

Vorhersage von Wirtschaftsindikator mittels Sentimentenanalyse von Nachrichtenartikeln und maschinellem Lernen

DIPLOMARBEIT

zur Erlangung des akademischen Grades

Diplom-Ingenieur

im Rahmen des Studiums

Wirtschaftsinformatik

eingereicht von

Christoph Hämmerle, BSc
Matrikelnummer 01226577

an der Fakultät für Informatik

der Technischen Universität Wien

Betreuung: Univ.Prof. Dr. Allan Hanbury

Mitwirkung: Dr. Navid Rekabsaz

Wien, 5. März 2019

Christoph Hämmerle

Allan Hanbury

Prediction of an Economic Indicator using Machine Learning and Sentiment Analysis of News Articles

DIPLOMA THESIS

submitted in partial fulfillment of the requirements for the degree of

Diplom-Ingenieur

in

Business Informatics

by

Christoph Hämmerle, BSc

Registration Number 01226577

to the Faculty of Informatics

at the TU Wien

Advisor: Univ.Prof. Dr. Allan Hanbury

Assistance: Dr. Navid Rekabsaz

Vienna, 5th March, 2019

Christoph Hämmerle

Allan Hanbury

Erklärung zur Verfassung der Arbeit

Christoph Hämmerle, BSc
Weyringergasse 31/8, 1040 Wien

Hiermit erkläre ich, dass ich diese Arbeit selbständig verfasst habe, dass ich die verwendeten Quellen und Hilfsmittel vollständig angegeben habe und dass ich die Stellen der Arbeit – einschließlich Tabellen, Karten und Abbildungen –, die anderen Werken oder dem Internet im Wortlaut oder dem Sinn nach entnommen sind, auf jeden Fall unter Angabe der Quelle als Entlehnung kenntlich gemacht habe.

Wien, 5. März 2019

Christoph Hämmerle

Danksagung

Zuerst möchte ich mich bei meiner Familie bedanken. Sie machte es mir überhaupt möglich das Studium zu absolvieren, nicht nur aus finanzieller Sicht sonder auch die ständige Unterstützung bei auftretenden Problemen. Weiters möchte ich mich bei allen bedanken die mich während dieser Zeit begleitet und vor allem unterstützt haben, sei es fachlich oder emotional.

Ganz besonders will ich mich auch bei meinen Betreuern Allan und Navid bedanken. Mit Ihrer Hilfe wurde diese Arbeit zu dem was sie jetzt ist. Durch zahlreiche Rückmeldungen und Vorschläge die Arbeit zu verbessern und auf eine höhere Stufe zu heben, habe ich sehr viel gelernt.

Christoph

Acknowledgements

First, I want to thank my family. Without them my studies would not have been possible. They not only supported me financially but also emotionally. Their constant assistance made this possible. Further, I would like to thank everyone who accompanied and supported me on this path.

Allan and Navid, I would like to thank you for the constant help during this thesis. The countless feedback and advice made this thesis to what it is now. Through your support I was able to learn a lot.

Christoph

Kurzfassung

Das Ziel dieser Diplomarbeit ist es den Zusammenhang zwischen Meinungen und Stimmungen in Zeitungsartikeln und dem Kursverlauf des Austrian Traded Index (ATX) zu untersuchen. Dafür wenden wir bewährte Methoden aus den jeweiligen Forschungsgebieten an. Mithilfe eines Webscrapers extrahieren wir Zeitungsartikel von der Website einer österreichischen Tageszeitung. Die Polarität und die Stimmungsbilder in diesen Artikeln werden automatisiert untersucht und ermittelt. Dafür erstellen wir ein Lexikon, welches positiv und negativ assoziierte Wörter im Hinblick auf die Wirtschaft und die wirtschaftliche Entwicklung enthält. Relevante Wörter aus einem bekannten deutschen Lexikon, das SentiWS Lexikon, werden von uns extrahiert und zusätzlich durch Wörter, welche als positiv bzw. negativ im Bezug auf wirtschaftliche und finanzielle Entwicklung interpretiert werden, ergänzt. Wir erstellen weiters ein Model um Preisschwankungen und Volatilität des ATX vorherzusagen. Die Arbeit bestätigt, dass die Zeitungsartikel nützliche Informationen, um den Kursverlauf des ATX zu erklären, enthalten. Das erstellte Lexikon ist nicht eindeutig besser als das generelle SentiWS Lexikon. In einigen Testszenarien hat es jedoch durchaus Vorteile. Weiters bestätigt diese Arbeit die Ergebnisse anderer Arbeiten auf diesem Gebiet. Vor allem negative Aussagen werden erkannt und werden als einflussreich in den verwendeten Modellen erachtet. Im Vergleich dazu können positive Aussagen kaum mit positiver Entwicklung des Kurses in Verbindung gebracht werden. Die Resultate der getesteten Modelle sind stark durch diese Negativität beeinflusst.

Abstract

The analysis of textual data and their predictive quality has gained the interest of many researchers, especially in the financial domain. This thesis investigates whether newspaper articles contain information to describe the changes of the Austrian Traded Index (ATX). We apply state of the art methods to extract newspaper articles from the online platform of an Austrian newspaper, to perform sentiment analysis of the articles and to build machine learning models in order to predict price and volatility developments of the ATX. As the newspaper articles are written in German, we create a new sentiment lexicon, called German Financial Sentiment Lexicon (GFSL), by extracting sentiments from the SentiWS, a general German sentiment lexicon, and adding financial sentiment words to the lexicon. Our findings show the newspaper articles contain information which allow predictions of price and volatility movements. The GFSL does not clearly outperform the SentiWS lexicon, although in some scenarios it clearly has an advantage over the general lexicon. The results confirm the findings of previous studies such that negative sentiments highly influence the outcome of the model while positive sentiments are hardly relatable to positive development of the index.

Contents

Kurzfassung	xi
Abstract	xiii
1 Introduction	1
1.1 Research Questions	1
1.2 Contributions	2
1.3 Structure of the Thesis	4
2 Background	5
2.1 Machine Learning	5
2.1.1 Types of Machine Learning	6
2.1.2 Machine Learning Process	8
2.1.3 Algorithms	10
2.2 Sentiment Analysis	12
2.2.1 Development and Applications	13
2.2.2 Common Sentiment Analysis Approaches	13
2.2.3 Levels of Analysis	15
2.2.4 Text Processing	16
2.2.5 Challenges in Sentiment Analysis	16
3 Related Work	17
3.1 Sentiment Analysis in the Financial Domain	17
3.2 Sentiment Lexicons	20
4 German Financial Sentiment Lexicon	23
4.1 Motivation	23
4.2 Used Resources	23
4.3 Creating GFSL	24
5 Sentiment Analysis of the Newspaper Articles	27
5.1 Data Accumulation	27
5.1.1 Web Scraping	28
5.1.2 ATX-Quotation Data	29

5.2	Sentiment Analysis	29
5.2.1	Approach	29
5.2.2	Sentiment Analysis Algorithm	30
5.3	Defining the Label	34
5.4	Experiment Design	36
6	Results and Discussion	37
6.1	Results of the Regression	38
6.2	Results of the Classification	40
6.3	Qualitative Analysis	42
6.4	Discussion	44
7	Conclusion and Future Work	45
	List of Figures	47
	List of Tables	49
	List of Algorithms	51
	Bibliography	53
	Attachments	57
	GFSL	57

Introduction

Economic development is considerably influenced by sentiments, opinions, expectations and events happening around the world. The behavior and actions of individuals may be related to some force or effect which influenced them. Often newspapers are the source of information in the daily lives of people. Some articles communicate positive feelings and other negative feelings to the reader. Of course different people are not influenced the same by one article and some articles might not cause any reaction in an individual. However, the general mood, communicated by the article, can be captured via sentiment analysis. In particular, the behavior of participants in the stock market may be influenced by newspaper articles and the aggregated effects they take may be correlated to the sentiments communicated via the newspaper articles. In this case, the aggregated effects reflect the development of the economic indicators, and especially in this thesis the Austrian Traded Index (ATX).

The ATX, like every stock market index, is a general indicator of growth or decrease of the economy it represents and reflects the general mood of stock market participants regarding future developments. This thesis investigates the extent to which the communicated sentiments of newspaper articles are relatable to the development of the ATX.

Section 1.1 describes the research questions of this thesis. In Section 1.2 we present the contributions of this thesis, followed by a brief description of the methods we use. In the last part of this chapter, Section 1.3, we provide an overview of the thesis.

1.1 Research Questions

The analysis of textual data and their predictive quality has gained the interest of many researchers, especially in the financial domain. The research in this area is mostly based on the analysis of texts written in English. However, there is a lack of research on German text in the financial domain.

The first research question investigates whether the extracted sentiments correlate with

the ATX observations. We test if the sentiment weights add relevant information to the machine learning models in order to explain the development of the ATX.

For the second research question we investigate which kind of data, extracted from the ATX observations, is best described by the sentiment of the newspaper articles. We consider volatility and price movements. We also classify these time series into two groups, up and down movements, and investigate the performance of classifying price and volatility movements by training supervised machine learning models with the sentiment weights of the newspaper articles. With the models we predict the time series and compare the predictions to the actual observations of the index. We answer which kind of data, extracted from the price observation of the ATX, is best described by the sentiment weights of the newspaper articles. Do the sentiments of the newspaper articles allow qualitative predictions?

For our third research question we compare the performance of the two lexicons according to the accuracy of their predictions, resulting from the supervised machine learning models. Does the domain specific lexicon provide advantages over the general sentiment lexicon? Previous studies show the importance of domain-specific lexicons in sentiment analysis of financial data. We assume the domain specific lexicon has an advantage over the general lexicon as it provides domain specific expressions and therefore is enabled to capture the communicated sentiments more accurately resulting in higher predictive performances.

1.2 Contributions

This thesis contributes to the research area of sentiment analysis of German financial texts in two ways.

First, the creation of the German Financial Sentiment Lexicon (GFSL), a sentiment lexicon which contains words that are associated to a sentiment in the financial domain. Creating a lexicon in the German language, considering the financial context, is an important issue. The quality of the predictions from supervised machine learning models highly depends on the lexicon used to obtain the sentiment weights as the lexicon is the essential component to capture expressions in the newspaper articles. We create the GFSL by extracting words, associated with the financial domain, from the SentiWS lexicon, a well established general German sentiment lexicon [RQH10]. Further, we extend the GFSL with expressions from General Inquirer dictionaries by manually translating economic and financial expressions. We also manually add sentiment words from articles published on the Vienna stock exchange website. The GFSL lexicon is evaluated by comparing the predictions, obtained from the machine learning algorithms, to the predictions computed using the general SentiWS lexicon for sentiment analysis.

In the second contribution, we explore the applicability of newspaper articles written in German as a source of information to explain the changes of the Austrian Traded Index (ATX). The outcomes contain results of several machine learning algorithms applied to data extracted from the ATX observations, considering price differences and volatility. We investigate whether the sentiments communicated in the newspaper articles allow

predictions of the stock market price and the volatility by comparing various machine learning algorithms according to the accuracy of their predictions. We apply state of the art approaches to our problem of extracting newspaper articles, performing sentiment analysis and building machine learning models. Each step is explained in the following enumeration.

1. Data Accumulation

To obtain relevant newspaper articles we implement a web scraper. The web scraper extracts data from websites, in this case newspaper articles from the online archive of *derStandard*, an Austrian daily newspaper. The need for searching online archives and solely obtaining articles which may have an impact on market development, especially articles categorized as economic, requires the implementation of a web scraper. Further, we process the extracted data and store it in a database in order to continue with the sentiment analysis.

2. Sentiment Analysis

We use sentiment analysis to extract the attitude provided by the newspaper articles gathered in the first step. In order to do so we compare the words of an article to the GFSL and the SentiWS lexicon, a general German sentiment lexicon. Further, we process every newspaper article, count the occurrences of the words in the lexicons in the newspaper articles and weight these occurrences according to common weighting schemes. We use the resulting time series of weights in the third step.

3. Machine Learning

Finally, we select supervised machine learning algorithms to examine the relationship between the obtained sentiment weights and the Austrian Traded Index (ATX). We investigate various settings to find the data which promises the most accurate results and allows to draw conclusions about the relationship between the sentiment weights and the ATX. We test classification and regression scenarios like price changes, volatility changes and volatility directions. The quality of the classification models is evaluated by comparing the results to the most frequent group classifier. The regression cases are tested regarding their significance. The quality of the predictions is evaluated using cross-validation techniques.

The aim of this thesis is not to provide an exact estimation of the ATX with the aid of newspaper articles and the sentiments extracted from them since the index is influenced by various factors. The work should be rather seen as an attempt to verify whether this approach leads to insights in the extraction of daily newspaper articles, written in German, from the online archive of an Austrian newspaper and whether their sentiments are relatable to economic development.

1.3 Structure of the Thesis

In Chapter 2, we describe the **Background** for this thesis. We explain the concepts of machine learning and sentiment analysis. The various approaches are briefly discussed and illustrated with examples. The examples also provide an insight into the wide area of applications.

In Chapter 3, we explain **Related Work**. Other research and its most important results are presented.

In Chapter 4, we explain the creation of the **German Financial Sentiment Lexicon** used for the sentiment analysis of the newspaper articles. We explain the process of creating the GFSL and our motivation behind it.

Chapter 5 explains the details of our approach to **Sentiment Analysis of the Newspaper Articles**. First, we explain how we accumulated the data. We then introduce our approach to perform the sentiment analysis of the newspaper articles. In the last section of this chapter we describe the design of our experiments.

In the first part of Chapter 6, **Experiment Results**, we explain and illustrate the outcomes of the experiments. In the second part we qualitatively analyze the results of the experiments. Finally, we discuss the main findings.

The last chapter, **Conclusion and Future Work**, summarizes the outcomes and conclusions of this thesis and presents further research possibilities in this area.

Background

Chapter 2 provides the background needed for the analysis of the newspaper articles. Section 2.1 provides an introduction to machine learning. Different types of machine learning and their applications will be discussed. Furthermore, the general approach to build and evaluate a machine learning model is briefly discussed. In Section 2.2 we explain sentiment analysis. We define the term sentiment analysis and introduce various areas of application. Finally, we introduce different approaches and levels of analysis to illustrate common challenges in this research area.

2.1 Machine Learning

Machine learning is a research field which addresses extracting knowledge from data. Machine learning algorithms aim to learn patterns in the provided input data in order to generalize decision making processes to unseen data. Manually designing rule-based systems is feasible for applications in which humans have a good understanding about the rules and the general functionality which they aim to describe. Nevertheless, there are disadvantages regarding rule-based systems. For example the defined rules are domain specific. Changing the domain slightly might result in a complete rework of the defined decision rules. Further, defining rules requires a deep understanding about the domain. Machine learning provides the tools to automate the decision-making process by generalizing from input data [MG16].

In the first part of this section we describe the different types of machine learning. Examples illustrate their area of application. The second part of this section describes the general machine learning process to build a machine learning model. Finally, we briefly explain the machine learning algorithms used for this thesis.

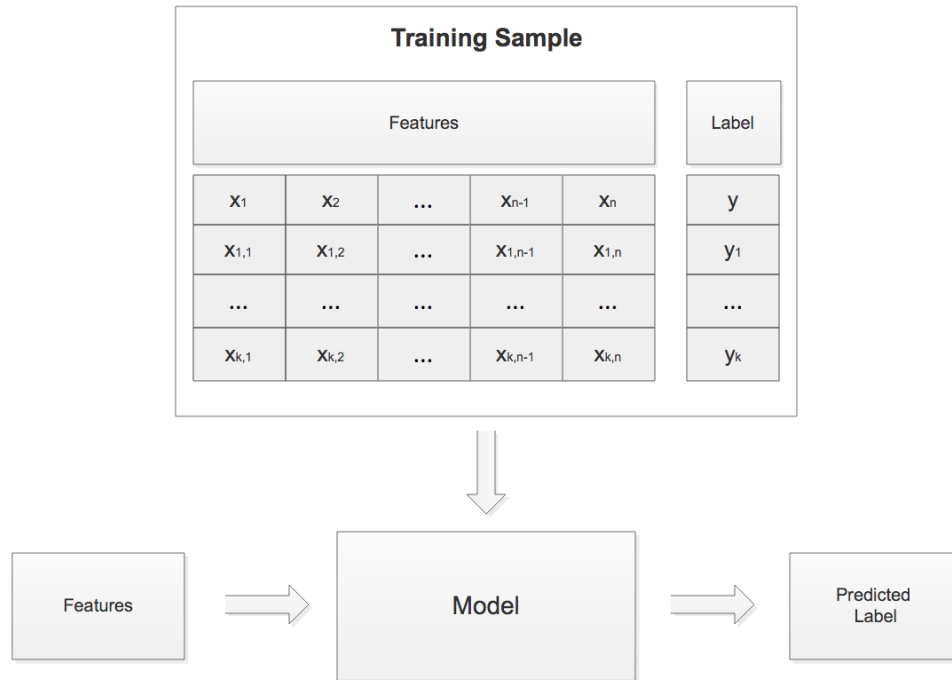


Figure 2.1: Supervised Learning

2.1.1 Types of Machine Learning

Raschka [Ras15] distinguishes between three types of machine learning: supervised learning, unsupervised learning and reinforcement learning. The fundamentals and the differences between these types are explained in this section. Examples will illustrate the application area of the different types of machine learning.

Supervised Learning

The goal in supervised learning is to learn a model with given data and given labels. The features describe the given input data and the label represents the output of the model. This means the model is trained with input values and the corresponding given output value. The trained model is then used to process unseen data and predict the output. The term supervised refers to a set of samples where the desired output is already known and used to train the model [Ras15]. Figure 2.1 visualizes the approach. The first row describes the features names X_1, \dots, X_n and the label name y . The consecutive rows display the features of different observations. The model is trained with the feature set and the corresponding label set. Each observation consists of features and the corresponding output. After the model is trained the same features are extracted from

unseen input data. The supervised machine learning model is then enabled to predict the corresponding label.

As explained by Raschka [Ras15], the supervised machine learning methodology consists of two major subcategories. The first is classification. The goal is to predict categorical class labels of new data based on past data. An example for multi-class classification would be the recognition of letters from handwriting. The task is to classify handwritten letters into letters from the alphabet. To do so an algorithm is trained with several examples of handwritten letters and the corresponding correct classification. The algorithm should then be able to classify new observations with a certain accuracy. An example of a binary classification problem is the identification of spam emails. An algorithm is trained with emails and the corresponding categorization, either is spam or is not spam. According to the defined features the algorithm should then be able to classify future observations.

The second type of supervised machine learning is regression. The goal of regression analysis is to predict a continuous outcome from a given number of explanatory variables by finding a relationship between the explanatory variables and the response variable. An example for a regression model is the prediction of someones annual income with the level of education, the age and the place of living. The procedure to do so is the same, except for the algorithm, as for classification. The algorithm is trained with the feature data and the corresponding label, in this case education, age and place of living and the corresponding annual income. The algorithm is then able to predict from future observations [MG16].

In this work, we apply both supervised classification and regression machine learning methods to investigate the relationship between the newspaper articles and the changes of the ATX.

Reinforcement Learning

Reinforcement learning deals with the continuous improvement of a system based on the interactions of the system with the environment. The system has a current state and a reward signal. Through interaction with the environment the system learns a series of actions, via exploratory trial-and-error approach or planing, in order to maximize the reward signal. An example for such a system would be a chess player simulation. The system learns to perform moves according to the current state and the likely reward, in this case to win the game [Ras15].

Unsupervised Learning

Unsupervised learning is used to explore the structure of the data and to extract meaningful information without training the algorithm with a known outcome or defining a reward signal. Clustering and dimensionality reduction are two major applications.

Clustering is a technique to organize observations into groups which share a degree of similarity and are dissimilar to other groups. This technique is often used for marketing purposes to group users and target the different groups individually [Ras15].

Dimensionality reduction is another important application of unsupervised learning. Data

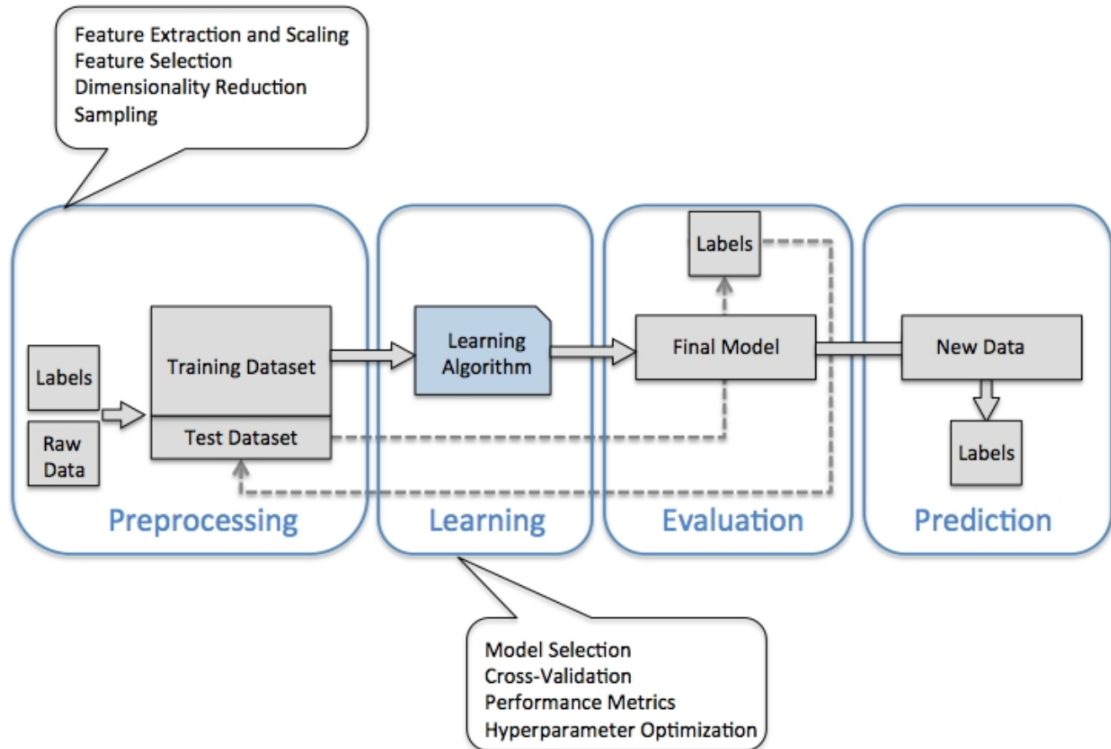


Figure 2.2: Machine Learning Process [Ras15]

often consist of multiple features. Due to limitation of computational power, storage capacity, computational time or to reduce noise in data the reduction of dimensions is helpful. The projection of multiple dimensions into two or three dimensions is useful to visualize data [Ras15].

2.1.2 Machine Learning Process

This section is about the process to train and evaluate a machine learning model in order to be able to accurately predict new data. We summarize the machine learning process defined by Raschka [Ras15], Figure 2.2, in the following.

Preprocessing

The preprocessing is about preparing the input data. As described in the supervised machine learning approach the data consist of the feature set and the label set. The feature set is used to describe and explain the behavior of the labels. The label set is the data which is intended to be described by the model. Feature extraction deals with the creation of features based on the input data. The features may differ regarding their value range and or distribution. In order to prevent unequal influence of features, based on their value range, on the machine learning algorithm, especially algorithms which

compute distances between features, the features might need to be scaled.

Features might not always contain useful information or the information of a feature is described by other features. In such a case it is desired to only use the features which indeed are useful to explain the label. Further, the feature set is often a high dimensional structure. Features often correlate with each other or do not provide information in order to explain the corresponding label. Therefore, it is important to reduce the dimension of the features. This reduction results in a less complex structure and allows future steps, like training the machine learning algorithm, to use less computational power and time. The reduced amount of dimensions may be varied to see whether a higher number of dimensions will increase the accuracy significantly or if the chosen reduction justifies the possible loss in accuracy.

In order to train and evaluate the machine learning model, the available data is split into two sets, called training- and testing set. The training set is used to train the algorithm and the testing set is used to evaluate the quality of the model. In order to test different settings a part of the training set may be further split. The resulting set is referred to as the validation set. The validation set is used to optimize hyper-parameters. Once the most promising setting of hyper-parameters is found, the model is applied on the testing data. The result of the testing data is the indicator of the quality of the model.

Learning and evaluation

The learning process of the algorithm highly depends on the previously processed data. Different feature sets and labels may lead to many input combinations for a machine learning model. Testing various feature sets with the corresponding label, which together with a suitable algorithm, provide accurate predictions is challenging and may take several attempts.

In order to compare the performance of machine learning models we use common metrics. For the regression algorithms we use the R-squared metrics to compare different experiments. The R-squared value provides information about how close the predictions compared to the actual observations are. It is defined as follows:

$$R^2 = 1 - \frac{\sum_{i=1}^n (\hat{y}_i - y_i)^2}{\sum_{i=1}^n (y_i - \bar{y})^2}$$

Where \hat{y}_i are the predicted values and y_i the actual observations. The numerator is referred to as the residual sum of squares, it describes the difference between the actual observations and the predictions. The value of the R^2 is between 0 and 1. If the residual sum of squares is zero, the predictions are equal to the observations, the fraction will be zero and the R^2 has its maximum value of 1. The denominator describes the variance of the actual observations.

In order to compare different classification approaches we use the prediction accuracy which is defined as the number of correct classifications divided by the total number of classifications.

Machine learning algorithms often consist of many hyper-parameters which affect the

behavior of the machine learning algorithm. These hyper-parameters have to be optimized to the given data. In order to optimize the hyper-parameters to our data we use grid search with cross-validation. Grid search is a technique to test all combinations of predefined parameter settings. For each parameter we define the settings we want to investigate. The machine learning algorithm is trained and evaluated with every combination of the parameter settings resulting in the optimal parameter values. It is crucial to avoid overfitting the model. This means to overly optimize the parameters to the specific input set. If the results of the validation and the testing set are highly different the model may be overfitted, meaning the hyper-parameters are tuned too much. Applying the algorithm on new data may then result in poor predictions although the algorithm provided good results on the validation set. Testing the final model on the test set indicates whether the model is able to adapt to new data. Raschka [Ras15] mentions in this context the term *generalization performance*. This means the algorithm has to have the ability to not only fit the training data but also to process new data and receive comparable results. If the performance is not satisfying or the model is overfitted the procedure starts from the beginning.

We retrain the resulting model with the optimized hyper-parameters on the training and validation set to increase the amount of training data. We compare the predictions to the test set of the label and compute the performance metrics. Furthermore, we provide information about the significance of the used model, answering if the features add useful information to the model. For the regression case the information about the significance of the results may be obtained via the p-value with a defined significance level. The classification results are evaluated using the different splits of the input set. The results are compared to a baseline. A baseline is a trivial classifier, for example a classifier which always predicts the most common group in the training set.

Prediction

After a suitable model is found and the performance is satisfying the model may be applied on new future data.

2.1.3 Algorithms

In this section we briefly explain the algorithms applied in this work. We provide sources which explain the algorithms in detail.

- **Support Vector Regression (SVR):** In the following we will briefly explain the basic functionality of SVR with the example in Figure 2.3. In this example the problem is reduced to a two dimensional linear regression problem. The goal is to find a function $f(X)$ that has the highest ϵ value, including many observations, and at the same time is as flat as possible. So observations with an error less than ϵ , between the tube $+\epsilon$ and $-\epsilon$ and the function $f(X)$, will be tolerated. Observations outside the tube, with a higher error than ϵ , will be considered in the optimization function. *"The constant $C > 0$ determines the trade-off between the*

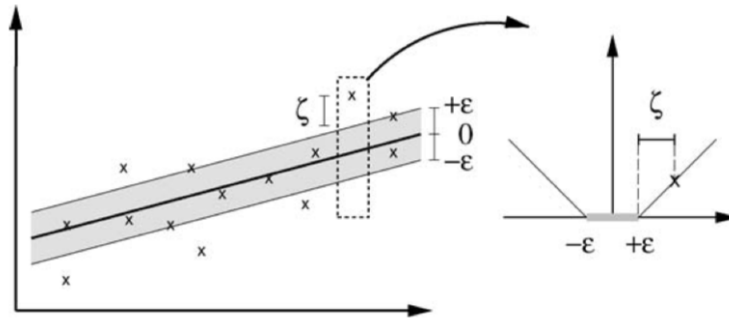


Figure 2.3: Linear SVR, Smola and Schölkopf [SS04]

flatness of f and the amount up to which deviations larger than ϵ are tolerated." So C can be seen as the penalty parameter for the error term. The observations outside the tube are penalized and contribute to a cost function which is minimized in the optimization problem [SS04].

This example is very simple but sufficient to explain the basic functionality. The data is not always linear, as in this example. To cover non-linear data the kernel can be varied. The kernel is used to map the input data into higher dimensional spaces in which the data is separable. The line may then be a curve but the general idea is the same. Further, the example is only two dimensional. Increasing the dimensions increases the complexity of the optimization problem. The function $f(X)$ then describes a high dimensional construct instead of a line. For a detailed description of SVR we refer to Smola and Schölkopf [SS04].

The SVR algorithm has various parameters to control the learning process. First, the kernel types like linear, polynomial, sigmoid and rbf. The penalty parameter, C , the parameter defining the size of the tube, ϵ , in which no penalty is assigned and the parameter regulating the tolerance, tol , for the algorithm to stop.¹

- **Support Vector Classification (SVC):** SVC is very similar to the previously described SVR algorithm. As mentioned by Hsu et al. [HCL03] the goal is to find a linear separating hyperplane with the maximal margin. Figure 2.4 illustrates the desired outcome in a two dimensional example. The different shapes of the observations signalize the group affiliation. The goal is to find a hyperplane which separates the groups, in this example the dashed line with the additional maximal margin. Similar to the SVR algorithm different hyper-parameter settings can be selected, like the kernel and again the penalty parameter C .
- **Random Forest Regression (RFR):** Breiman [Bre01] provides a detailed description of the Random Forest algorithm. We will briefly explain the functionality

¹<http://scikit-learn.org/stable/modules/generated/sklearn.svm.SVR.html>, last accessed 05.09.2018

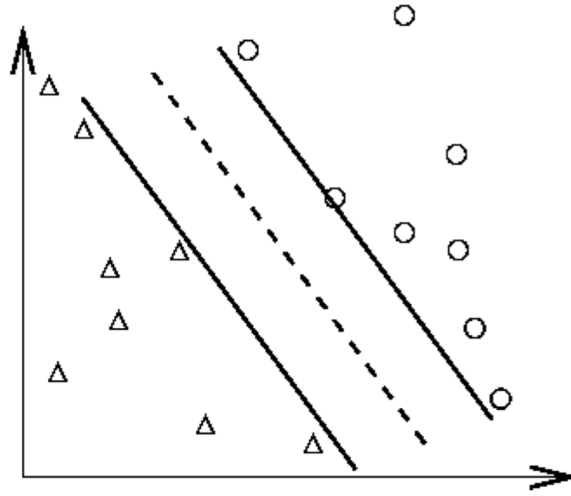


Figure 2.4: SVC Example [HCL03]

according to the documentation of the used implementation.² The random forest algorithm follows an ensemble approach. The goal is to combine the predictions of several base estimators, in this case decision tree estimators, to improve generalizability and robustness of the estimation. Each tree in the ensemble is built from a sample drawn with replacement from the training set. Randomness is added to the algorithm by choosing the best split among a random subset of the features. In the regression case the final estimation is calculated by averaging the output of all trees in the ensemble.

- **Random Forest Classification (RFC):** The base functionality of the algorithm is similar to the regression case. The final decision to which class the observation is assigned to is conducted via averaging the probabilistic prediction of all classifiers.²
- **Principal Component Analysis (PCA):** Jolliffe [Jol11] describes PCA as probably the most common technique used for dimensionality reduction. PCA reduces dimensionality by finding linear combinations, called principal components, which have maximum variance for the data. Furthermore, the components are uncorrelated in order to minimize loss in information and maximize the dimensionality reduction. Jolliffe [Jol11] provides a detailed mathematical explanation.

2.2 Sentiment Analysis

Liu [Liu15] defines sentiment analysis as follows: *"Sentiment analysis, also called opinion mining, is the field of study that analyzes people's opinions, sentiments, appraisals,*

²<http://scikit-learn.org/stable/modules/ensemble.html#forest>, last accessed 23.10.2018

attitudes, and emotions toward entities and their attributes expressed in written text. The entities can be products, services, organizations, individuals, events, issues, or topics. [...] In a nutshell, sentiment analysis or opinion mining aims to identify positive and negative opinions or sentiments expressed or implied in text and also the targets of these opinions or sentiments [...]."

Since the primarily focus of sentiment analysis is on analyzing textual data it has been a research area of natural language processing (NLP). Sentiment analysis has also been widely studied in the context of data mining, web mining, and information retrieval. Due to the increasing importance of sentiment analysis to businesses and society, the research area has spread from computer science to management and social science [Liu15].

In the following we provide an overview of the broad research area of sentiment analysis and opinion mining. First, the development and some applications are presented. We further explain common approaches for sentiment analysis, different levels of analysis and the processing of text prior to sentiment analysis. At the end of this section we highlight common challenges in sentiment analysis to show the complexity of this research area.

2.2.1 Development and Applications

The increasing amount of data and more importantly the availability of data in digital form led to an increase in research on sentiment analysis. Opinions of individuals are widely spread among blogs, reviews, social media, comments, and every form of participation or engagement in a topic. Sentiment analysis not only focuses on the attitude which a text transmits but also on the individual who wrote this text. Based on the statements in the text a sentiment profile of the author may be established. This information provides various possibilities like product recommendations or the determination of the political orientation of the individual.

Not only businesses and organizations are interested in the opinion about their products and services but also governments have interest in how the public thinks about policies and policy changes. It is no longer necessary to conduct a survey about how users think about a specific topic. This information may already be publicly available on social media or other platforms. Further, there are various research topics which use sentiment analysis. For example during elections, the posts on Twitter provide information about candidates approval. Another popular application is stock market prediction with the aid of public opinions, CEO - letters, board posts, and other publicly available information concerning this domain [Liu15].

2.2.2 Common Sentiment Analysis Approaches

This section is dedicated to introduce common approaches to perform sentiment analysis. In the following we briefly summarize the approaches listed by Cambria et al. [CDBF17]:

- **Knowledge-based technique:** Text is classified according to unambiguous affect words. Basically, a predefined set of words is matched to a text, trying to determine its sentiment. The major weakness of the knowledge-based approach is poor

recognition when linguistic rules are involved. For example, the sentence *"Today was a happy day"* can be correctly classified by a knowledge base while *"Today wasn't a happy day at all"* will likely be incorrectly classified. The accuracy of knowledge-based approaches highly depends on the employed resources. Without a comprehensive knowledge base it is difficult to determine the sentiments associated with natural language.

- **Statistical methods:** Statistical methods intend to learn grammatical constructs and word co-occurrence frequencies. Techniques like deep learning use large training corpora of affectively annotated text to learn the valence of affect keywords and word co-occurrence frequencies. Due to the amount of data statistical methods need they are only accurate when applied to users' text on the page or paragraph level.
- **Hybrid approaches:** Hybrid approaches exploit both knowledge-based techniques and statistical methods to better grasp conceptual rules and language constructs to accurately classify text. These approaches try to combine the advantages of recognizing semantics of knowledge-based approaches and the ability to recognize linguistic patterns of statistical methods in order to better understand the sentiments transmitted.

Bing Liu [Liu15] distinguishes between unsupervised and supervised methods for sentiment analysis. The summarized descriptions of both approaches are listed below.

- **Supervised Classification:** Sentiment classification is a text classification problem. Researchers applied machine learning algorithms to classify text. The key for accurate classification are the features. The most common used features are terms and their frequency. These features are individual words and their frequency counted in the text. Furthermore, the part of speech (POS) is considered as another class of feature. For example, adjectives are important expressions of opinions. Phrases which are used in natural language are challenging to automatically capture. There are also words which completely shift the sentiment of a text, the classical example is *not*, or *don't*. Several studies use different combinations of these features for their machine learning algorithms.
- **Unsupervised Classification:** The unsupervised approach uses phrases and sentiment words to classify text. Syntactic patterns are analyzed in the text and interpreted. Therefore, a sentiment lexicon containing words and phrases which are relatable to the investigated sentiments is used. A sentiment score is then computed for each document according to the expressions found in the text and the scores which were assigned to each expression. The simplest form to compute the score for a document is to add up all scores for the expression and check whether this score is positive or negative. Of course, there are many variations to this approach, not only in calculation of the overall score but also the score of each expression in the text.

2.2.3 Levels of Analysis

In this section we highlight the different textual levels of analysis. The different levels are described by [Liu15].

Document level

The document level tries to classify whether a document as a whole expresses a positive or negative sentiment. This analysis implies that each document describes the opinion of one entity, for example a product review describing the opinion about one product. Further, there exists one opinion holder, for example one person which wrote the product review. These assumptions are necessary, multiple opinion holders in one document may have different opinions about the underlying entity it describes. This example already describes a limitation of document analysis, for example if a text contains sentiments describing multiple entities a more fine-grained analysis is needed. The foundation for this kind of analysis is described in the sentence level and the aspect level.

Sentence level

The previous example described the need of a more detailed analysis of a document. This level analyzes each sentence in a document. The general goal is the same, classifying sentences into sentiments. Each sentence can be seen as a very short document. The challenge is the reduced amount of information. Therefore, sentences may also contain no specific sentiment at all and need to be seen as neutral, for example facts. Furthermore, the difference between subjectivity and objectivity has to be considered in the analysis. Subjective sentences express sentiments of an opinion holder whereas objective sentences are expressed by someone but display a general statement, this statement although can have some implicit sentiments. For example, the sentence *"The earphone broke in two days."* expresses a fact but implicitly describes an unsatisfying condition.

The sentence level analysis assumes there is one sentiment per sentence. The target of the sentiment is not considered, which displays some restrictions in this approach. For example, the sentence *"The picture quality of this camera is amazing and so is the battery life, but the viewfinder is a little small for such a great camera"* expresses positive sentiments for the picture quality and the battery life whereas the viewfinder is viewed as too small. Sentences may contain different sentiments for different targets. The aspect based level tries to correctly classify opinions in a text and further determine also the target of these sentiments.

Aspect level

The aspect level further tries to extract the target of the sentiment expressed by the holder of the opinion. Or to whom/what does the expressed sentiment belong? For example, it does not make sense to declare the sentence *"Apple is doing very well in this poor economy."* as either positive or negative. This sentence contains two contrary sentiments. The aspect level focuses on extracting the sentiment and matching it to

the right target. In this example *doing very well* to *Apple* and *poor* to *economy*. This example shows the complexity of sentiment analysis.

2.2.4 Text Processing

Preprocessing text is an essential task prior to performing sentiment analysis. Pröllochs et al. [PFN15] list the following steps:

1. **Cleaning:** Punctuation marks and other symbols are removed from text.
2. **Tokenization:** Text is mostly stored as a string. The analysis of the texts is simpler if the text is stored as single words. Each sentence is split into its consistent words. These single words are referred to as tokens.
3. **Removal of stop words:** Stop words are words without a deeper meaning, such as the three articles of German *der*, *die*, *das*. Due to their missing sentiment and irrelevance in the context of sentiment analysis they can be omitted.
4. **Stemming:** Stemming describes the process of reducing an inflected word to its stem. The goal is to capture the meaning of inflected words even if the stem is not a valid root form.

2.2.5 Challenges in Sentiment Analysis

Language nuances, like language-specific phrases and the use of grammatical constructs, which may seem easy in daily use, intensify the complexity an algorithm faces to perform sentiment analysis. This section briefly describes further challenges in analyzing language constructs.

Negations and even double negations completely change the sentiment of a text. For example the sentence "*Das Produkt ist nicht unbrauchbar.*" uses the word *nicht* to shift the sentiment, initially negative, of the word *unbrauchbar* to an overall positive sentiment of the sentence.

Further, sarcasm and ironic expressions, which can be easily understood by humans, are a complex task for an algorithm to identify. Imagine a product review of a customer. The customer describes the incompetence of the service team and writes the following sentence: "*Bei der hohen Qualität des Kundenservices muss das Produkt mindestens genau so gut sein.*" Humans easily understand the customer claims the product has to be equally bad as the service provided by the company. For an algorithm this is not trivial to understand.

There are many texts which require sensitivity to capture essential and important parts. Some expressions are not as relevant as others which reveal themselves only by reading and understanding the full text.

Combining all the mentioned language-specific phrases in a single text requires the algorithm to not only correctly identify these constructs but also to interpret them accordingly in order to accurately classify expressions.

Related Work

Chapter 3 provides an overview about related work and contributions of researchers in the area of sentiment analysis, focusing on the financial domain. The first section, Section 3.1, summarizes the outcomes of other research in this area. In Section 3.2 we list German resources to perform sentiment analysis.

3.1 Sentiment Analysis in the Financial Domain

Understanding the market dynamics, which tend to explain the behavior of markets, has gained significant interest in research and industry. Newspapers and other forms of media transmit information to consumers. Consumers are often not aware of how these underlying sentiments of news may affect their behavior. Tetlock [Tet07] investigates dependencies between media and stock market activities. Negative sentiments of media induce downward pressure on prices. Furthermore, pessimism leads to temporarily high trading volume. The temporary decrease in returns, which can be predicted by pessimism, are reversed within a few days.

Schumaker and Chen [SC09b] distinguish two different trading philosophies, the fundamentalists and the technicians. The first group tries to determine the price of a stock from financial numbers of the overall economy or the specific sector in which the stock is traded. Contrary supporters of the technician philosophy try to derive the price from historical data. They believe arbitrage opportunities can be found by sensitively analyzing the historical data and volume movements. Furthermore, Schumaker and Chen compare different approaches to predict prices from financial news articles. They claim the most common technique used to classify articles is the bag of words approach. This technique simply searches for the occurrences of words in the article and assigns a sentiment to it. In their approach, grammar or specific word combinations, which combined establish a different meaning, are ignored.

The importance of grammar and phrases are addressed by Chan and Chong [CC17].

They recognize phrases in news articles by a machine learning technique. They conclude the sentiments expressed through text streams are helpful to analyze trends in a stock market index.

Schumaker et al. [SZHC12] try to answer if positive and negative subjectivity influence the prediction of stock prices with news articles. Especially negative subjectivity of news articles have an impact on trading behavior. This conclusion supports the findings of Akhtar et al. [AFOS13]: *"We document asymmetric announcement effects of consumer sentiment news on United States stock and stock futures markets. While a negative market effect occurs upon the release of bad sentiment news, there is no market reaction for the counterpart good news."*

Nopp and Hanbury [NH15] analyzed risk attitudes of banks with sentiment analysis of CEO letters and outlook sections of annual reports. They recommend sentiment analysis for detecting risks to be used in an aggregated form since the evaluation on the individual level led to inaccurate predictions contrary to the aggregated analysis which revealed significant correlations.

While the previous studies focus on the sentiment, extracted from text, Rekabsaz et al. [RLB⁺17] explores the effect of combining factual market data and sentiment scores. They evaluate different fusion techniques to combine factual data and text resources. Further, they predict the volatility in financial markets with factual market data and sentiment scores of annual disclosures of companies.

Several studies address the possibilities of stock market prediction using textual data. Henrique et al. [HSK18] use Support Vector Regression (SVR) to predict stock prices for small and large capitalizations by using daily and up-to-the-minute price information. They claim the model has predictive power especially when using a strategy to periodically update the model. Their findings also indicate an increased accuracy of their predictions during lower volatility periods.

Atkins et al. [ANG18] used two stock indices and two equities to empirically show that information extracted from textual news sources is better at predicting the volatility direction of the market than the price movement. *"Though the inability to predict close price movement any better than random contradicts previously published results, our results suggest that information in news, influencing markets via sentiment-driven behavior essentially affects second order statistics of the financial system."* Regarding second order statistics of the financial system they mention volatility and trade volumes. For their analysis they extracted news articles describing the market, politics, business and world news. The extracted articles are then preprocessed. Further, the semantics of the articles are modeled using topic models. These topic models and the volatility movements were used to fit a Naive Bayes prediction model. They managed to achieve an average directional prediction accuracy for volatility for new data of 56% while the prediction of the asset close price performed no better results than random classification.

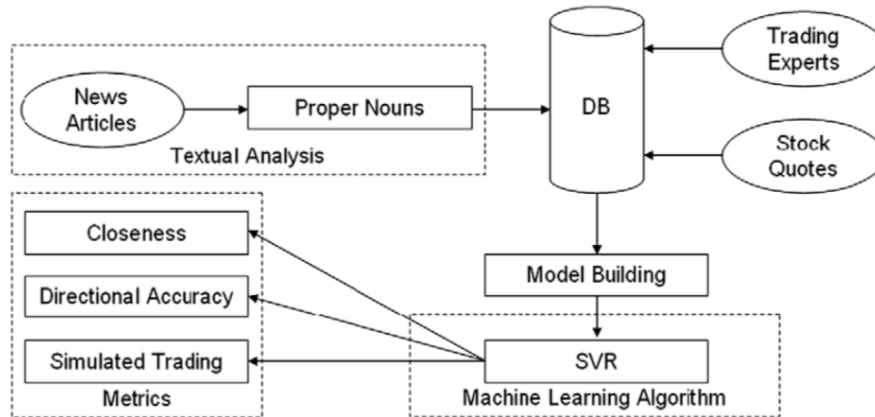


Figure 3.1: System Design of Schumaker and Chen [SC09a]

Schumaker and Chen [SC09a] provide an illustration, Figure 3.1, regarding the functionality of the approach they used to predict stock prices. They extract news articles from Yahoo! Finance and represent them using Proper Nouns. Proper Nouns retain only nouns and noun phrases, which fall in defined categories, from the news articles. Further, minutely stock price data and information regarding buy/sell recommendation from trading experts are gathered and stored. They built and tested different models according to the Global Industry Classification Standard¹ (GICS) distinguishing between sectors, industry groups, industries, and sub-industries. The models are built with only the relevant information regarding the classifications. They use the SVR algorithm with a binary representation of the proper nouns, either being present in the financial news articles or not, for their models. Training the model based on sectors resulted in a score of directional accuracy of 71.18%. Directional accuracy measured if the prediction of the price (up/down) matched the actual price movement. Further, they simulated trading and achieved 8.5% return on invest. For our approach we perform a sentiment analysis of the newspaper articles and represent each article as a vector of sentiment weights. We build a model and apply the SVR algorithms to these weights to predict price and volatility movements. As a metrics we compute the R2 value and the classification accuracy, the percentage of correctly classified observations, in order to compare the results of various models.

Luss and D'Aspremont [LD15] predict the direction of returns and abnormal returns, absolute returns greater than a predefined threshold, of intraday price movements using SVM classification. They apply a bag-of-words approach to press releases and weight the word occurrences with the TFIDF weighting scheme. Although they are unable to predict the direction of returns, abnormal returns appear to be predictable using absolute returns or textual data.

¹<https://www.msci.com/gics>, last accessed 09.09.2018

3.2 Sentiment Lexicons

This section lists common German lexicons in the literature. Each of them is briefly described.

- **SentiWS:** The SentiWS [RQH10] lexicon is a German language resource for sentiment analysis provided by the University of Leipzig. They provide this lexicon as part of their project "Deutscher Wortschatz".² This project provides information about the German language since 1990. They gather and edit data from publicly accessible documents and combine them to a corpora collection. These collections are among the most comprehensive publicly available data in German.³ The SentiWS lexicon consists of approximately 1650 positive and 1800 negative word stems. Furthermore, the lexicon contains different derivations of each word resulting in approximately 16000 positive and 18000 negative words.⁴
- **German Emotion Lexicon:** Klinger et al. [KSR16] create German emotion lexicons to automatically detect emotion in literary studies. Currently they provide lexicons for seven emotions.⁵
- **SePL:** Sentiment Phrases List (SePL) is a German list of phrases which are associated with a sentiment in German. The list is provided by the institute of information systems from the university of Hof and contains over 14000 entries.⁶
- **GermanPolarityClues:** is a German sentiment lexicon created by Waltinger [Wal10] by combining several lexicons. The lexicon has over 10000 polarity features associated to three polarity scores. The scores provide a direction, positive, negative or neutral, to the terms.⁷
- **OpenThesaurus:** A platform which provides synonyms and associations for German words. Further, they provide wordlists for specific domains like geography, biology, physics, politics, economy and more.⁸
- **BPW Dictionary:** Bannier et al. [BPW19] provide a German dictionary for business communication. The dictionary is based on a commonly used English dictionary to examine finance and accounting specific text. The BPW consists of more than 10000 words and is publicly accessible on the website of University of

²<http://wortschatz.informatik.uni-leipzig.de/de>, last accessed 23.07.2018

³<http://wortschatz.informatik.uni-leipzig.de/de/documentation>, last accessed 23.07.2018

⁴<http://wortschatz.informatik.uni-leipzig.de/de/download>, last accessed 21.03.2019

⁵<https://bitbucket.org/rklinger/german-emotion-dictionary/src>, last accessed 06.02.2019

⁶<http://www.opinion-mining.org/SePL-Sentiment-Phrase-List>, last accessed 06.02.2019

⁷<http://www.ulliwaltinger.de/sentiment/>, last accessed 06.02.2019

⁸<https://www.openthesaurus.de>, last accessed 06.02.2019

Giessen.⁹ The lexicon was published after conducting our experiments therefore we do not use it in this thesis. We include it in this list to provide comprehensive sources to perform sentiment analysis for German text.

⁹https://www.uni-giessen.de/fbz/fb02/forschung/research-networks/bsfa/textual_analysis, last accessed 21.04.2019

German Financial Sentiment Lexicon

In Chapter 4 we explain the process of creating the German Financial Sentiment Lexicon, referred to as GFSL. In Section 4.1 we explain our motivation to create the GFSL. Section 4.2 describes the sources to build the lexicon. In the last section, Section 4.3, we provide information about the creation of the GFSL.

4.1 Motivation

We believe a domain specific lexicon enables a more accurate extraction of sentiments, considering the financial domain, communicated in the newspaper articles. Domain specific expressions, which are not considered in a general sentiment lexicon, allow us to extract the sentiments more precisely. This increase in quality of the obtained information may provide the possibility to gain more accurate predictions, assuming there exists a correlation between the sentiments in the newspaper articles and the changes of the ATX. Although we found several resources to perform sentiment analysis for German, as listed in the **Related Work** chapter, we identified a lack of research in this area because there is no German sentiment lexicon for the financial domain available. Motivated by this research opportunity we have decided to create the GFSL.

4.2 Used Resources

The basis for the GFSL is taken from the SentiWS lexicon. Remus et al. [RQH10] created this lexicon. SentiWS lists positive and negative sentiment bearing words, their part of speech tag, the positive and negative polarity of the word in an interval of -1 and 1, and some of the inflections of the word. It contains adjectives, adverbs, nouns and verbs

which express a positive or negative sentiment. The SentiWS lexicon is built with three major resources:

- General Inquirer: Remus et al. [RQH10] use the categories *Pos* and *Neg* and translate them via Google Translator into German. These categories provide the basis of the lexicon. Additionally, they manually added some words of the finance domain since the lexicon was initially developed for sentiment analysis of financial blog and newspaper articles.
- Co-occurrence Analysis: The second source is a co-occurrence analysis of rated product reviews. The authors identified words which occur significantly often and manually chose the 200 most significant ones.
- German Collocation Dictionary: The third source is the German Collocation Dictionary. The words extracted by the two previous sources were used to distinguish between semantic groups. The German Collocation Dictionary contains semantic groups from which additional words are extracted.

The latest version of the SentiWS lexicon consists of approximately 1650 positive and 1800 negative word stems. Including the inflections of the words, the SentiWS lexicon consists of 16000 positive and 18000 negative word forms.¹

4.3 Creating GFSL

The initial SentiWS lexicon contains words which are uncommon in the financial and economic domain. These words are used to capture the general sentiment of text. For the GFSL we manually extract approximately 1500 basic word forms, without the inflections, related to sentiments in the financial domain by our personal judgement. Further, the polarity of the words is removed since the lexicon is adapted for a different purpose and therefore of limited use.

We extend the GFSL by context specific dictionaries from the General Inquirer Categories.² The General Inquirer is a computer-assisted approach for content analysis. It is basically a mapping tool which categorizes texts into domains according to predefined dictionaries which contain common words for the specific domain.³ The words of the following categories are translated and added to the GFSL. The dictionaries can be found at the General Inquirer homepage.²

- Econ@: containing words of financial, economic and industrial context.
- Exch: words concerning buying, selling and trading.

¹<http://wortschatz.informatik.uni-leipzig.de/de/download>, last accessed 21.03.2019

²<http://www.wjh.harvard.edu/~inquirer/homecat.htm>, last accessed 23.07.2018

³<http://www.wjh.harvard.edu/~inquirer/3JMoreInfo.html>, last accessed 23.07.2018

- TrnGain: words describing transaction gains, general words for accomplishment.
- TrnLoss: words describing transaction loss, general words for not accomplishing
- WltPt: words for wealth in business and commerce.
- WltTran: words for pursuit of wealth.

Further, we manually added expressions describing trends of stock market indexes, development of companies and phrases about future developments by reading financial and economic newspapers and news of the Vienna Stock Exchange.⁴ We added approximately 200 words like *Kursgewinn*, *Kursverlust* and *insolvent* to the GFSL.

The GFSL contains only stemmed word forms because in German words may have many inflections. Including and finding all inflections of the words in the GFSL is a major effort. We decided to stem the words in the newspaper articles as well in order to compare the words in the GFSL to the words in the articles. To stem the words we use the German stemmer in the NLTK python library.⁵ The NLTK - package uses the Snowball stemming algorithm for German language.⁶ The stemmed words are further checked for uniqueness. The final GFSL contains a total of 2267 positive and negative stemmed words describing positive and negative sentiments in the financial and economic domain. The GFSL is attached at the end of the thesis.

⁴<https://www.wienerborse.at/news/>, last accessed 11.10.2018

⁵<https://www.nltk.org>, last accessed 18.08.2018

⁶<http://snowballstem.org>, last accessed 18.08.2018

Sentiment Analysis of the Newspaper Articles

In Chapter 5 we explain our approach to perform the sentiment analysis of the newspaper articles. Section 5.1 describes the extraction of the newspaper articles. In Section 5.2 we explain the process of preprocessing and weighting the newspaper articles. Section 5.3 lists several label vectors we create to perform the machine learning. In the last part of Chapter 5, Section 5.4, we introduce our experiment design.

5.1 Data Accumulation

The news articles are gathered from *derStandard*.¹ It is an Austrian newspaper and has all of its articles available in an online archive.² The archive consists of articles since 2007. This online archive is the source for the data. The structure of the archive provides a possibility to navigate through the articles.

[https://derstandard.at/archiv/\(year\)/\(month\)/\(day\)](https://derstandard.at/archiv/(year)/(month)/(day))

The parenthesized expressions, e.g. year, month and day, display the structure of the URL used to navigate through the archive. By inserting the defined values a webpage containing a list of links, referring to every single article which was published at this specific date, is provided. The following section explains the algorithm used to extract the relevant articles from the archive and the challenges which occurred throughout the process.

¹<https://derstandard.at>, last accessed 23.07.2018

²<https://derstandard.at/archiv>, last accessed 23.07.2018

5.1.1 Web Scraping

Web scraping is the automated process of gathering data from the internet. This is mostly accomplished by writing a program which accesses web pages and extracts the desired information [Mit15].

The following pseudocode explains the functionality of the algorithm used to retrieve the articles.

Algorithm 1 Web Scraper

```
1: startDate  $\leftarrow$  2007/01/01
2: endDate  $\leftarrow$  2017/12/31
3: dates  $\leftarrow$  getDates(startDate, endDate)
4: for each date in dates do
5:   html  $\leftarrow$  getHTML("https : //derstandard.at/archiv/" + date.year + "/" +
   date.month + "/" + date.day)
6:   links  $\leftarrow$  extractAllLinks(html)
7:   for each link in links do
8:     article  $\leftarrow$  getHTML(link)
9:     if validateArticle(article) then
10:       headline  $\leftarrow$  getHeadline(article)
11:       text  $\leftarrow$  getText(article)
12:       storeArticle(date, headline, text)
13:     end if
14:   end for
15: end for
```

The web scraper accesses every date in the archive which lies in between the *startDate* and *endDate*. The function *getDates()* returns all relevant dates, including the starting date and the last date of which the articles should be extracted. The URL for every date is accessed and the HTML code is saved in *html*. This HTML code contains a list of links to the articles which were published at this date. All these links are extracted to further iterate through them. Each link is accessed and the HTML of this page is stored in *article*. Further, the function *validateArticle()* checks if the article is listed in the "Wirtschaft" - section. The headline and the actual text of the article are extracted in line 10 and 11. The date, headline and the text are then stored in a database.

The extraction of the specific parts of the HTML document is accomplished by the BeautifulSoup HTML parser.³ The parser allows to navigate through the Document Object Model (DOM) of the HTML document.

³<https://www.crummy.com/software/BeautifulSoup/bs4/doc/>, last accessed 24.03.2019

ATX Data	
Date	Closing Price
02.01.2007	4558.96
03.01.2007	4565.91
04.01.2007	4515.17
05.01.2007	4421.55
08.01.2007	4381.50
09.01.2007	4405.16

Table 5.1: Structure of ATX-Data

5.1.2 ATX-Quotation Data

The data describing the historical development of the ATX can be downloaded from the Vienna stock exchange website.⁴ Daily measurements about the opening price, the daily maximum, the daily minimum, the closing price and the development in percentage are provided. For the purpose of this thesis the date and the closing price are extracted. The time series of the ATX consist of missing dates, as noticeable in Table 5.1 . Trading opens on Monday and closes on Friday. Further, Austrian holidays are days without trading. The time series, starting from 01.01.2007 until 31.12.2017, consists of 2728 observations where trading took place.

5.2 Sentiment Analysis

This sections explains how the news articles are further processed and analyzed. Furthermore, we describe the preprocessing of the articles and illustrate this process with an example.

5.2.1 Approach

So far the text and the date of issue of the news articles are stored in the database. For the further analysis we group the articles by their publication date. For each day we appended the newspaper articles to retrieve a time series of articles since the data of the ATX is as well measured on a daily basis. Trading on the stock market exchange does not take place on weekends and Austrian holidays. Therefore, only articles published on a date on which trading took place are further considered in order to utilize them to build a model to describe the ATX observations.

Prior to perform the sentiment analysis we preprocess the articles according to the process described by Pröllochs et al. [PFN15], explained in the **Background** chapter. The actual sentiment analysis is accomplished via the lexicon-based approach. The lexicon is

⁴https://www.wienerborse.at/indizes/aktuelle-indexwerte/historische-daten/?ID_NOTATION=92866&ISIN=AT0000999982, last accessed 16.08.2018

Sentiment weights					
Date	$Word_1$	$Word_2$...	$Word_{n-1}$	$Word_n$
d_1	w_1	w_2	...	w_{n-1}	w_n
\vdots	\vdots	\vdots	\vdots	\vdots	\vdots
d_t	w_1	w_2	...	w_{n-1}	w_n

Figure 5.1: Post Weighting Structure of the Newspaper Articles

used to capture the sentiments in the preprocessed articles. The lexicon as well as the articles consist of stemmed words. With the GFSL lexicon the articles, published on a trading day, are analyzed. This means for every trading date there exist an observation consisting of the weighted sentiment words of the lexicon. How we perform the weighting of the words is explained in the **Weighting** section. Figure 5.1 shows the weights and the representation of the news articles at a specific date. The words $Word_i$ represent the words in the lexicon. The weights w_i represent the weights of $Word_i$ of the news articles at the specific date d_t . The weights only represent the score computed according to the weighting method. They do not allow for conclusions on the sentiments communicated at this day. They are used to build a supervised machine learning model.

5.2.2 Sentiment Analysis Algorithm

The sentiment analysis of the articles is split into two parts. First, the articles are preprocessed as described in the **Background** chapter. In the second part the weights are computed, according to the weighting schemas, and stored. The following two sections explain these two steps. The functionality of the algorithms to perform these steps is displayed in an abstract form and explained.

Preprocessing

First the dates on which trading took place are extracted by reading the file which contains observations of the ATX-course. The next step is to load the database with the scraped newspaper articles and set up a new database to further store the preprocessed articles. We load the stopwords for the German language. The stopwords and the stemmer are received from the Natural Language Toolkit (NLTK).⁵ The NLTK - package uses the Snowball stemming algorithm for German language.⁶ For every date, of the previously loaded dates, the newspaper articles are extracted. The articles published on the same date are appended and returned as text. This text is then tokenized with the *WordPunctTokenizer* from the NLTK package. It splits the text into words and punctuations. The following example illustrates the functionality of the *WordPunctTokenizer*.

⁵<https://www.nltk.org>, last accessed 18.08.2018

⁶<http://snowballstem.org>, last accessed 18.08.2018

Algorithm 2 Preprocessing of the Articles

```

1: dates  $\leftarrow$  getDates()
2: db  $\leftarrow$  loadDB()
3: db_processed  $\leftarrow$  setUpDB()
4: stopwords  $\leftarrow$  stopwords.words("german")
5: for each date in dates do
6:   text  $\leftarrow$  getArticles(date, db)
7:   tokens  $\leftarrow$  WordPunctTokenizer().tokenize(text)
8:   tokens_processed  $\leftarrow$  []
9:   for each word in tokens do
10:    if word.isalpha() AND word not in stopwords then
11:      tokens_processed.append(stem(word))
12:    end if
13:  end for
14:  writeProcessedTokens(tokens_processed, date, db_processed)
15: end for

```

Input:

"An der Wiener Börse stieg der wichtigste österreichische Aktienindex am späten Freitagvormittag um über 1,3 Prozent bis auf 3.003,78 Punkte."

Output:

['An', 'der', 'Wiener', 'Börse', 'stieg', 'der', 'wichtigste', 'österreichische', 'Aktienindex', 'am', 'späten', 'Freitagvormittag', 'um', 'über', '1', ',', '3', 'Prozent', 'bis', 'auf', '3', ',', '003', ',', '78', 'Punkte', '']

The algorithm further iterates through every word in the resulting list of words. List elements which consist of non alphabetic characters and punctuations are removed. The remaining words are stemmed and stored in a database. The continuation of the example displays the outcome after the preprocessing step.

Input:

['An', 'der', 'Wiener', 'Börse', 'stieg', 'der', 'wichtigste', 'österreichische', 'Aktienindex', 'am', 'späten', 'Freitagvormittag', 'um', 'über', '1', ',', '3', 'Prozent', 'bis', 'auf', '3', ',', '003', ',', '78', 'Punkte', '']

Output:

['wien', 'bors', 'stieg', 'wichtig', 'osterreich', 'aktienindex', 'spat', 'freitagvormittag', 'prozent', 'punkt']

Algorithm 3 Weighting of the Articles

```

1: dates  $\leftarrow$  getDates()
2: db  $\leftarrow$  loadDB()
3: sentiments  $\leftarrow$  getSentiments()
4: sentiments_dict  $\leftarrow$  dict.fromkeys(sentiments, 0)
5: df  $\leftarrow$  calculateDF()
6: for each date in dates do
7:   tokens  $\leftarrow$  getTokens(date, db)
8:   wordcount  $\leftarrow$  Counter(tokens)
9:   for each word in wordcount do
10:    if word in sentiments_dict then
11:      sentiments_dict[word] = wordcount[word]
12:    end if
13:  end for
14:  writeRowToCSV(weighting(sentiments_dict, df))
15:  sentiments_dict  $\leftarrow$  dict.fromkeys(sentiments, 0)
16: end for

```

We apply each of the weighting schemes to the features. Like in the preprocessing algorithm the first line loads the dates on which trading took place. In line two the database, where the previously processed articles are stored, is loaded. Further, the sentiment lexicon is loaded. To count the occurrences of the sentiments in the articles we create a dictionary. The dictionary basically consists of two lists, the first list contains the words of the lexicon and the second list the corresponding occurrences. The dictionary is initialized with zeros since some words of the lexicon may not be present at every date. The function *calculateDF*() returns the document frequency of all words in the sentiment lexicon as a dictionary. The document frequency is used for the TCIDF and the TFIDF weighting scheme as explained previous in this section. The algorithm further iterates through all dates and loads the tokens, the list of words resulting from the preprocessing step, for each date. Words in this list are then counted resulting in a dictionary of words and their occurrences. This list contains words which are not in the sentiment lexicon. The words in the lexicon and the tokens are both stemmed from the previous steps enabling the comparison and extraction of only the words which are in the sentiment lexicon.

The occurrences of these words are then stored in the *sentiments_dict*, which was initialized at the beginning with zero values. The dictionary *sentiments_dict* now contains the occurrences of the sentiments of all articles published at the date which is being iterated.

Finally, the sentiment weights, considering the four different weighting schemas, are computed and stored in separate files. The occurrences of the *sentiments_dict* dictionary are then set to zero for the next iteration step.

In this section we established the fundamentals for the machine learning part. The newspaper articles have been cleaned, tokenized, stop words and punctuation have been

removed and each word is reduced to its stem, which makes it possible to compare the words in the newspaper articles to the words in the sentiment lexicon. Further, the weighting process has been performed for the lexicons, resulting in four data sets for each sentiment lexicon. Together with the ATX observations these files provide the foundation for the machine learning part.

5.3 Defining the Label

The feature set consists of the sentiment weights, currently four different feature sets per lexicon according to the different weighting schemes calculated previously: TC, TSIDF, TF and TFIDF. The corresponding label vector consists of the closing price of the ATX. The rows of the feature set correspond to the rows in the label vector. They match trading days of the ATX between 01.01.2007 and 31.12.2017. Throughout the process the label vector will be varied in order to find the combination which provides the most accurate model. We test the described label vectors in this section. The y_i represent the label entry and the p_i the closing price observations of the ATX. We investigate regression and classification scenarios. The feature set, the sentiment weights, always stays the same. Only the label set is varied. We perform the experiments with each weighting schema. In the following we define several labels for the regression and the classification experiments.

Labels for Regression:

- For the first attempt to find a model, the label vector is defined as the difference between the closing price of the current date and the date before.

$$y_i = p_i - p_{i-1}$$

- Another possible option for the label vector is the change in volatility. Negative or positive statements in the articles may cause volatility changes. The measurement of volatility is adapted from Rekabsaz et al. [RLB⁺17].

$$v_{[i,i+t]} = \sqrt{\frac{\sum_{l=i}^{i+t} (p_l - \bar{p})^2}{t}}$$

For t explaining future volatility changes:

$$y_i = v_{[i,i+t]} - v_{[i-t,i]}$$

For t explaining past volatility changes:

$$y_i = v_{[i-t,i]} - v_{[i-2*t,i-t]}$$

The parameter t is varied to investigate different spans of volatility changes.

Labels for Classification:**Directions:**

Atkins et al. [ANG18] predict closing price movement and volatility directions. The following labels test whether the sentiment weights allow to classify directions of different values. The resulting label values are classified as positive or negative depending on the sign.

- Closing price directions:

$$y_i = p_i - p_{i-1}$$

- Return directions, as defined in the paper of Rekabsaz et al. [RLB⁺17] and Atkins et al. [ANG18].

$$y_i = r_i - r_{i-1}, \text{ where } r_i = \ln(p_i) - \ln(p_{i-1})$$

- Volatility directions, using the volatility measurement mentioned in Rekabsaz et al. [RLB⁺17] and previously described:

$$y_i = \ln(v_{[i,i+t]}) - \ln(v_{[i-t,i]})$$

Again the parameter t is varied.

Abnormal values:

Luss and D'Aspremont [LD15] predict abnormal returns. They classify returns exceeding a threshold into abnormal returns and try to predict these cases with news articles. The following classifications are designed similar to their approach.

In this thesis the threshold in order to classify abnormality is the distance from the mean of the observations. The following measurement is used.

$$\text{mean}(v) \pm x * \text{STD}(v)$$

The x defines the range for normality, everything outside of this range is defined as abnormal. Different values of x are investigated. Furthermore, classifications of two and of three groups are tested. The two groups are the normal values and the abnormal values. The three group labels distinguish between abnormal low values, normal values and abnormal high values. The different settings are tested for the closing price movements, the returns and the volatility changes as defined previously.

5.4 Experiment Design

This section describes our experiment design. The settings for every experiment is the same. The regression and classification scenarios only differ by the method of evaluation and the tested algorithms. The single steps are performed for each label and each weighting scheme.

1. **Loading:** The first step is to load the sentiment weights and the labels.
2. **Preprocessing:** The features are normalized using the L2 normalization. Further, dimensionality reduction is applied to reduce the time needed for training. We test an amount of 100, 400 and 800 features for each scenario.
3. **Fitting the Model:** For the regression cases we test support vector regression (SVR) and random forest regression (RFR). For the classification cases we test SVC and RFC. In order to find the best hyper-parameter setting, grid search with cross-validation is applied for each algorithm to automatically find the best model according to predefined parameter ranges. Grid search automatically executes every combination of parameter setting and returns the best performing parameter space according to the cross-validation score.
4. **Testing:** The best performing model from step 3 is then applied on the testing set. The resulting predictions are used in the next step.
5. **Evaluation Metrics:**
 - Regression: In the regression case the R2-value and the corresponding p-value are calculated.
 - Classification: To compare classifications we establish a baseline with a dummy classifier. This classifier always predicts the most common group of the training set. The difference between the accuracy of the baseline and the accuracy of the predictions indicates whether the model is capable to capture relevant information from the sentiment weights.

We conduct the experiments with the RFC and SVC algorithms for the classification cases and the RFR and SVR algorithms for the regression cases. We use all four weighting schemas for each experiment. Further, we perform each experiment with the GFSL and the SentiWS lexicon to evaluate the GFSL.

Results and Discussion

In this chapter we present the results of the conducted experiments. The tables display the most interesting results. Some of the tested labels do not provide any useful information, they are not further investigated. In Section 6.1 we present the results of the regression experiments, followed by the results of the classification scenarios in Section 6.2. Section 6.3 is about the qualitative analysis. We discuss the results in Section 6.4.

The following list explains the notation and abbreviations used in this chapter.

- Label: The label used for the model.
- SVR: The Support Vector Regression algorithm.
- SVC: The Support Vector Classification algorithm.
- GFSL: The created German Financial Sentiment Lexicon is used for the sentiment analysis.
- SWS: The SentiWS lexicon is used for the sentiment analysis.
- RFR: The Random Forest Regression algorithm.
- RFC: The Random Forest Classification algorithm.
- R2: The R2 score of the testing set multiplied by 100. Stars indicate the significance of the regression coefficients to the model. The different levels for the p-value are displayed as follows: p-value < 0.05: *, p-value < 0.01: ** and p-value < 0.001: ***
- Baseline: The baseline is defined as the accuracy of the predictions of a dummy classifier which always predicts the most frequent class of the training set. All other settings are compared to this baseline. The + and – signs indicate the difference between the accuracy of the experiment and the baseline.

SVR		RFR	
GFSL	SWS	GFSL	SWS
2.44***	0.99*	0.77*	0.48

Table 6.1: R2 * 100 Closing Price Difference

		SVR		RFR	
Direction	# of days	GFSL	SWS	GFSL	SWS
Future	3	0.03	0.14	0.91*	0.63
	7	0.32	0.51	0.34	0.76*
	14	0.28	1.08*	0.70*	1.31**
	30	0.14	0.20	0.66	1.18*
	91	0.00	0.19	0.49	1.71**
	182	2.58***	6.30***	1.54**	0.54
	365	9.22***	14.98***	3.09***	2.86***
Past	3	0.07	0.97*	0.28	0.69
	7	0.63	1.40**	0.89*	0.59
	14	0.00	0.21	0.56	1.06*
	30	0.15	0.26	0.94*	1.04*
	91	0.90*	1.97***	0.47	1.70**
	182	2.63***	1.47**	0.78*	0.37
	365	0.01	0.00	0.33	0.56

Table 6.2: R2 * 100 Volatility Changes

- Bold values: In each row of the tables one value is displayed in bold. This setting of algorithm and sentiment lexicon achieved the best results among the other settings.

6.1 Results of the Regression

Table 6.1 displays the performed R2 scores, multiplied by 100, of the two algorithms together with the different lexicons. As observable the sentiment weights contain useful information to predict the development of the closing price. Especially, the case with the SVR algorithm and the financial lexicon performed value of 2.44 with significant influence of the sentiment weights. Both algorithms performed better with the financial lexicon. Support vector regression provide better results in this scenario.

The investigation of the volatility changes, Table 6.2, confirms the support vector algorithm performs better than the random forest algorithm in the regression case. What is further observable, the higher the span the better the results, especially when investigating future volatility changes. This might be due to fewer fluctuations with the higher span. It is also noticeable that the SentiWS lexicon outperforms the financial lexicon in the future significant results. The label, investigating volatility shifts of 365 days into the past, performs surprisingly badly. The label is computed using observations

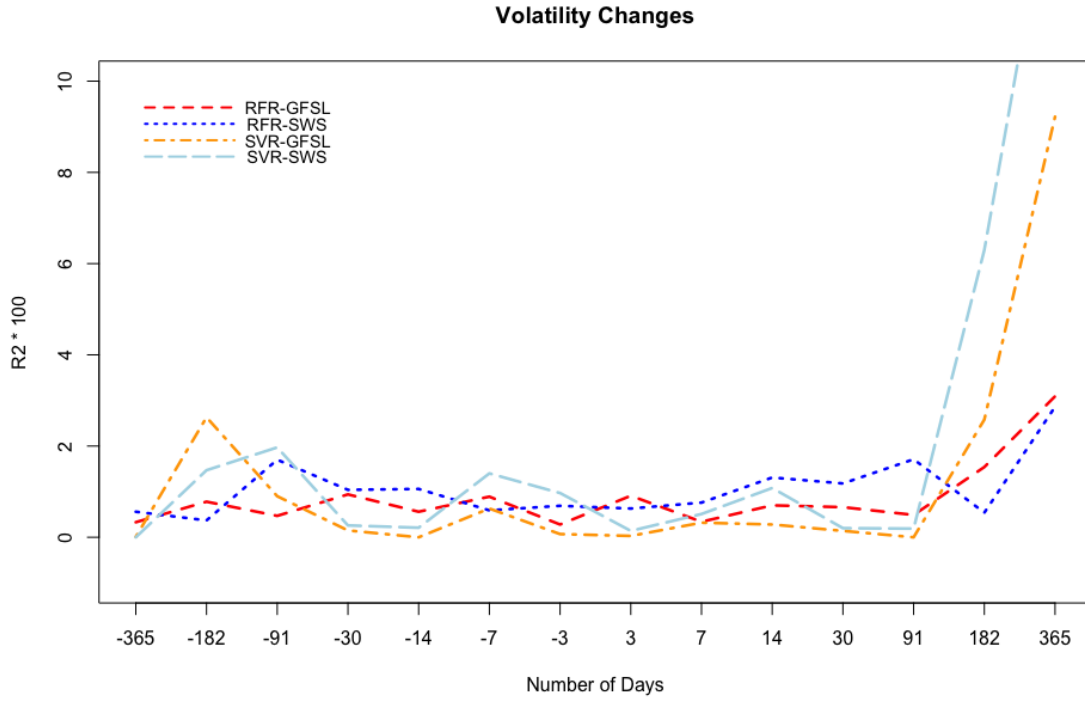


Figure 6.1: $R^2 * 100$ of Volatility Changes based on the results in Table 6.2. The x-axis indicates the span used to compute the volatility.

two years into the past. This means the values are not calculated as intended for the first two years as the data is not available. This explains the fall in performance.

Figure 6.1 illustrates the results of Table 6.2. As mentioned the $R^2 * 100$ value significantly increases with spans higher than 91 days into the future.

	SVC		RFC	
Baseline	GFSL	SWS	GFSL	SWS
56.12%	-2.38	-2.74	+1.28	-0.37

Table 6.3: Accuracy Closing Price Direction

			SVC		RFC	
Direction	# of days	Baseline	GFSL	SWS	GFSL	SWS
Future	3	49.36%	-1.65	+3.47	+5.86	+4.38
	7	51.91%	-0.54	+0.55	+2.93	+3.11
	14	53.38%	-1.47	+2.37	+3.65	+1.09
	30	46.98%	0.00	0.00	+3.11	+1.46
	91	46.25%	0.00	0.00	+5.48	+2.74
	182	42.23%	+4.02	+6.39	+5.85	+6.58
	365	37.84%	+7.86	+12.25	+5.12	+10.97
Past	3	49.36%	+2.58	+3.47	+3.29	+3.10
	7	51.55%	0.00	+2.74	+3.29	+2.92
	14	53.19%	0.00	0.00	+0.74	+0.18
	30	45.70%	0.00	0.00	+2.92	+1.46
	91	60.14%	0.00	0.00	+0.18	+0.18
	182	42.59%	0.00	0.00	0.00	0.00
	365	28.33%	+6.40	+9.93	+17.55	+11.52

Table 6.4: Accuracy Volatility Directions

6.2 Results of the Classification

Predicting direction of the closing price resulted in only one model exceeding the baseline, as seen in Table 6.3. The RFC algorithm works best with the GFSL in this test case achieving an increase of 1.28% percentage points over the baseline. All other tested cases are below the baseline.

In contrast to the regression results the classification of volatility changes performs better for shorter time spans, as observable in Table 6.4. The general SentiWS lexicon performs better with the support vector classification algorithm. The random forest algorithm achieved the best results. The GFSL contains words and phrases which are characteristic for positive and negative economic times. The random forest classifier uses these phrases in order to make specific decision rules. This might explain the good results with the financial lexicon, especially in the shorter time spans. We observe a shift in volatility directions. As the baseline always predicts the most frequent group of the training set we can clearly see there is a shift in the test set because the accuracy is lower than 50%. We observe the wider the span the stronger the shift, except the label explaining the volatility directions 91 days into the past. The algorithms are able to detect this shift and perform better results than the baseline. Still these results are

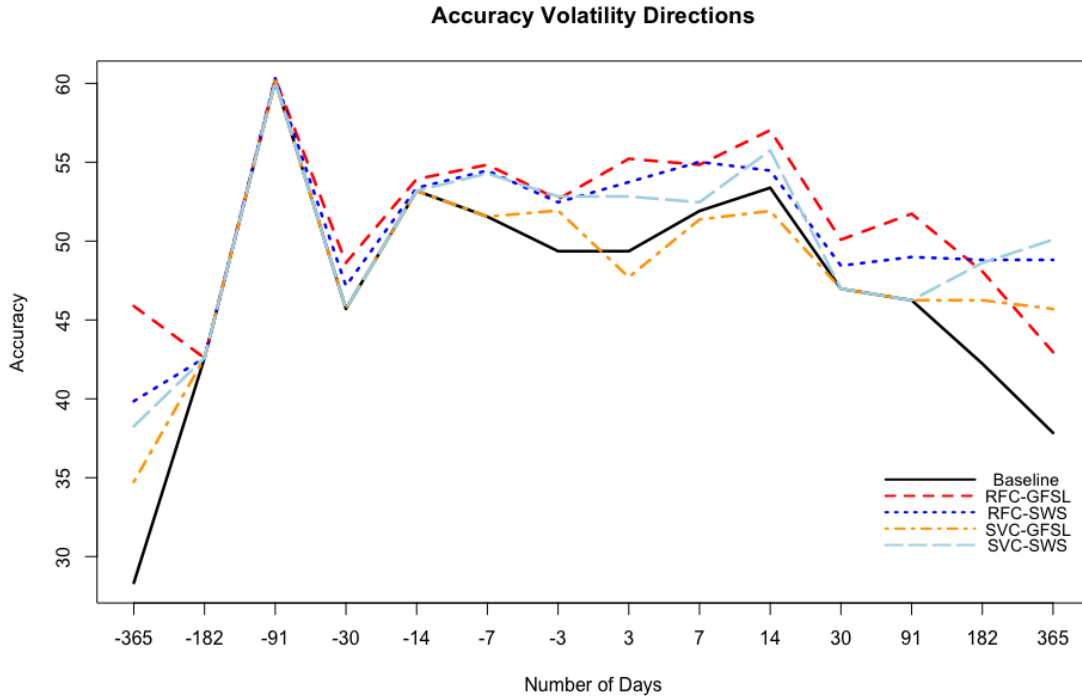


Figure 6.2: Prediction accuracy of Volatility Directions based on the results in Table 6.4. The x-axis indicates the span used to compute the volatility.

often lower than 50% which makes them irrelevant as random classifying would lead to approximately 50% in the long run. The spans up to 14 days into the future and into the past provide satisfying results.

Figure 6.2 illustrates the results presented in Table 6.4. We can see the red line, describing the results of the RFC algorithm with the GFSL, almost has the highest accuracy throughout the spans. As mentioned the predictions work well for short spans into the future and into the past which is observable in the concave trend of all five lines. Table 6.5 describes the results of classifying abnormal values. We investigate different settings for x , regulating the amount of values which are classified as abnormal, ranging from 0.1 up to 2. We achieve the best results with a x value of 0.1 and classification into three groups; abnormal low, normal and abnormal high values. The random forest algorithm, in general, provided better results, despite the cases listed in the table. In our experiments volatility directions with a span of more than 14 days into the past do not provide any useful outcome.

Direction	Algorithm	Label	Baseline	GFSL	SWS
-	SVC	Closing Price	50.09%	+2.01	-2.56
Future	RFC	Volatility 3d.	43.87%	+1.28	+2.38
	RFC	Volatility 7d.	46.62%	+0.54	+2.55
	SVC	Volatility 14d.	45.52%	+0.91	+3.10
	RFC	Volatility 30d.	34.73%	+3.48	+2.56
	RFC	Volatility 91d.	38.21%	+2.37	+3.10
	RFC	Volatility 182d.	36.56%	+3.11	+3.65
	RFC	Volatility 365d.	35.10%	+11.33	+11.04
Past	RFC	Volatility 3d.	43.69%	+1.46	+4.75
	RFC	Volatility 7d.	46.25%	+1.46	+3.11
	RFC	Volatility 14d.	45.52%	+2.93	+0.91
	RFC	Volatility 30d.	75.13%	0	0
	SVC	Volatility 91d.	70.38%	0	0
	RVC	Volatility 182d.	34.00%	0	0
	RFC	Volatility 165d.	26.14%	0	0

Table 6.5: Accuracy Abnormal Values, with 3 Classes



Figure 6.3: Feature Importance GFSL



Figure 6.4: Feature Importance SWS

6.3 Qualitative Analysis

Figure 6.3 shows the most important features of the classification of the closing price directions using the GFSL with the random forest classifier. Figure 6.4 displays the feature importance with the same setting but using the SWS lexicon for the sentiment analysis. The most important features of the RFC are displayed as the model allows to extract the most important features. The size of the words indicate their importance.

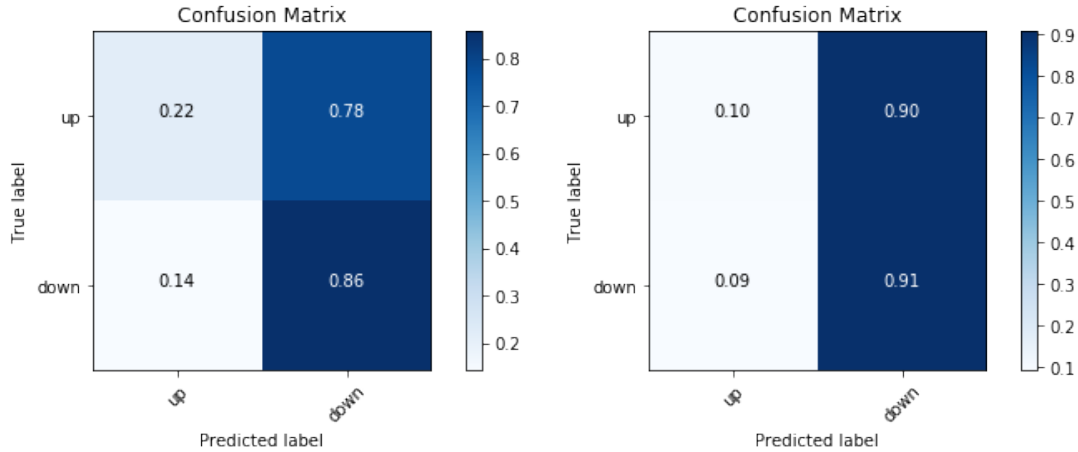


Figure 6.5: Closing Price Difference GFSL Figure 6.6: Closing Price Difference SWS

The bigger the words the higher their importance in the random forest model. Although the GFSL contains positive and negative associated words, with respect to the economic and financial domain, the words in Figure 6.3, which are immediately recognized, are negatively associated when thinking about a stock market index or an economy. For example "fiel", "verlor" and "tief". The feature importance of the words, derived from the model using the SWS lexicon, Figure 6.4, are more general, which is not surprising as it is the intention of the lexicon to provide a general corpora of German sentiment words.

The confusion matrix of the predicted labels of the RFC algorithm and the true labels, Figure 6.5 for the GFSL and Figure 6.6 for the SWS lexicon, show the predictions are heavily biased towards the negative side, predicting a negative movement of the ATX, for both lexicons. These observations are similar to the ones by Akhtar et al. [AFOS13]. As explained in their work, they find no market effect upon the release of good news while negative news have an impact on the stock market.

The algorithm is able to connect negative sentiments to decreasing movements on the stock market but unable to connect positive sentiments in the newspaper articles to upward movements on the stock market. Further, Schumaker et al. [SZHC12] claim that especially negative subjectivity of news articles have an impact on trading behavior.

We expected the combination of the two lexicons to combine the strengths of the GFSL in some cases and the strengths of the SentiWS in other test cases. This is not the case. Combining the lexicons leads on the one hand to improvements in some experiments but on the other hand to worse results than the other lexicons. There is no clear pattern observable.

6.4 Discussion

This list summarizes the results and main outcomes of this thesis. The following insights are gained:

- In general the sentiments extracted from the newspaper articles contain useful information to build regression and classification models to predict closing price changes, closing price movements, changes in volatility and volatility movements.
- The quality of the predictions varies with the algorithm and lexicon used for the analysis. The sentiments alone do not allow for qualitative predictions. However, they contain useful information which may be used in more advanced models especially when it comes to the prediction of volatility and changes in volatility.
- In general the support vector algorithm performs better when it comes to the regression scenarios while the random forest algorithm provides better results when classifying volatility changes and abnormal values.
- The random forest algorithm performs better with the GFSL in the classification of volatility changes with a small span.
- Regression models, describing volatility changes, are more accurate with higher spans, contrary to classifying volatility changes which work better with a lower span.
- Although the created GFSL does not clearly outperform the SentiWS lexicon, it has advantages in specific test cases. In our experiments the GFSL outperformed the SentiWS in the classification of closing price directions and the prediction of closing price changes.
- The outcomes are highly influenced by negativity in the sentiments. The tested models are able to relate negative sentiments to a decrease of the ATX, positive market behavior is however hardly relatable to positive sentiments.

Conclusion and Future Work

Sentiment analysis of financial news, statements, and text, describing the financial or economic domain, has gained increasing interest of researchers and practitioners.

This thesis focuses on analyzing newspaper articles of an Austrian daily newspaper as a source to predict the development of the Austrian Traded Index (ATX). We recognized a lack of resources to perform sentiment analysis of German financial texts. We therefore contribute to this research area by creating a fin.-specific lexicon, referred to as German Financial Sentiment Lexicon (GFSL), insights in the extraction of daily newspaper articles, and extensive analysis on the possibilities in explaining the behavior of the ATX by the sentiment analysis of news articles.

Predicting the ATX with newspaper articles is a challenging task. We conduct several experiments to study the correlation between the sentiments, extracted from the articles, and the ATX. We train supervised machine learning models to predict closing price changes, volatility changes, closing price movements, volatility movements and abnormal values. We evaluate the models with cross-validation techniques and baselines.

Our experiments show the newspaper articles do contain useful information, which positively correlate with the development of the ATX. Further, we observe the GFSL does not clearly outperform the general SentiWS lexicon, although it has advantages in predicting closing price changes and classifying closing price movements. The outcomes of our experiments are highly influenced by negativity in the sentiments. The tested models are able to relate negativity in the sentiments to negative market behavior while positivity is hardly relatable to positive market behavior confirming the results of other studies.

For future work we consider the following directions:

- The GFSL is one point to improve the results. In our experiments we show the potential of the GFSL. Creating a domain specific sentiment lexicon requires expertise in linguistics. Extending the GFSL with more domain specific words,

expressions and language constructs will further increase the quality of future outcomes.

- We recommend to test more machine learning algorithms. An interesting direction is trying Recurrent Neural Networks (RNN), as it considers the development of the sentiments over time. To include a history of sentiments into a model enables to identify phases of positive and negative market behavior.
- This analysis focuses on the articles of one newspaper. We see potential in combining several sources, for example news published by companies listed on the stock exchange as well as using the contents of several online newspapers. Several studies combine sentiment analysis and factual market data to build prediction models. Other economic influences add useful information to such a model and the combination of textual and factual market data enables new possibilities.

List of Figures

2.1	Supervised Learning	6
2.2	Machine Learning Process [Ras15]	8
2.3	Linear SVR, Smola and Schölkopf [SS04]	11
2.4	SVC Example [HCL03]	12
3.1	System Design of Schumaker and Chen [SC09a]	19
5.1	Post Weighting Structure of the Newspaper Articles	30
5.2	Wordcloud Feature Occurrences GFSL	32
6.1	R2 * 100 of Volatility Changes based on the results in Table 6.2. The x-axis indicates the span used to compute the volatility.	39
6.2	Prediction accuracy of Volatility Directions based on the results in Table 6.4. The x-axis indicates the span used to compute the volatility.	41
6.3	Feature Importance GFSL	42
6.4	Feature Importance SWS	42
6.5	Closing Price Difference GFSL	43
6.6	Closing Price Difference SWS	43

List of Tables

5.1	Structure of ATX-Data	29
6.1	$R^2 * 100$ Closing Price Difference	38
6.2	$R^2 * 100$ Volatility Changes	38
6.3	Accuracy Closing Price Direction	40
6.4	Accuracy Volatility Directions	40
6.5	Accuracy Abnormal Values, with 3 Classes	42

List of Algorithms

1	Web Scraper	28
2	Preprocessing of the Articles	31
3	Weighting of the Articles	33

Bibliography

- [AFOS13] Shumi Akhtar, Robert Faff, Barry Oliver, and Avanidhar Subrahmanyam. Stock salience and the asymmetric market effect of consumer sentiment news. *Journal of Banking and Finance*, 37:4488—4500, 2013.
- [ANG18] Adam Atkins, Mahesan Niranjan, and Enrico Gerding. Financial news predicts stock market volatility better than close price. *The Journal of Finance and Data Science*, 4(2):120 – 137, 2018. ISSN: 2405-9188.
- [BPW19] Christina Bannier, Thomas Pauls, and Andreas Walter. Content analysis of business communication: introducing a german dictionary. *Journal of Business Economics*, 89(1):79–123, Feb 2019.
- [Bre01] Leo Breiman. Random forests. *Machine Learning*, 45(1):5–32, Oct 2001. ISSN: 1573-0565.
- [CC17] Samuel W.K. Chan and Mickey W.C. Chong. Sentiment analysis in financial texts. *Decision Support Systems*, 94:53–64, 2017.
- [CDBF17] Erik Cambria, Dipankar Das, Sivaji Bandyopadhyay, and Antonio Feraco. *A Practical Guide to Sentiment Analysis*. Springer Publishing Company, Incorporated, 1st edition, 2017. ISBN: 3319553925.
- [HCL03] Chih-Wei Hsu, Chih-Chung Chang, and Chih-Jen Lin. A practical guide to support vector classification. Technical report, Department of Computer Science, National Taiwan University, 2003.
- [HSK18] Bruno Miranda Henrique, Vinicius Amorim Sobreiro, and Herbert Kimura. Stock price prediction using support vector regression on daily and up to the minute prices. *The Journal of Finance and Data Science*, 4(3):183 – 201, 2018. ISSN: 2405-9188.
- [Jol11] Ian Jolliffe. *Principal Component Analysis*, pages 1094–1096. Springer Berlin Heidelberg, Berlin, Heidelberg, 2011. ISBN: 978-3-642-04898-2.
- [KSR16] Roman Klinger, Surayya Samat Suliya, and Nils Reiter. Automatic Emotion Detection for Quantitative Literary Studies – A case study based on Franz

- Kafka’s “Das Schloss” and “Amerika”. In *Digital Humanities 2016: Conference Abstracts*, pages 826–828, Kraków, Poland, July 2016. Jagiellonian University and Pedagogical University.
- [LD15] Ronny Luss and Alexandre D’Aspremont. Predicting abnormal returns from news using text classification. *Quantitative Finance*, 15(6):999–1012, 2015.
- [Liu15] Bing Liu. *Sentiment analysis: Mining opinions, sentiments, and emotions*. Cambridge University Press, 2015. ISBN: 9781139084789.
- [MG16] A.C. Müller and S. Guido. *Introduction to Machine Learning with Python: A Guide for Data Scientists*. O’Reilly Media, 2016. ISBN: 9781449369897.
- [Mit15] Ryan Mitchell. *Web Scraping with Python: Collecting Data from the Modern Web*. O’Reilly Media, Inc., 1st edition, 2015. ISBN: 1491910291.
- [NH15] Clemens Nopp and Allan Hanbury. Detecting risks in the banking system by sentiment analysis. *Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing*, pages 591—600, 2015.
- [PFN15] Nicolas Pröllochs, Stefan Feuerriegel, and Dirk Neumann. Enhancing sentiment analysis of financial news by detecting negation scopes. In *2015 48th Hawaii International Conference on System Sciences*, pages 959 – 968. IEEE, 2015. ISBN: 9781479973675.
- [Ras15] Sebastian Raschka. *Python Machine Learning*. Packt Publishing, 2015. ISBN: 9781783555130.
- [RLB⁺17] Navid Rekabsaz, Mihai Lupu, Artem Baklanov, Allan Hanbury, Alexander Dür, and Linda Andersson. Volatility prediction using financial disclosures sentiments with word embedding-based IR models. *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics*, pages 1712–1721, 2017.
- [RQH10] R. Remus, U. Quasthoff, and G. Heyer. Sentiws – a publicly available german-language resource for sentiment analysis. In *Proceedings of the 7th International Language Resources and Evaluation (LREC’10)*, 2010.
- [SC09a] Robert P. Schumaker and Hsinchun Chen. A quantitative stock prediction system based on financial news. *Information Processing & Management*, 45(5):571 – 583, 2009. ISSN: 0306-4573.
- [SC09b] Robert P. Schumaker and Hsinchun Chen. Textual analysis of stock market prediction using breaking financial news: The azfintext system. *ACM Transactions on Information Systems*, 27, 2009.
- [SS04] Alex J. Smola and Bernhard Schölkopf. A tutorial on support vector regression. *Statistics and Computing*, 14(3):199–222, 2004. ISSN: 0960-3174.

- [SZHC12] Robert P. Schumaker, Yulei Zhang, Chun-Neng Huang, and Hsinchun Chen. Evaluating sentiment in financial news articles. *Decision Support Systems*, 53:458—464, 2012.
- [Tet07] Paul C. Tetlock. Giving content to investor sentiment: The role of media in the stock market. *The Journal of Finance*, 62:1139–1168, 2007.
- [Wal10] Ulli Waltinger. Germanpolarityclues: A lexical resource for german sentiment analysis. In *Proceedings of the Seventh International Conference on Language Resources and Evaluation (LREC)*, Valletta, Malta, May 2010. electronic proceedings.

Attachments

GFSL

abbau	aktivitat	anspannet	aufstieg
abbrech	aktivst	anspannt	aufstock
abbruch	akzeptabel	anspanntet	aufstocket
abfall	akzeptabl	ansprech	aufstockt
abgab	akzeptanz	anspruch	aufstocktet
abgemacht	akzepti	anspruchsvoll	auftrag
abgesichert	akzeptiert	ansteig	auftrieb
abgesturzt	alarm	anstieg	aufwand
abgewertet	alarmbereitschaft	anstreng	aufwart
abhäng	alptraum	anteil	aufwartstr
ablehn	andrang	antreib	aufwend
abmach	androh	antrieb	aufwert
abmachet	anfall	arbeitslos	aufwertet
abmacht	angebot	arbeitslosst	aufwertetet
abmachtet	angekurbelt	arm	aufwertt
abn	angemess	armst	aufwerttet
abnahm	angemessen	armut	ausbau
abnehm	angemessn	atemberaub	ausbaut
abrutsch	angenahert	attraktiv	ausbaust
abschaff	angereichert	attraktivitat	ausbaut
abschliess	angesehen	attraktivst	ausbautet
abschluss	angesehen	aufbess	auseinandersetzt
abschreck	angesehn	aufbesser	ausfall
abschwach	angespannt	aufbesserst	ausfuhr
absenk	angriff	aufbessert	ausgab
absich	angstlich	aufbessertet	ausgebaut
absicher	anheb	aufgebessert	ausgeg
abstieg	ankurbel	aufgebracht	ausgeglichen
abstoss	ankurbeln	aufgehob	ausgeglichn
absturz	ankurbelt	aufgemuntert	ausgeschaltet
absturzet	ankurbeltet	aufgeregt	ausgeweitet
absturzs	ankurb1	aufgestockt	ausgewog
absturzt	anlag	aufgewertet	ausgewogen
absturztet	anleg	aufheb	ausgewogn
absurd	annah	aufhebt	ausgezeichnet
abwartstr	annaherst	auflos	ausgleich
abweich	annahert	aufmunt	ausreich
abwert	annahertet	aufmunternd	ausschalt
abwertet	annehmbar	aufmunterst	ausschaltet
abwertetet	annehmbarst	aufmuntert	ausschaltetet
abwertt	annullier	aufmuntertet	ausschaltt
abwerttet	anpass	aufreg	ausschalttet
abzeptabl	anpassungsfah	aufreget	ausscheid
abzock	anreich	aufregt	ausserordent
affar	anreicher	aufregtet	aussichtslos
aggression	anreicherst	aufruhr	aussteig
aggressivitat	anreichert	aufschrei	austeigt
agil	anreichertet	aufschreis	ausweit
agilitat	anreiz	aufschwing	ausweitert
akquisition	ansehn	aufschwung	autonomi
aktiv	anspann	aufsteig	auweitet

bahnbrech	beisteuertet	besorgniserreg	dankbar
bankrott	beitrag	besorgt	dankbarst
barri	beitret	bess	danket
beachtenswert	beitritt	besser	dankt
beachtlich	bekannt	besserst	danktet
beauftrag	bekraft	bessert	dauerhaft
beauftraget	bekraftiget	bessertet	deal
beauftragt	bekraftigt	besserungn	defekt
beauftragtet	bekraftigtet	best	defizit
bedarfsgerecht	belast	bestat	defizitar
bedarfsorientiert	belastbar	bestatiget	defizitarst
bedau	belastbarst	bestatigt	depression
bedenk	beleb	bestatigtet	depressiv
bedeut	belebet	bestmog	depressivst
bedeutsam	belebt	bestraf	desast
bedroh	belebtet	bestrafet	deutlich
bedrohet	beliebt	bestraft	diplomat
bedroht	belohn	bestraftet	diskrediti
bedrohtet	belohnet	beteil	diskreditier
beeindruck	belohnt	beteiliget	diskreditieret
beeindrucket	belohntet	beteiligt	diskreditierst
beeindruckt	bemerkenswert	beteiligtet	diskreditiert
beeindrucktet	bemuh	betracht	diskreditiertet
beeintracht	bemuht	betreff	diss
beeintrachtiget	berausch	betrifft	dissens
beeintrachtiget	berechn	betroff	disziplin
beeintrachtiget	berechnet	betrog	diszipliniert
beend	bereich	betrug	divid
beendet	bereicher	betrugt	dominanz
beendetet	bereicherst	beunruh	domini
beendt	bereichert	beunruhiget	dominier
beendetet	bereichertenbereit	beunruhigt	dominieret
beford	bereichertet	beunruhigtet	dominierst
beforder	bereit	bewahrt	dominiert
beforderst	bericht	bewill	dominiertet
befordert	beruh	bewilligt	drama
befordertet	beruhiget	bewirk	dramat
befried	beruhigt	bewirkt	drang
befriediget	beruhigtet	bewirkt	dranget
befriedigt	beruhmt	bewirktet	drangt
befriedigtet	beschad	bewund	drangtet
befrist	beschadigt	bewundernswert	drastisch
befristet	beschafft	bewunderst	droh
befurcht	beschäftigt	bewundert	drohet
befurchtet	beschäftigungslos	bewundertet	droht
befurchtetet	beschäftigungslosst	bezahl	drohtet
befurchtt	bescheid	bezahlt	drohung
befurchttet	bescheiden	blendend	drossel
begehrt	bescheidn	blockad	drosseln
begeist	beschleun	blocki	drosselt
begeister	beschleuniget	blockier	drosseltet
begeisterst	beschleunigt	bluh	drossl
begeistert	beschleunigtet	bluhet	druck
begeistertet	beschränk	bluht	druckend
begrenz	beschränket	bluhtet	druckt
begrenzt	beschränkt	bonitat	durchdacht
begrenztet	beschränktet	bonus	durchhalt
begunst	beschuld	bonuss	durchschlag
begünstigt	beschuldigt	boom	dynam
begünstigt	beschw	breitgefachert	effektiv
begünstigtet	beschwerd	brockeln	effektivst
behag	beseit	brockelnd	effektivvoll
beherrscht	beseitiget	brockelt	effizient
behind	beseitigt	bussgeld	effizienz
behindert	beseitigtet	chaos	ehrgeiz
behindert	besitz	chaotisch	eif
behindertet	besitzs	comeback	eiferst
beigesteuert	besond	crash	eifert
beileg	besonder	dampf	eifertet
beisteu	besondere	dampfet	eifrig
besteuerst	besonderer	dampft	eigenkapital
besteuert	besondererst	dampftet	eign
	besorgnis	dank	eignet

eignetet	entwirr	erschütternd	feieret
eignt	entwirret	erschütterst	feiern
eigntet	entwirrst	erschüttert	feierst
einbrech	entwirrt	erschüttertet	feiert
einbruch	entwirrtet	ersparnis	feiertet
einbuss	erarbeit	erstaun	fertig
einbusst	erarbeitet	erstaunet	fertiget
eindrucksvoll	erarbeitetet	erstaunt	fertigt
einflussreich	erarbeit	erstauntet	fertigtet
eingekauft	erarbeitet	erstklass	fest
eingeschränkt	erb	erstrebenswert	fiel
eingespart	erbet	ertrag	finanzi
eingestürzt	erbst	ertragreich	finanzier
einhalt	erbt	erweit	finanzierungskost
einheit	erbtet	erweiter	finanzkris
einkassi	erdruck	erweiterst	finanzspritz
einkauf	erfahr	erweitert	fit
einkaufet	erfahren	erweitertet	fitness
einkauft	erfahrn	erwünscht	fitt
einkauftet	erfolg	erzeug	flaut
einkomm	erfolgreich	erzeuget	fleiss
einnahm	erfolgserlebnis	erzeugnis	fleissig
einschränk	erfreu	erzeugniss	flexibel
einschränket	erfreuet	erzeugt	flexibilität
einschränkt	erfreulich	erzeugtet	flexibl
einschränktet	erfreust	erzielt	fliess
einschuchter	erfreut	erzielet	fliessend
einsink	erfreutet	erzielt	fliesst
einspar	ergebnislos	erzieltet	flori
einsparet	ergebnisreich	etabliert	florier
einsparst	ergiebt	euphor	florieret
einspart	erheb	euphori	florierst
einspartet	erhoh	existenzbedroh	floriert
einsteig	erhohet	exklusiv	floriertet
einsturz	erhoht	exklusivst	fluch
einsturzet	erhohtet	exorbitant	folgenschw
einsturzs	erhol	expandi	folgenschwer
einsturzt	erholet	expandier	folgenschwerst
einsturztet	erholsam	expandiert	ford
eintrag	erholt	expandirt	forder
einwandfrei	erholtet	expansion	forderst
einwandfreist	erklimm	expansionskur	fordert
energ	erklimmt	explosiv	fordertet
engagement	erleicht	explosivst	fortdauernd
engagi	erleichter	export	fortschreit
engagier	erleichterst	exporti	fortschritt
engagieret	erleichtert	exportiert	fragil
engagierst	erleichtertet	exzellent	fragwurd
engagiert	erlos	fabelhaft	freiheit
engagiertet	erloset	fahig	freu
enorm	erlost	fahrlass	freud
enthusiasmus	erlostet	fall	freudig
enthusiast	ermog	fallend	freuet
entlass	ermoglicht	fallt	freundlich
entlast	ermoglicht	falsch	freust
entlastet	ermoglichtet	faszini	freut
entlastetet	ermunternd	faszinier	freutet
entlastt	ermut	faszinieret	friedlich
entlasttet	ermutiget	faszinierst	frist
entlohn	ermutigt	fasziniert	fruchtbar
entschad	ermutigtet	fasziniertet	fruchtbarst
entschloss	erneu	fatal	frust
entschlossen	erneuerst	fehl	frustration
entschlossn	erneuert	fehlend	frustri
entschluss	erneuertet	fehleranfall	frustrier
entspann	ernst	fehlerfrei	frustrieret
entspannet	ernuchter	fehlerfreist	frustrierst
entspannt	ernuchternd	fehlerhaft	frustriert
entspanntet	erreich	fehlermeld	frustriertet
enttausch	erreicht	fehlkauf	fuhr
enttauschet	erreicht	fehltritt	fuhrend
enttauscht	erreichtet	fehlverhalt	fuhret
enttauschtet	erschutt	feier	fuhrst

fuhrt	gemindert	haft	hochtreib
fuhrtet	genehm	handel	hochwert
fuehrung	genehmiget	handeln	hoff
funktioni	genehmigt	handelsaufnahm	hoffet
funktioni	genehmigtet	handelt	hoffnung
funktionieret	genotigt	harmon	hoffnungslos
funktionierst	genutzt	harmonisi	hoffnungsvoll
funktioniert	gerecht	harmonisier	hofft
funktioniertet	gerechtfertigt	harmonisieret	hofftet
funktionsfah	gerettet	harmonisierst	hoh
furcht	gering	harmonisiert	hohepunkt
furchtbar	gerutscht	harmonisiertet	hundertprozent
furchtbarst	geschafft	hart	ideal
furchterreg	geschafft	hartnack	idyll
furchtet	geschäftsauflös	haushaltsdefizit	illegal
furchtetet	geschick	haushaltsdefizit	imm
furchtlos	geschrumpft	heb	immens
furchtt	geschuldet	hebt	imponier
furchttet	gesenkt	heftig	import
garanti	gesichert	heikel	importi
garantienlgarantiert	gesorgt	heikl	importiert
garantiert	gespart	heit	imposant
gebessert	gespendet	heiter	inspirier
gebluht	gestarkt	heiterst	individuell
gebuhrt	gesteigert	heitr	ineffizient
gedampft	gestieg	hektik	ineffizienz
gedankt	gestockt	hektisch	inflation
gedrängt	gestort	hellig	inkompetenz
gedroht	gestottert	hemm	innovation
gedrosselt	gestrahlt	hemmet	innovativ
gedruckt	gestreikt	hemmt	innovativst
geduld	gesturzt	hemmtet	inspirier
geeifert	gesunk	hemmung	insolvent
geeignet	getrennt	herabgesetzt	insolvenz
geerbt	getrostet	herabsetz	inspiri
gefahr	gewach	herabsetzet	inspirier
gefahrtd	gewachs	herabsetzt	instabil
gefahrtdet	gewackelt	herabsetztet	instabilitat
gefahrtdetet	gewagt	heranwachs	intakt
gefahrtdt	gewährleist	heraufsetz	intelligent
gefahrtdtet	gewalt	herausford	intensiv
gefall	gewarnt	herausforder	intensivist
gefeiert	gewinn	herausforderst	interess
gefertigt	gewinnbring	herausfordert	interessant
gefestigt	gewinnzon	herausfordertet	interessi
gefordert	gezielt	herausgefordert	interessier
gefrenut	gezittet	herausrag	interessieret
gefuhrt	gezogert	herausraged	interessierst
gefurchtet	glaubwurd	hervorrag	interessiert
gegluckt	gleichwert	highlight	interessiertet
gehaltsszulag	gluck	hilf	investor
gehemmt	gold	hilflos	investi
gehofft	golden	hilfreich	investier
gejubelt	goldig	hilfsprogramm	investieret
gekippt	goldn	hilfszahl	investierst
geklagt	grandios	hind	investiert
geklektert	gratulation	hindernis	investiertet
gekostet	gravier	hinderniss	investition
gekundigt	greifbar	hindert	involvi
gekürzt	greifbarst	hinreich	involviert
gelahmt	grenz	hinzufug	jahresgewinn
gelass	grenzenlos	hinzufuget	jahresverlust
gelassen	gross	hinzufugt	jubel
gelassn	grossart	hinzufugtet	jubeln
geldgeb	grossspur	hinzugefugt	jubelt
geldstraf	grosst	hoch	jubeltet
gelernt	grosstmog	hochattraktiv	jubl
gelindert	grosszug	hochgestellt	kapital
gelohnt	grundleg	hochgrad	kapitalerhoh
gelung	grundlich	hochklass	kapitalertragssteu
gelungen	gunstig	hochrang	kapitalverbrech
gelungn	gut	hochstmog	kapitulation
gemeinschaft	gutst	hochststand	kassi

kassiert	konsistenz	last	meisterhaft
katastroph	konsolidi	laun	mild
katastrophal	konsolidier	lebensfah	milder
kauf	konsolidieret	lebhaft	mildernd
kauft	konsolidierst	leer	mildert
kein	konsolidiert	legal	mind
kipp	konsolidiertet	legalitat	minder
kippet	konstant	leicht	minderst
kippst	konstanz	leichtglaub	mindert
kippt	konstruktiv	leichtig	mindertet
kipptet	konstruktivst	leichtsinn	minus
klag	konsulti	leid	mis
klaget	konsultier	leidend	miserabel
klaglos	konsultieret	leidet	miserabl
klagt	konsultierst	leistung	misstrau
klaget	konsultiert	leistungsfah	mittelmass
klar	konsultiertet	leistungsstark	mobil
klarst	konsum	lern	mobilisi
klett	konsumi	lernet	mobilisier
kletterst	konsumiert	lernt	mobilisieret
klettert	konterproduktiv	lerntet	mobilisierst
klettertet	kontinui	limit	mobilisiert
knack	konzertiert	limiti	mobilisiertet
knackt	kooperation	limitier	mobilitat
knapp	kooperativ	limitiert	mod
knappheit	kooperativst	lind	modern
knappst	kooperi	linderst	modernisi
kollabi	koordini	lindert	modernisier
kollabier	koordiniert	lindertet	moiglich
kollabieret	korrekt	liquid	monopol
kollabierst	korrektur	liquidation	moral
kollabiert	korrigi	liquiditat	motivation
kollabiertet	korrigier	lob	motivi
kollaps	korrigiert	lock	motiviert
kollidi	korrupt	locker	muhelos
kollidier	korrupcion	lockerst	muhsam
kollidieret	kost	lockr	mut
kollidierst	kostbar	lohn	nachfrag
kollidiert	kostbarst	lohnend	nachhalt
kollidiertet	kostengunst	lohnet	nachlass
kollision	kostenintensiv	lohnsteu	nachtei
komfortabel	kostenintensivst	lohnt	naiv
komfortabl	kostenlos	lohtet	naivitat
kompatibilitat	kostenlosst	loyalitat	negativ
kompensation	kostet	lukrativ	negativbescheid
kompensi	kraft	lukrativst	negativitat
kompensier	kraftig	luxorios	negativst
kompensieret	kraftvoll	luxorios	nennenswert
kompensierst	kris	macht	nervos
kompensiert	kritik	machtig	nervositat
kompensiertet	kritisch	mag	neuordn
kompetent	kundig	mager	nicht
kompetenz	kundiget	magerst	niedergang
komplett	kundigt	magr	niederlag
kompromiss	kundigtet	makellos	niedrig
konflikt	kurseinbruch	mangel	not
konfronti	kursgewinn	mangelbehaftet	notfall
konfrontiert	kursverli	mangelhaft	notig
konjunkturabschw	kursverlier	mangelnd	notiget
konjunkturaufschw	kursverlust	manipulation	notigt
konjunkturaufschwung	kursziel	manipuli	notigtet
konjunkturruckgang	kurz	marod	notstand
konkret	kurzet	massiv	nutz
konkur	kurzt	massivst	nutztet
konkurrenz	kurztet	maximal	nutzlich
konkurrenzfah	kurzung	maximi	nutzt
konkurrenzkampf	labil	maximier	nutztet
konkurri	lahm	maximieret	nzureich
konkurs	lahmet	maximierst	offensiv
kons	lahmt	maximiert	ohn
konsens	lahmtet	maximiertet	optimalitat
konsequent	lahmung	maximum	optimismus
konsistent	langsam	meist	optimist

ordent	qualitativ	riesig	schwach
ordnungsgemass	qualitativst	risiko	schwachung
panik	qualitätsverbesser	riskant	schwer
panisch	quartalsgewinn	riski	schwerfall
pann	quartalsverlust	riskier	schwerst
paradi	rach	riskieret	schwerwieg
paradies	rachend	riskierst	schwierig
paradiess	racht	riskiert	schwung
partn	rational	riskiertet	sehr
partnerschaft	realisi	rivalitat	senk
pech	realisiert	robust	senket
perfektion	realist	rosig	senkt
perfektionismus	recht	ruckgang	senktet
perfektionist	rechtfert	rucklauf	senkung
pessimismus	rechtfertiget	ruckschlag	sensation
pessimist	rechtfertigt	ruckstell	sensationell
planlos	rechtfertigtet	rucktritt	sensibel
planlosst	rechtlich	ruh	sensibl
planmass	rechtmass	ruhig	serios
planvoll	rechtsgult	ruin	seriositat
pleit	rechtswidr	ruini	sich
plus	reduktion	ruinier	sicher
poplaritat	reduzi	ruinos	sicherst
popular	reduzier	rutsch	sichert
popularitat	reduzieret	rutschet	sichertet
positiv	reduzierst	rutscht	sichr
positivitat	reduziert	rutschtet	simpel
potent	reduziertet	sani	simpl
potenz	regulation	sanier	sink
prachtig	reguli	saniert	sinnlos
prachtig	regulier	sank	sinnvoll
prachtvoll	reguliert	schad	skandal
praferenz	reibungslos	schadig	skandalos
praktikabel	reich	schadlich	skeptisch
praktikabl	reichhalt	schaff	solid
praktisch	reichlich	schaffet	sonnig
prami	reichtum	schafft	sorg
prazision	reinfall	schafftet	sorgenfrei
preisgunst	reizend	scheit	sorgenfreist
preissturz	reizvoll	scheiternd	sorget
preissturzs	rekord	scheitert	sorgfalt
prekar	rekordergebnis	schenkung	sorglos
prekarst	rekordhoch	schlecht	sorgsam
privilegiert	rekordtief	schliessung	sorgt
probl	rekordverlust	schlimm	sorgetet
problem	rekordwert	schmelz	souveran
problemат	relevant	schmelzend	spar
problemlos	relevanz	schmelzt	sparet
problemlosst	rendit	schmolz	sparkur
produktiv	renovier	schnappch	sparkurs
produktivitat	rentabel	schnell	sparprogramm
produktivst	rentabilitat	schnellig	sparsam
professionell	rentabl	schock	sparst
profit	reparatur	schockier	spart
profitabel	respektabel	schockn	spartet
profitabl	respektabl	schonungslos	spektakular
profiti	respekti	schrumpf	spektakularst
profitiert	respektier	schrumpfet	spekulation
progression	respektvoll	schrumpft	spekuli
projekt	respektier	schrumpftet	spekuliert
prot	rett	schub	spend
protesti	rettet	schuld	spendet
protestier	rettetet	schuldet	spendetet
protestiert	rettung	schuldetet	spendt
ptimal	revanchi	schuldhaft	spendtet
qualifikation	revanchier	schuldig	spitzenverdi
qualifizi	revanchieret	schuldlos	sprung
qualifizier	revanchierst	schuldn	spurbar
qualifizieret	revanchiert	schuldt	spurbarst
qualifizierst	revanchiertet	schuldtet	stabil
qualifiziert	rezession	schutz	stabilisi
qualifiziertet	richtig	schutzmassnahm	stabilisier
qualitat	riesengross	schutzs	stabilisieret

stabilisiert	stürzt	übersteig	unprofitabel
stabilisiert	sturzs	übersteigt	unprofitabl
stabilisiertet	stürzt	ubertreff	unqualifiziert
stabilität	sturztet	ubertrifft	unrealist
stagnation	stutz	ubertroff	unregelmass
stagni	stutzend	umbruch	unrentabel
stagniert	stutzt	umfangreich	unrentabilität
standhaft	stutzt	umgang	unrentabl
stark	subvention	umstritt	unruh
starket	subventioni	umstritten	unschlagbar
starkt	subventionier	umweltschad	unschlagbarst
starkt	subventionieret	unabhäng	unschon
steuerhinzuzieh	subventionierst	unangefocht	unschuld
steig	subventioniert	unangefochten	unserios
steigend	subventioniertet	unangenehm	unsich
steiger	tatkraft	unantastbar	unsicher
steigert	tauglich	unantastbarst	unsicherst
steigert	tendenz	unattraktiv	unsichr
steigertet	tendi	unattraktivität	unsolidar
stellenabbau	tendiert	unaufholbar	unstet
stellenabbaus	teu	unbefried	unt
stetig	teuerst	unbefristet	unterbrech
steu	teufelskreis	unbegrenzt	unterbricht
steuererhöh	teufelskreiss	unbeliebt	unterbroch
steuerhinterzieh	teur	unbeschwert	unterdruck
stieg	tief	unbesiegbar	unterdrucket
stiftung	tilg	unbesiegbarst	unterdrückt
stillleg	tilgt	unbezahlbar	unterdrückt
stillstand	tilgung	unbezahlbarst	unterentwickelt
stimuli	trag	unbrauchbar	untergang
stimulier	trage	undurchsicht	untergedrückt
stimuliert	tragheit	unein	untergeh
stimuliert	traum	uneingeschränkt	unternehmer
stimuliert	trauma	unerbitt	unterstütz
stimuliertet	traumat	unerklar	untreu
stock	traumatisi	unermess	unubertreff
stocket	trenn	unermud	unubertroff
stockt	trennet	unerreicht	unubertroffen
stocktet	trennt	unerschrock	ununterbroch
stor	trenntet	unerschrocken	ununterbrochen
storet	triumph	unerschrockn	unverantwort
storst	triumphi	unersetz	unvergleichbar
stort	triumphier	unertrag	unvergleichbarst
stortet	triumphieret	unerwartet	unverhältnismass
störung	triumphierst	unerwünscht	unverhofft
störungsfrei	triumphiert	unfah	unverkrampt
störungsfreist	triumphiert	unfair	unvermeid
stott	trost	unfairst	unverstand
stotterst	trostet	unfall	unvorhergesehen
stottert	trostet	unfehlbar	unvorhergesehen
stottert	trostt	unfehlbarst	unwicht
straf	trostet	ungeahnt	unwidersteh
strafbar	trub	ungebroch	unwirksam
strafbarst	trubend	ungebrochen	unwirtschaft
strafverfahren	trubt	ungedeckt	unwiss
strafzins	trugschluss	ungerecht	unzufried
strahl	turbolenz	ungeschickt	unzufrieden
strahlend	turbulent	ungesetz	unzufriedn
strahlet	uber	ungewollt	unzulass
strahlt	überdurchschnitt	unglaub	unzumutbar
strahlt	überholt	ungleich	unzumutbarst
strapaz	überglück	unglück	unzurechnungsfah
streich	überhol	ungunst	unzureich
streicht	überholt	universell	unzuverläss
streik	überholt	unklar	vehement
streiket	überholtet	unklarst	verangst
streikt	überlastet	unklug	verangstigt
streiktet	überrasch	unkompliziert	verantwort
strukturel	überrascht	unlaut	verantwortungsbewusst
strukturier	überschreit	unlauter	verantwortungsvoll
sturm	überschreitet	unlauter	verband
sturmisch	übersteh	unlautr	verbess
sturz	übersteht	unmog	verbesser

verbiets	verlass	verzweifels	widerspruch
verbind	verli	verzweifeln	widerstand
verblind	verlier	verzweifelt	widerstandsfah
verbluff	verlor	verzweifeltet	widrig
verbot	verlust	verzweifl	wiederhergestellt
verboten	verlustzon	viel	wiederherstell
verbotn	vermach	vielfalt	wiederherstellt
verbraucht	vermachtet	vielseit	wiederherstellt
verbraucht	vermachtet	vielversprech	wiederherstellt
verbundet	vermachtet	vielzahl	widerstand
verdacht	vermind	visionar	willkur
verderb	verminder	volum	wirksam
verdi	verminderst	voranbring	wirkungslos
verdien	vermindert	vorankomm	wirkungsvoll
verdient	vermindert	vorantreib	wirtschaft
verdientet	vernunft	vorgesorgt	wirtschaftskris
verdirbt	verrechn	vorsicht	wohlergeh
verdorb	verrechnet	vorsorg	wohlhab
veredel	verring	vorsorget	wohlstand
veredeln	versag	vorsorgt	wund
veredelt	versagt	vorsorget	zahl
veredeltet	verschieb	vorteil	zahl
veredl	verschiebt	vorteilhaft	zahlung
vereinbar	verschlecht	vorurteil	zahlungsunfah
vereinfach	verschlechter	vorwand	zahlungsverzoger
verfall	verschmutz	vorwart	zerschlag
verfallen	verschmutzt	vorwurf	zerschlagen
verfalln	verschmutzt	vorzeigbar	zerstor
verfalsch	verschmutztet	vorzeigbarst	zerstorer
verfalschet	verschwend	vorzug	zerstoret
verfalscht	verschwender	wach	zerstort
verfalschtet	verschwendet	wachs	zerstort
verfehl	verschwendet	wachsam	zerstortet
verfehlet	verschwendt	wachsend	zielgerichtet
verfehlt	verschwendtet	wachset	ziellos
verfehltet	versich	wachst	zielstreb
vergeb	versicher	wachstet	zitt
vergemacht	versicherst	wachstum	zitterst
vergeud	versichert	wachstumschanc	zittert
vergeudet	versichertet	wackel	zittertet
vergeudetet	verspiel	wackeln	zog
vergeudt	verspielt	wackelt	zogerst
vergeudtet	verspieltet	wackeltet	zogert
vergezogert	verständnis	wackl	zogert
vergross	verständnis	wag	zoll
vergrosser	verstärk	waget	zolln
vergrosserst	verstärket	wagst	zufried
vergrossert	verstärkt	wagt	zufrieden
vergrossert	verstärktet	wagtet	zufriedenstell
vergut	verstoß	warn	zufriedn
vergudet	verteid	warnend	zugelegt
verhandel	verteidigt	warnet	zukauf
verhandeln	verteidigt	wegweis	zukauf
verhandelt	verteidigtet	weitgeh	zukunftsweis
verhandeltet	verteuert	weitlauf	zulag
verhandl	vertrag	weitraum	zulass
verhangisvoll	vertrau	weitreich	zuleg
verhangnis	vertrauenerweck	weitsicht	zuleget
verhangnis	vertrauensvoll	wenig	zulegt
verhangnisvoll	vertrauenswurd	wert	zulegtet
verheer	verurteil	wertgeschätzt	zunahm
verkauf	verurteilt	wertlos	zunehm
verkauft	verurteilt	wertschatz	zuruckgeh
verklein	verurteiltet	wertschatztet	zuruckging
verkleinert	verwerf	wertschatzt	zuruckhalt
verkleinert	verwirr	wertschatztet	zuruckzieh
verkleinert	verzicht	wertsteiger	zusammenarbeit
verkleinert	verzog	wertverfall	zusammenbrech
verlangsam	verzogert	wertverlust	zusammenbruch
verlangsam	verzogert	wertvoll	zusammenhalt
verlangsam	verzogert	wettbewerbswidr	zusatz
verlangsam	verzogert	wichtig	zuschlag
verlangsam	verzogert	widerruf	zuschuss

zusich
zusicher
zusichernd
zustimm
zuverlass

zuversicht
zuvorkomm
zwangslag
zwangsverkauf
zwecklos

zweifel
zweifelhaft
übereinstimm
überfluss
überhoh

überlast
überleg
überschaubar
überschuss