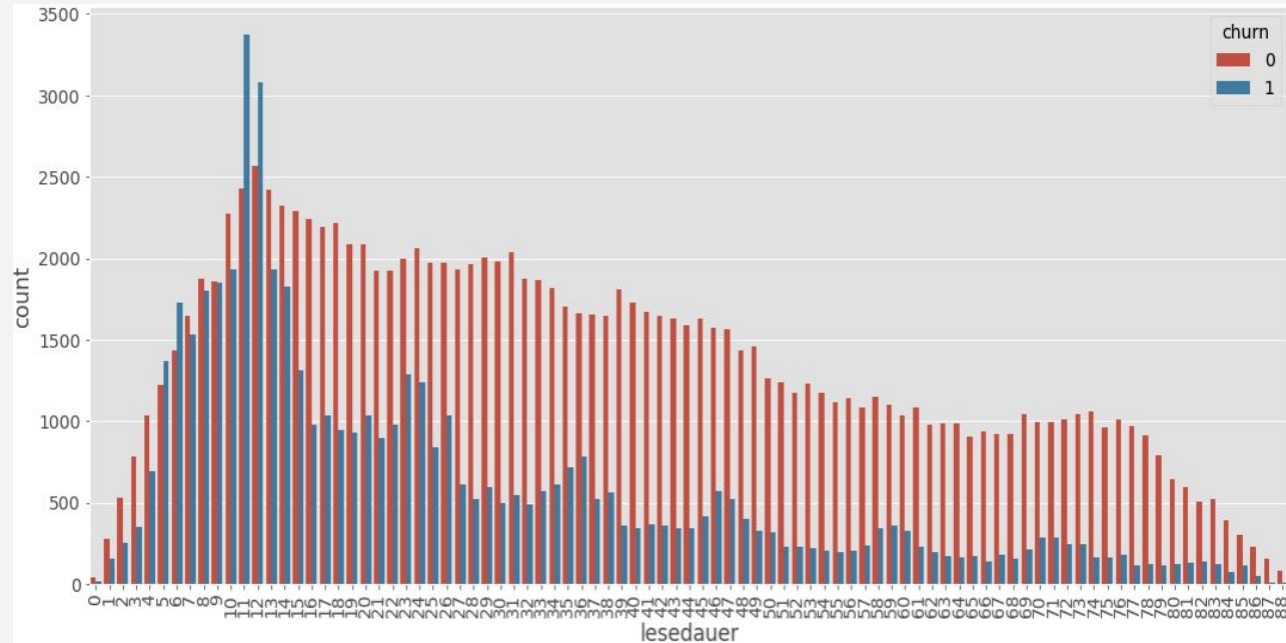
RETENTION

**30,22%**

CHURN

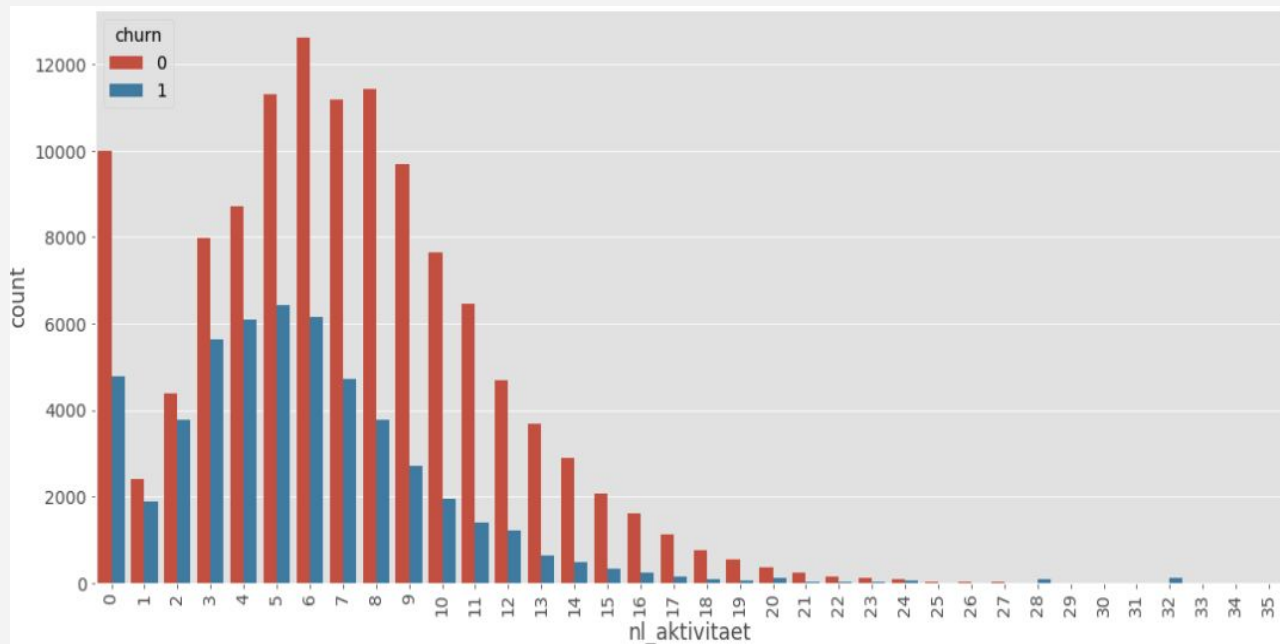
# READING TIME

- The churn rate decreases with increasing reading time
- The highest churn rate is after 11 or 12 months
- 12 months of reading seems to be a threshold towards loyal customers
- There are always smaller peaks in the churn rate, after two years, three years, etc.



## NEWSLETTER ACTIVITY

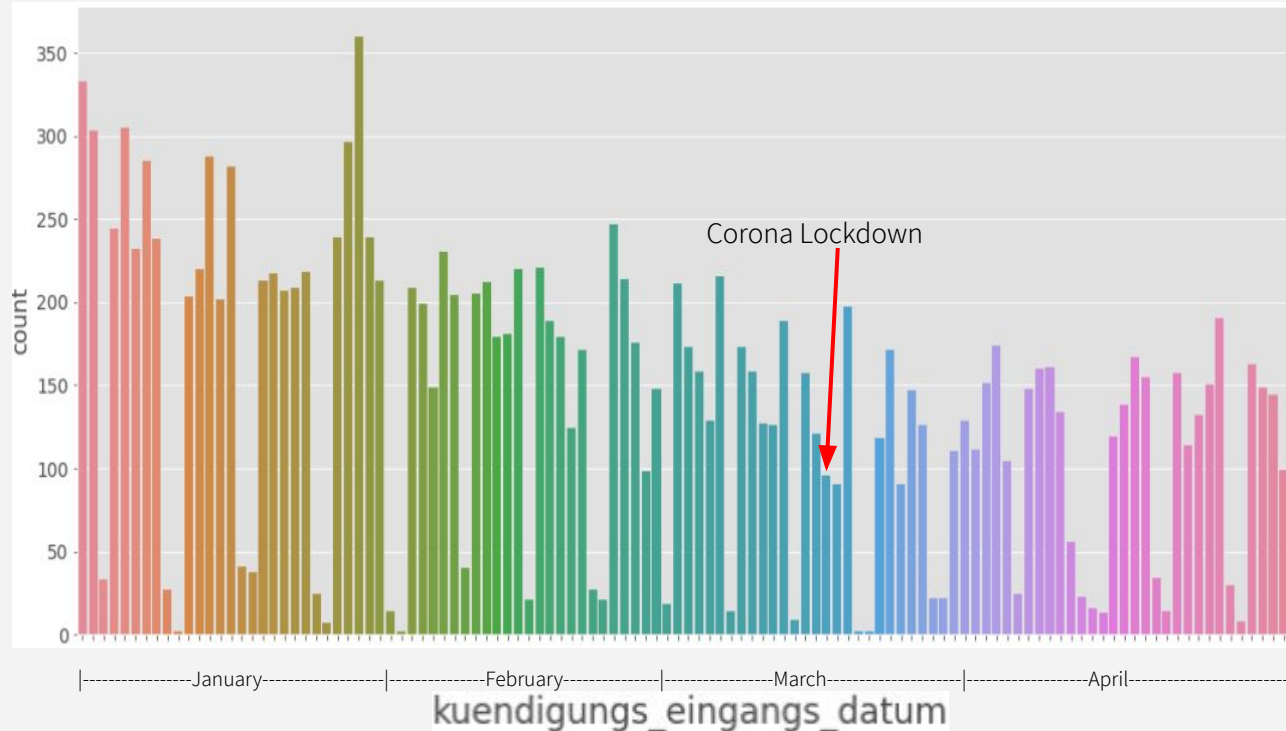
- The newsletter activity in the households is a strong predictor
- If the household in question has more than 5 active newsletters, its churn rate drops sharply.





## DOES CORONA HAVE AN EFFECT ON THE CHURN RATE?

- We do not have a comparable period, as we are only looking at one year
- At the end of the year, at the beginning of the next year, there is increased churn activity
- Two days after the lockdown on March 18, 2020, there is a small peak in the churn rate



**03**

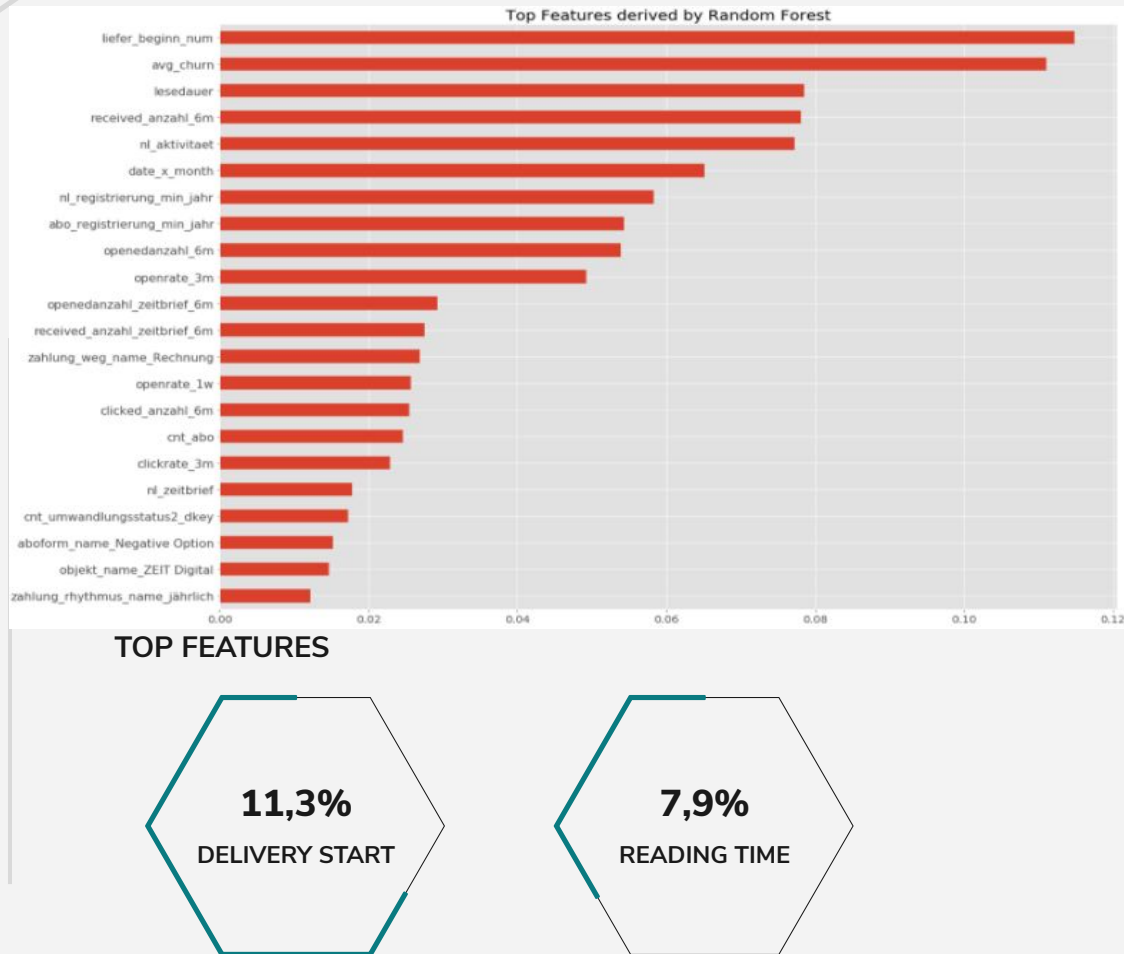
FEATURE IMPORTANCE





## FEATURE SELECTION RANDOM FOREST

We had to drop the top feature 'avg\_churn' because of data leakage. Delivery start, reading time, subscription registration and newsletter registration are top explanatory features.



# TOP 20 FEATURES

After checking the multicollinearity between the features, these top features remain

## Features

- zahlung\_weg\_name\_Rechnung: Payment method, in particular invoice payment method
- zahlung\_rhythmus\_name\_jährlich: Payment frequency, especially annually
- aboform\_name\_Negative Option: Subscription form especially negative option
- objekt\_name\_ZEIT Digital: Object especially ZEIT Digital
- nl\_zeitbrief: Newsletter Zeitbrief
- liefer\_beginn\_num: Delivery start
- nl\_aktivitaet: Newsletter activity
- abo\_registrierung\_min\_jahr: The year in which the first subscription was ordered
- nl\_registrierung\_min\_jahr: First registration of the newsletter
- received\_anzahl\_6m: Number of newsletters received in 6 months
- date\_x\_month: Day of termination and random date
- openedanzahl\_6m: Number of newsletters opened in 6 months
- openrate\_3m: How many times was the EMail opened in 3 months
- received\_anzahl\_zeitbrief\_6m: Number of newsletter Zeitbrief received within 6 months
- openedanzahl\_zeitbrief\_6m: Number of open newsletter Zeitbrief within 6 months
- clicked\_anzahl\_6m:: Number of clicks within 6 months
- cnt\_abo: How many additional subscriptions per household
- cnt\_umwandlungsstatus2\_dkey: Number of additional subscriptions that also converted (not just free trial period)
- openrate\_1w: Number of opened EMail in a week
- clickrate\_3m: Number of clicks within 3 months



# 04

## MACHINE LEARNING MODELS

## WHAT IS OUR GOAL?

Every model has a threshold that discerns positives from negative predictions. The lower the threshold the more instances get predicted positive. We can tweak the machine learning models by higher or lower the threshold.



### PRECISION

What is the percentage of correctly predicted positives? As a percentage: How often did the model predict positives?



### ACCURACY

How often the model has been right? When one class is very rare it leads to false conclusions.



### RECALL

How many of the positives did the model correctly predict? The recall is also called sensitivity.



### F1-SCORE

Jupiter is the biggest planet in our Solar System and the fourth-brightest object in the sky



## RESULTS ON TRAINING DATA SET

- Three of the machine learning models stand out in terms of the results on the training data set
- The PRECISION values are very high with up to 93% and even the RECALL values are astonishingly high with at least 81%, which could be a signal for overfitting
- On the next page we optimize three models via grid search and look at the results on the test data set

Model	Accuracy	Recall	Precision
0 logreg	0.66	0.67	0.66
1 tree	0.65	0.66	0.64
2 forest	0.92	0.94	0.90
3 XGB	0.68	0.66	0.69
4 AdaBoost	0.67	0.68	0.67
5 KNN	0.75	0.81	0.72
6 SVM	0.70	0.68	0.70
7 Stacking	0.92	0.89	0.93

# Machine Learning Models optimized vía Grid Search



## KNN

KNN classifies based on similarity measures. KNN is a non-parametric algorithm.



## XGBoost

XGBoost is an optimized gradient boosting library. XGBoost provides a parallel tree boosting and can solve problems beyond billions of examples.

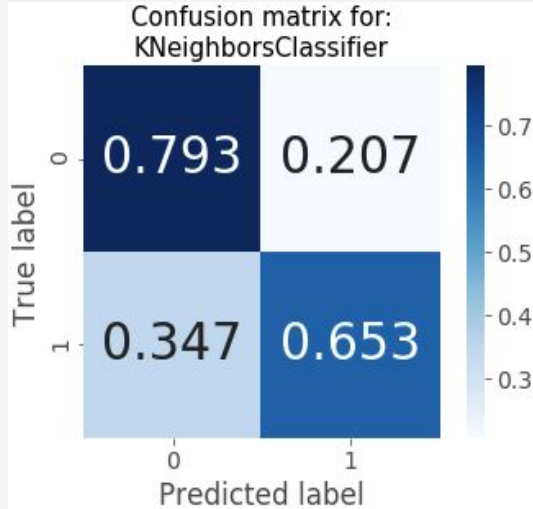


## Random Forest

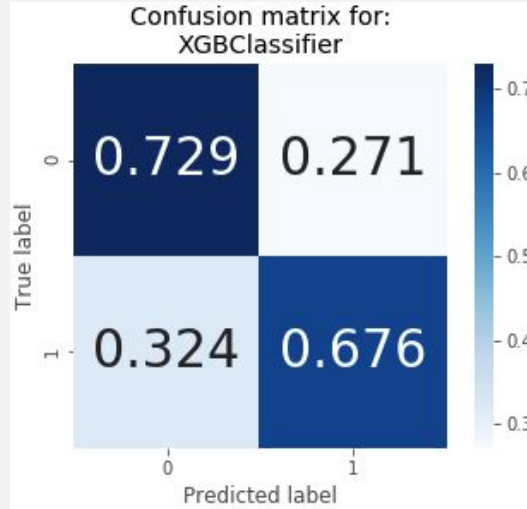
The random forest is a supervised learning algorithm that randomly creates and merges multiple decision trees into one forest.

## RESULTS ON TEST DATA SET

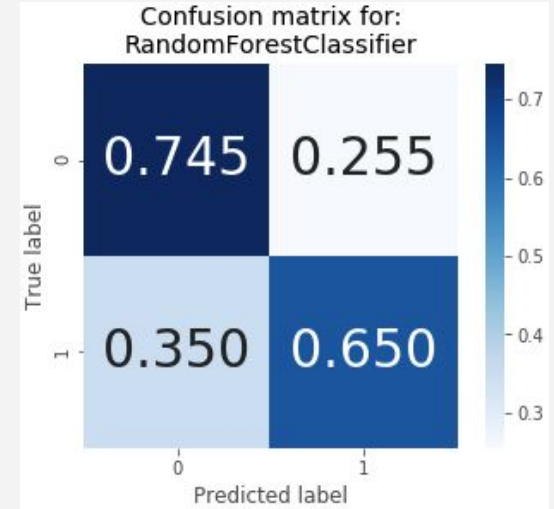
KNN



XGBoost



Random Forest



ROC AUC Score: 0,72

0,70

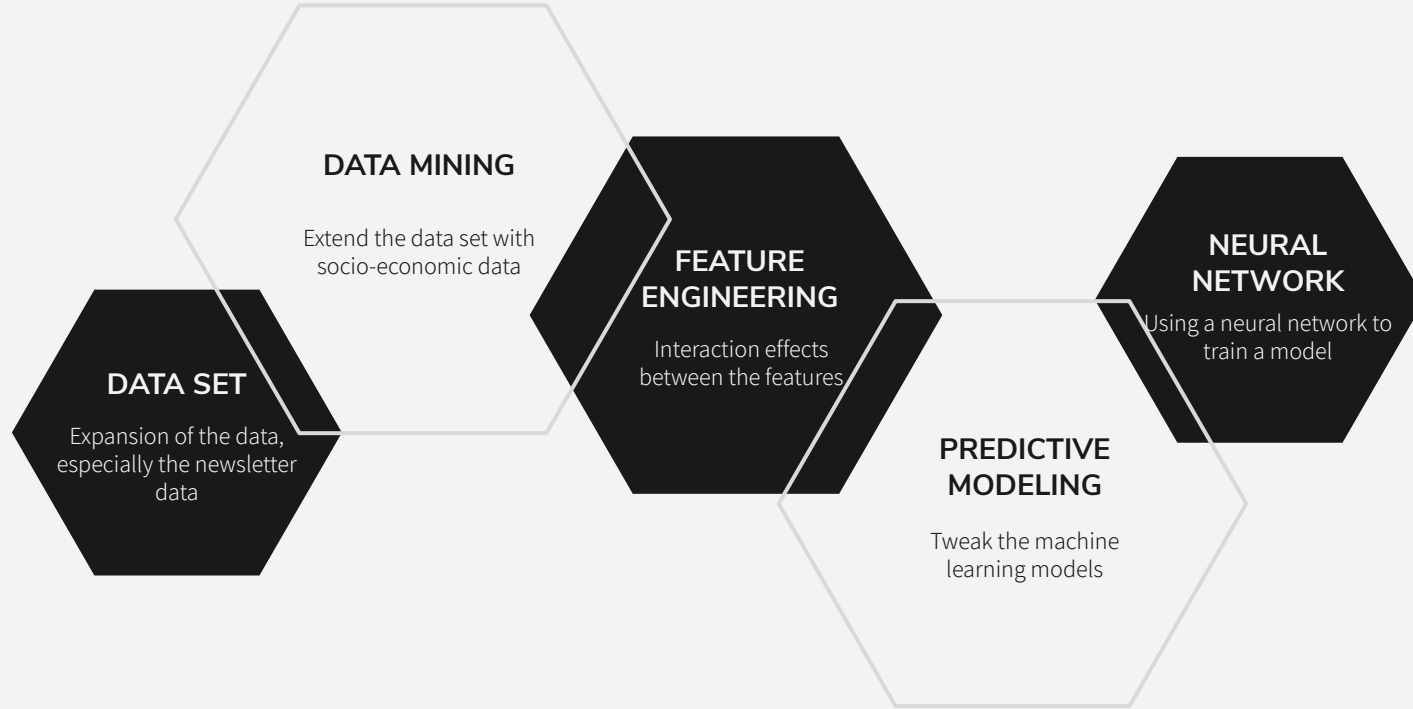
0,69

The image is a composite graphic. On the left, a large black triangle points towards the center. The background on the right is a soft-focus photograph of a sunset or sunrise, with a hand holding a white paper airplane. The text '05' is prominently displayed in the center-left area.

**05**

**FUTURE WORK**

## FUTURE WORK







# THANKS

Does anyone have any questions?

fuhrer\_sebastian@web.de

+49 174 6363886

<https://github.com/booty7>

