# Exercise 1 Classification

e12045110 Maria de Ronde
e12045110 Quentin Andre
e11921655 Fabian Holzberger

April 24, 2021

# Introduction

**Applied Algorithms**

**Performance Metrics**

# Amazon Reviews Dataset

**Dataset Description**

**Pre-Processing**

**Parameter-Tuning**

**Performance-Analysis**

# Congressional Voting Dataset

**Dataset Description**

**Pre-Processing**

**Parameter-Tuning**

**Performance-Analysis**

# Email Spam Dataset(link to dataset)

**Dataset Description**

The task of the email spam dataset is to predict if an email is spam or not. The link above contains three datasets from which we choose two, namely the `lingSpam.csv` and `enormSpamSubset.csv` for our project, since they have no missing values and the same layout. The dataset contains 12604 emails where 43.11
Every email has a binary target-label assigned, such that a 0 marks non-spam and a 1 marks spam emails. In figure 2 the distribution of the characters per email is shown. We see that most emails have a lenght in the range of 100 to 10.000 characters.

| Index | Body | Label |
|-------|------|-------|
| 100 | Subject: inexpensive online medication here pummel wah springtail cutler bodyguard we ship quality medications overnight to your door !... | 1 |
| 6006 | Subject: organizational changes we are pleased to announce the following organizational changes : enron global assets and services in order to increase senior management focus on our international businesses... | 0 |

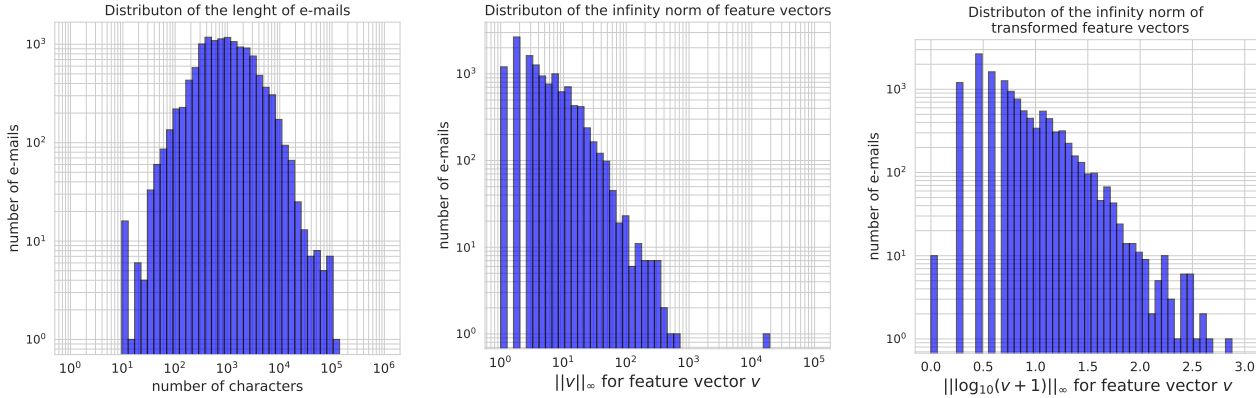Figure 1: Structure of the Email-Spam Dataset



Figure 2: left: Distribution of e-mail lenghts, middle: Distribution of maximum norm in extracted word vectors with lenght 8000, right: Distribution of maximum norm in extracted word vectors with lenght 8000 after removing outlyers and applying logarithmic transformation

## Pre-Processing

In the dataset are 213 duplicate emails that we first remove and then peform the train/test-split where the trainset size is 20% of the original dataset. Next we apply the Bag of Words feature extractor to each email. The algorithm converts every email to a vector $v \in \mathbb{Z}^N$ of intergers. First we create a list of all words and count their occurences in all emails. Then we take the $N$ most common words and count the occurences of the most common words in every email to get $v$. Before applying the Bag of word extractor we pre-process emails by the following steps:1. remove links (http...), 2. remove all characters exept alphabetical chars and numbers, 3. convert uppercase to lowercase, 4. split text-bodys into separate words, 5. lemmatize all words, 6. remove stopwords. For the steps 5., 6. we use nltk python package. By the preprocessing we reduce the number of distinct words from 126019 to 103759 words.

In figure 2 middle we see the distribution of the maximum norm of the extracted vectors. One can identify that the maximum norm spans several orders of magnitude from 0 to $10^5$. Especially there is only one vector $v$ with $||v||_\infty > 10^4$. Outlyers with $||v||_\infty > 10^3$ are therefor removed in the testset. Additionally we apply the logarithmic transformation $\log_{10}(x + 1)$ to all the elements of a vector and obtain a $|| \cdot ||_\infty$ distrubution that is bounded by the maximum magnithude. Note that we add 1 to all components of a vector since this component is 0 after the logarithmic transformation. The distribution after the transformation is shown in figure 2 right.

## Parameter-Tuning

For the parameter Tuning we further split the trainingset into a validation and trainingset. The training is done on the trainingset and the parametertuning with regard to the performance on the validation set. We

select parametervalues by the cross validation performance, since we assume this performance if more stable when compared to the holdout validation.

**Perceptron**: We compare cross-validation with 10 splits to a holdout-validation. In figure 3 upper left we see the influence of the scaling method on the f1-score of the perceptron algorithm. The binary scaling means that the transformed vectors have value 1 in a component if the original vectors component is nonzero and 0 otherwhise. Binary and logarithmic sclaing have the best performance in the cross validation, where the logarithmic scaling has slightly better performance on the trasining set. Therefore we choose the logarithmic scaling for the dataset. In 3 upper middle the influence of the extracted features on the perceptron algorithm with default parameters is shown. The performance for the perceptron is optimal for 2000 features since the f1-score saturates for this feature number. Therefore we adapt this number of features and investigate the learing rate in figure 3 upper right. Here the optimal learing rate is 0.1 if we consider the cross-validation performance on the validaton set. Note that the holdout performance is not optimal for this value. The influence of the tolarance can be seen in 3 bottom left. For values lower than 0.001 it has a constant f1-score that dropt for the cross validation on the validation set when we increase it. We can therefore fix the tolarance at 0.001 and investigate the maximum number of terations in 3 bottom middle. Surprisingly a low number of iterations causes the f1-score to increase.
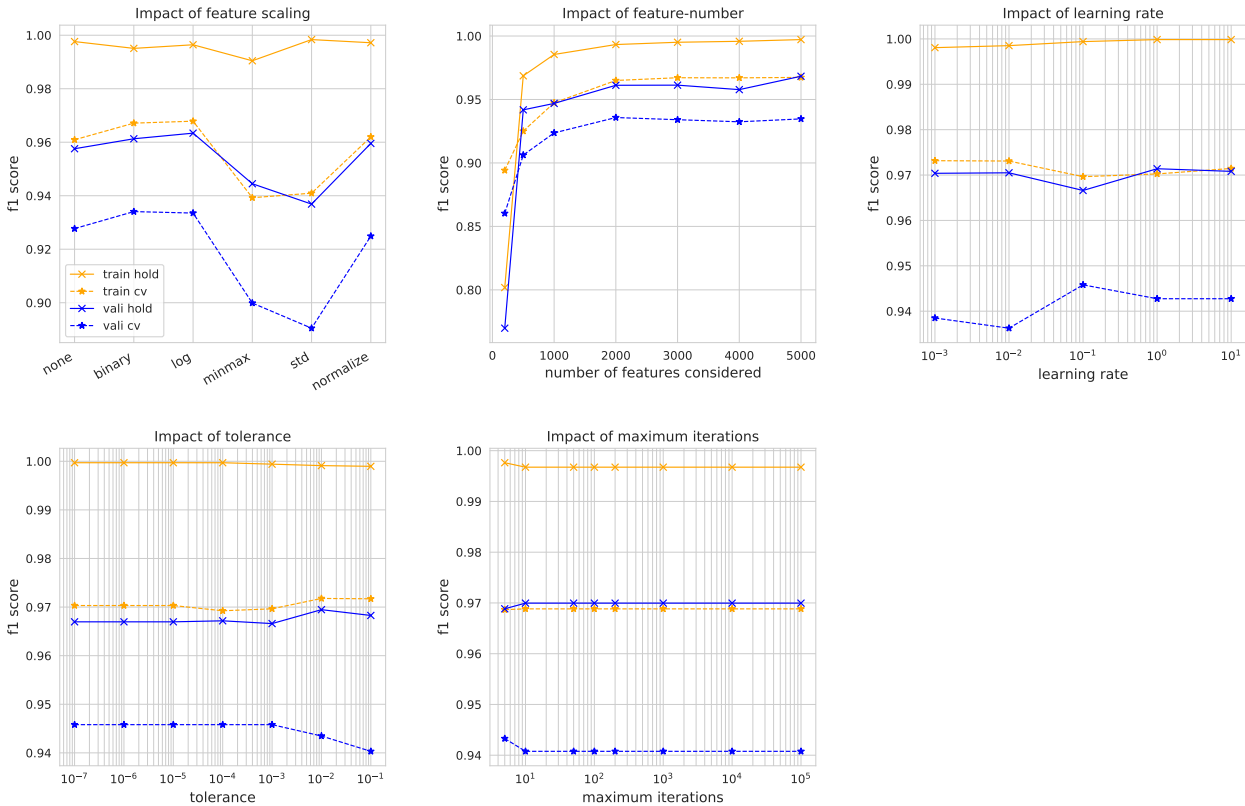


Figure 3: F1-score of Perceptron for different pre-processing parameters and model-parameters

For the perceptron the cross validation yields always a more pessimistic performance when compared to the holdout-validation. As optimal parameters we have chosen: scaling method: logarithmic, extracted features: 2000, learning rate: 0.001, tolerance: $1e - 6$. The f1-score on the testset is for this parameters: 0.97 for the holdout validation and 0.95 for the cross validation.

**Random forrest**: Since the random forrest has much larger runtime compared to the other algorithms in this

project we will not evaluate its performance on the training set and further set the number of validation steps in the cross validation to 4.
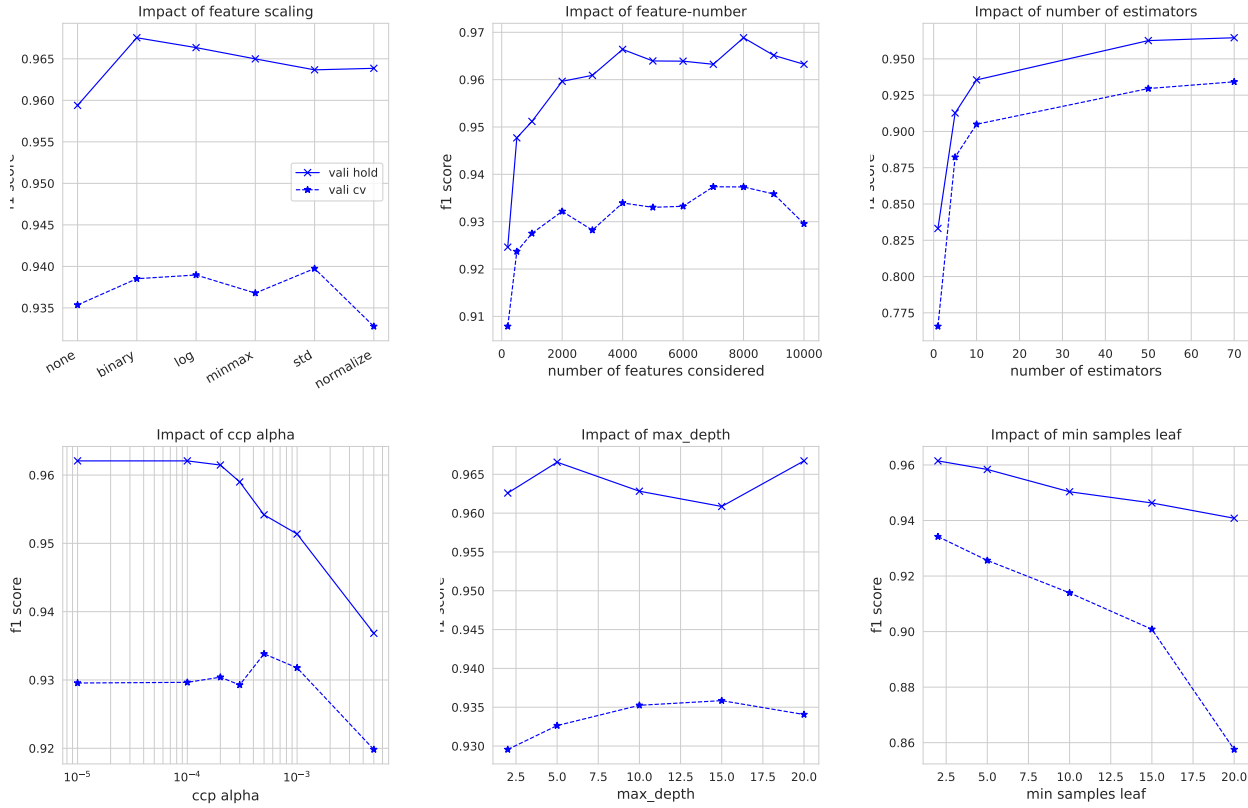


Figure 4: F1-score of Random-Forrest for different pre-processing parameters and model-parameters

The random forrest algorithm is usually not affected in its performance by scaling. Nevetheless we inspect the influence of the scaling parameters in at top left with 8000 features extracted. We see that no scaling gives even for this algorithm less performance. Since the difference for the cross-validation is small we keep the unscaled data for this algorithm. When considering the impact of the number of extracted features in top middle the f1-score in the cross validation improves till 7000 features are reached and then drops again. Therefore to decrease the runtime as much as we can we select 7000 features even if the holdout f1-score here is clearly lower. Another important parameter for performance and runtime in the random forrest is the number of trees that are build. The performance of this parameter can be seen in top right. The performance increases when we build more decision trees but at 50 trees this preformance increase is not significant anymore. Therefore to decrease the runtime we select 50 trees as optimum. So far the decision trees are grown to full lenght and might overfit largly. In bottom left we perform cost-complexity pruning on the tree. Larger values of the parameter will post-pune more tree branches. At $5e-3$ we reach an optimum in the cross validation f1-score which doesent differ significantly from the other performances achived for different values of this parameter. We still want to select parameter as large a possible to simplify the tree and prevent overfitting for the generqal case such that we select $5e-3$ for this parameter. We can also prevent the tree from overfitting by bounding its depth. The impact of this method is shown in bottom right. Here we obtain a optimum for a maximal tree depth of 15. The last parameter to prevent overfitting we investigate is the minimum number of spamples required to split a leaf into a new branch. This parameters performance impact can be seen in bottom right. Here increasing the parameter values drastically decreases the f1-score and therefore we dont make use of it.

In our analysis we obtained the best performance with : **Naive Bayes**:
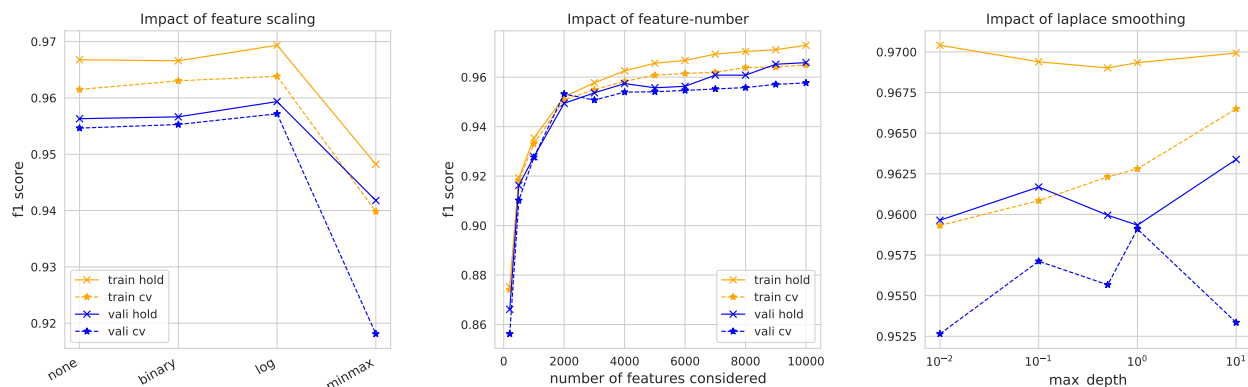


Figure 5

**Performance-Analysis**

# Bridges Dataset([link to dataset](#))

**Dataset Description**

**Pre-Processing**

**Parameter-Tuning**

**Performance-Analysis**

# Conclusion