# MODUL VIENNA UNIVERSITY

## WKO WIEN PRIVATE UNIVERSITY

# Opinion Mining and Sentiment Analysis using Rapidminer

Bachelor Thesis for Obtaining the Degree

Bachelor of Science (BSc) in

International Management

Submitted to Mr. Christian Weismayer

Parishek Singh Chauhan

1321030

Vienna, 12 December 2016

## Affidavit

I hereby affirm that this Bachelor's Thesis represents my own written work and that I have used no sources and aids other than those indicated. All passages quoted from publications or paraphrased from these sources are properly cited and attributed.

The thesis was not submitted in the same or in a substantially similar version, not even partially, to another examination board and was not published elsewhere.

_____                    _____

Date                                        Signature

# Abstract

In the recent years, a vast amount of research has been conducted on the topic of sentiment analysis and opinion mining. Businesses and organizations understand the potential benefits of developing sentiment analysis and opinion mining systems. In this study, Rapidminer as a solution is proposed for analyzing online product reviews. A corpus of 200 reviews for the 25hours hotel in Vienna, Austria was collected from the Tripadvisor website. Overall sentiment analysis as well as aspect-based sentiment analysis was performed on the reviews. The results were then compared with the star rating provided by the reviewer using SPSS software to check the accuracy of the results. Finally, using the results from aspect-based sentiment analysis, linear regression was used to predict the sentiment using the most frequently appearing aspects. The purpose of this study is to show that open source software like Rapidminer can be used effectively to calculate the sentiment from online reviews as well as for aspect-based sentiment analysis. The results show that Rapidminer is an effective tool. Aspect-based sentiment analysis can be used to predict sentiment and thereby business can use it to improve overall customer satisfaction by focusing on enhancing certain aspects of their products and services.

# Table of Contents

## List of Tables

- Table-1, Polarity confidence of the review text
- Table-2: Polarity confidence of the review title
- Table-3: Frequency of the aspects

## List of Figures

- Figure-1: Star rating
- Figure-2: Traveller type
- Figure-3: Aspects and their occurrence in reviews
- Figure-4: Spearman's correlation for star rating and sentiment
- Figure-5: Spearman's correlation for sentiment of the review text and of the title text
- Figure-6: Regression analysis

## List of Abbreviations

- SPSS: Statistical Package for the Social Sciences
- API: Application Program Interface
- E-WOM: Electronic Word of Mouth

# 1    Introduction

## 1.1   Introduction and Background

People all make decisions in their daily lives based on information that they gather from different sources. In recent years, there has been exponential growth of the Internet. A significant number of people gain access to the Internet everyday. As per the website (www.internetlivestats.com, 2016), about 40% of the population of the world accesses the Internet on a daily basis. It is largely believed that this number is expected to keep growing in the future. As the Internet becomes a normal in people's daily lives, they are also more likely to leave a larger footprint on the Internet. The footprint in question here is the plethora of raw data that people add to the Internet every day. Every minute of Internet usage is equivalent to 640 terabytes of data transferred (Burgess, 2013). A large chunk of this data is generated by social media websites like Facebook, Twitter, Instagram, etc. The Internet has become an essential part of people's lives. People also use the Internet to express their opinions towards certain things. These can range from the recent election, to the latest movie or the latest smartphone to hit the market. These opinions contain a treasure trove of data, which can be analyzed for meaningful insights into peoples' minds.

Hence, it is no wonder that companies, organizations, etc. would be interested in finding out what people think about their products, campaigns and so on. This brings us on to the topic of opinion mining and sentiment analysis. These two terms are interchangeable and have the same meaning. It becomes clear that with the vast amounts of raw data available on the Internet, the significance of sentiment analysis has increased in the recent years. Customers are more likely to do online research about the products before they make their purchasing decisions. 'In the modern era of information and communication technologies, it has become quite common that customers create their opinion not just by talking to friends or reading expert reviews in magazines but also reading reviews of other customers on the Internet' (Prichystal, 2016, p. 373). This online research mostly involves reading online reviews. This information is a potential treasure trove of information for businesses looking to gain meaningful insights into how customers feel about their products.

Hence, it becomes clear that there is a need to have a method of extracting meaningful information from this data.

## 1.2 Objectives of the study

The objective of the study is to analyze online reviews, which are widely available on the Internet. Most reviews available online contain the text of the review, the associated star rating provided by the reviewer themselves and a small title for the review. This study involves collecting reviews from an online source, analyzing the text contained in the review and calculating the overall sentiment associated with the text. As a final step aspect-based sentiment analysis is also performed on the reviews. Rapidminer software was used to conduct the analysis. The intention of this study was to find out if such an analysis would generate meaningful results and show that aspect-based sentiment analysis can help businesses to fine-tune their offerings in order to increase customer satisfaction. This is done by comparing the output of the textual analysis against the already mentioned star rating provided by the reviewer. The aspects are then compared against the star rating. The goal was to check if the mention of certain aspects could be used to predict the star rating associated with the review. Another goal of the study was to compare the sentiment in the review text provided by the reviewer against the calculated sentiment in the title of the review. The aim was to find out whether the title of the review alone can be used to predict the sentiment of the entire review.

## 1.3 Structure of the study

The study involved various stages. In the first stage a corpus of reviews about the 25hours hotel located in Vienna, Austria from the Tripadvisor website was collected. In the second stage, the reviews were analyzed using text-mining software Rapidminer to calculate the sentiment associated with the text of the review and the title of the review. During this stage, aspect-based sentiment analysis on the text of the review was also performed. In the third stage, analysis on the gathered reviews was performed using IBM SPSS Statistics software. A series of tests were performed on the data, which are listed as follows:

a) Spearman's correlation on sentiment associated with the text of the review and the star rating provided by the reviewer.

b) Spearman's correlation on sentiment associated with the text of the review and the title of the review.

c) Linear regression is performed to find out if star rating can be predicted using certain aspects mentioned in the text of the review.

# 2 Literature Review

## 2.1 How opinion was gathered in the pre-internet era

Opinion mining and sentiment analysis stem from the need to gather public opinion. In the pre-internet era, public opinion was gathered by conducting polls, surveys, etc. This process usually took longer and was an enduring task. Moreover, the cost of conducting surveys and polls was also usually high. Another problem that exists with gathering public opinion using this method is that it when we talk about gathering the opinion of the public, the aim is to capture the opinion of the entire population, however this is not entirely feasible in reality. In reality, a representative sample population is chosen which will reflects the opinion of the overall population. The problem that then arises through this method is that it is mostly impossible to have a sample that can be used to fully represent the overall population. Methods of gathering public opinion have been researched since the early 1900's. However, only recently starting from the year 2000 onwards other methods of gathering public opinion have been used. Not only polls and surveys but also the so-called online chatter can be used to gather and measure public opinion. This method eliminates the need for choosing a representative sample, the sentiment of the population as a whole can be analyzed. Moreover, it does not involve manual data entry in any of the steps. It is not as expensive as the older methods as everything can be done using computers and from a single location. Another key feature of this method is that the process of collecting individual's opinions is easier in the sense that the all the data that needs to be analyzed is already available to the researcher and they do not need to make any special effort in collecting the data. This greatly increases the efficiency and the speed of the overall process of gathering public opinion.

Gathering public opinion has been of significance in many areas namely, politics, business, governance, scientific research, etc. Political parties can be considered in

this regard as the ones who greatly need to gather and correctly measure public opinion. A good example of this is the elections in the United States of America where, political parties invest heavily in opinion polls, which help them choose the absolute correct future strategy.  It can be safely said that without having an idea about what the general public thinks about the political parties objectives, the political party would have a significant disadvantage compared to the other well-informed opponents. This could be seen clearly in the recent 2017 presidential elections in the United States of America where the Republicans were able to correctly measure public opinion and their presidential candidate was able to connect with the population on a greater level than the opponent. Measuring public opinion with such accuracy was certainly not possible in the pre-internet era. The advent of opinion mining and sentiment analysis has changed the world of politics and elections. Websites like Twitter and Facebook generate huge amounts of data each day, much of it is online chatter, which can also be analyzed. Political parties can mine this data to find out which issues are of highest concern to the people and what is their opinion on such issues.

Businesses are another group, which could benefit greatly from being able to correctly measure public opinion. However in this case, businesses are not concerned with political issues but rather the performance of their products and services and the opinions held by their customers about them. They could also use it to find out how their brand is perceived and the preferences of their customers. Moreover, they could also use it to do research on their competitors' brand as well as their offerings. In the pre-internet era gathering peoples' opinion was not possible on such large scale and the process was difficult and expensive. Market research agencies were employed to conduct large scale and time-consuming research. Even after spending a significant portion of their capital on market research, business could not be sure of the results, as it's difficult to accurately choose a sample that is fully representative of the population. With the growing popularity of the Internet, businesses can now use online chatter to gather public opinion. Research in the form of polls and surveys can also be conducted online. It is not required to choose a representative sample, information is already readily available and the entire population can be used to gather public opinion. This has increased the accuracy as well as reduced the cost and time required to measure public opinion. Moreover, in

the pre-internet era, opinion could not be gathered in real-time whereas, now it is possible to gather opinion in real time.

Governments always have needed to be able to correctly measure public opinion. Policy-making is a particular field where governments could benefit from knowing the opinion of the public, currently hardly any government is able to fully integrate public opinion into the policy making process. Being able to correctly measure public opinion allows governments to fully represent the interests of the people. However, currently this technology is not completely secure and is prone to manipulation by hackers. More research needs to be done in this area to make it completely secure and allow for the correct measurement of public opinion. The ability to fully gather the opinion of the public could be considered the holy grail of democracy. The aim of any democracy is to be fully representative of its population. Such a system would revolutionize how democracies function. The current system where citizens vote for someone to represent themselves and then such a group of such representatives form the government is an archaic system, which needs to be updated. Many functions which democratic governments need to perform like law making and policy making usually take a long period of time. In a system where people could directly represent themselves in the government would allow such processes to become faster and more efficient.

The field of opinion mining and sentiment analysis is only in the very early stages and researchers have barely even scraped the surface. Moreover, this area of research has only emerged since the invention of the Internet and more specifically with the increase in the popularity of social networking websites, micro-blogging websites and online review websites. The future of this field of research is promising and a huge amount of research still needs to be conducted to fully tap its potential.

## 2.2   Online Reviews

A number of researchers have looked into the topic of opinion mining and sentiment analysis. Levy, Duan & Boo in 2013 analyzed only 1-star reviews from 86 Washington DC hotels to try to understand poor or negative reviews. They looked at the most common complaint areas and found that the most common complaints were related to billing, check-in, hotel appearance, Internet, restaurant, room

service, safety and front-desk. They were also able to rank these issues based on their occurrence in the reviews. Furthermore they found out that hotels, which responded to negative reviews, were found to have a higher overall rating. Their study suggests a need for an effective management plan for negative reviews, and appointing a person with strong writing abilities to respond to negative reviews. They further stressed on the need to incorporate a strong and efficient feedback system, which will eventually help the hotel increase its overall rating.

Electronic word of mouth or E-WOM as it is abbreviated is simply the concept of word of mouth but in regards to the online world. It is interesting to look into how E-WOM spreads and what kind of effect it has on customers' travel choices. Filieri & McLeay in 2013 looked into how online reviews help in spreading electronic word of mouth, and how travellers base their decision to stay in a hotel based on online reviews. They found that product ranking or the so-called star rating had the highest influence on travellers' decision. This is due to many reasons, for example, it allows travellers to reduce the time and effort required in their search for an accommodation. Besides this, travellers also look at information factors such as quality, completeness and value-added information, which is not attainable through regular marketing channels. However, they found that the accuracy or relevance of the information available in online reviews was not a major factor in influencing travellers' decisions. They propose that travellers should be considered as sort of co-marketers in today's world who influence the decision of other travellers through the feedback they provide in online reviews.

Another interesting topic is to look into how eager people are to post online reviews. In this regard, a study conducted by Dellarocas, Gao & Narayan in 2010 looked at peoples' eagerness to post online reviews. They looked at two different categories of niche products as well as highly popular and reviewed products. It was found that people are equally likely to post reviews for niche products as well as highly popular products. They did this by looking at reviews posted for movies, it was found that initially people are more likely to post reviews for newly released films, followed by a decline as the revenues of the movie increase, finally as the movie becomes even more popular the eagerness to post reviews increases further amongst the people. Their study suggested that companies should make the number of reviews already

posted as a less prominent feature for niche products, which makes people more likely to post a review.

A key problem in recent years has been the issue of fake online reviews. Fake reviews are a major nuisance and a threat to the credibility of the online rating and review websites. Moreover, it has become increasingly difficult to detect fake reviews and ratings due to the perpetrators becoming better at posting fake online reviews. In this regard, Yoo & Gretzel in 2009 looked at mechanisms to separate fake reviews from real ones. Fake reviews can be really harmful for a website, moreover they may lead to the customer making a wrong decision which in-turn leads to more dissatisfaction in the end. They found that it is extremely difficult to detect fake online reviews as the people who post them are increasingly becoming better at hiding such reviews by basing their structure on other real reviews. They suggest that more research be conducted to confirm their findings.

A number of researchers have looked into the impact of online reviews on sales figures. A study conducted by Hu, Liu, & Zhang in 2008 looked into the effect of online reviews on sales revenue. It is interesting to find out the empirical effects on hotel sales with respect to online reviews. They found positive correlation between product sales and online ratings, however the effect of high online ratings on sales diminished over time. Moreover it was found that customers are likely to pay attention not only to the review score but as well as the reviewer profile, which implies that how often the reviewer has reviewed a product online. Another key finding was that customers are more likely to be influenced by a review when there exist few reviews for the product, however in case there already exist many reviews for a product, a new review, even if it contains new information is not likely to have any effect on the customers perception of the product. In another similar work by Duan, Gu, & Whinston in 2008, a study was constructed to analyze the effect of online reviews on movie box office sales. There findings were interesting as they pointed out that the overall rating of online reviews had no impact on the box office collections however, the amount of online reviews posted about a particular movie had a significant positive effect on the box office collections. According to them this indicates the word of mouth effect of online reviews has a more significant effect on the box office collections than the overall rating in the reviews. A similar study by

Zhu & Zhang in 2010 analyzed the effect of online reviews on product sales in the video game industry. They found that online reviews had a greater impact on sales of less popular games and online games. Eventually, the impact of online reviews depends on the characteristics of the game as well as that of the customer. This further suggests that online reviews play a more prominent role where customers don't have access to other information sources from where they can get feedback on the games. Their study also showed that a negative review could have a greater impact on a niche game. According to them, this phenomenon has allowed the industry to shift towards having more niche products as compared to having few very popular products at the top of the industry. Moreover, they also found that online reviews had significantly less impact on sales of products for customers, which had more experience with the Internet. Judith & Mayzlin in 2006 looked into the impact of online reviews on book sales at two major online book retailers. Interestingly it was found that the increase of positive reviews on a website resulted in higher sales for that particular book on that website. An interesting finding was that the impact of one-star negative reviews was higher than the impact of five-star positive reviews. It was also found that customers believe that they were able to make better purchasing decisions through reading other customers' online reviews. Ye, Law, & Gu in 2009 conducted an empirical study, which indicated that there exists a significant relationship between the hotels' online performance in terms of positive of negative reviews and hotel room sales. Tsaoa, Hsieh, Shih, & Lin in 2015 found that positive reviews had a positive effect on room sales for a hotel and negative reviews had a negative impact. Moreover, a higher number of reviews also strengthened booking intentions in customers. Another key finding was that it is easy to influence conventional rather than non-conventional customers. Poor-service was found to be a key factor, which led to a customer posting a negative online review. Hence, hotels must prioritize providing high quality service in order to gain more positive online reviews.

Many studies have been conducted in order to understand the decision making process of customers, however since the advent of online review websites, researchers are more interested in finding the effect of online reviews on the customers' purchasing behavior. In this regard, Papathanassis & Knolle in 2009 looked into the role of online reviews in the travellers' decision-making process for

choosing holiday destinations. They found that online reviews play a supplementary role in the travellers' decision-making process. They suggest that online reviews be made a central part of the marketing strategy of any organization. Moreover, they stress the need for developing automated text-mining processes to harness the full potential of online reviews. This implies that having an effective strategy to focus on feedback from online reviews is a key factor of success for any kind of business organization. A similar study by Vermeulen & Seegers in 2008 looked directly into the impact of online reviews on customers' decision-making process. They found that irrespective of the review being negative or positive, reading online reviews increased the awareness of the customer about the hotel. Importantly, it was found that the impact of a negative review eventually is nullified. This effect is stronger for less popular hotels, which implies it is easier to increase customer awareness for smaller, newer hotels. Moreover, they found that reading a review from an expert or an average traveller had more or less equivalent impact on customers' minds. In continuation with kind of research Gretzel & Yoo in 2008 conducted a study supported by Tripadvisor, which aimed to look into the role of online reviews on the customers' planning process, that is how customers perceive a review's credibility and usefulness, what motivates travellers to post online reviews and understanding travellers who heavily rely on online travel websites. The study had many interesting findings for example, 97.7% of travelers read online reviews on Tripadvisor. Moreover it was found that online reviews were a key factor in deciding the choice of accommodation for the journey. Travellers trust online reviews posted on reputable websites, preferring a clear style of writing. They also look into the reviewer profile like frequency of posting, objectivity and prior experience. It was also found that travellers rely on reviews to reduce the chance of facing unpleasant surprises later in the journey and it makes the planning process more enjoyable. Another key finding was that 83% of the respondents had at least once posted an online review. Chatterjee in 2001 analyzed whether customers actually look into online reviews and how they influence their purchasing decisions. In particular the impact of negative reviews was analyzed. It was found that the effect of negative reviews depends on the relationship between the customer and the retailer. In some cases where the customer is familiar with the retailer, the impact is lower. Browning & Sparks in 2011 found out that positive reviews had a greater impact on customers

decision making process for choosing their next hotel. In general, negative reviews tend to lower customer trust and direct them away from the hotel. Interestingly they found that with a few recent positive reviews the negative effect of a set of negative reviews can be neutralized and the customer becomes more acceptable towards the hotel. This is interesting as it implies that customers trust recent reviews more than older ones.

In another piece of work, researcher looked into the possibility of extracting book features from online reviews. These features would then be used to recommend books for other readers (Sohail, Siddiqui, & Ali, 2016). Ullah, Zeb, & Kim in 2015 looked into the helpfulness of reviews. Their aim was to find the factors, which make a review have a greater impact on the user. It was found that reviews, which contained more emotional content, were of more help to the customer in making the purchase decision. A key finding was that negative reviews, which were full with emotional content had no impact on review helpfulness while on the other hand positive reviews with high emotional content had a positive effect in increasing the helpfulness of the review.

## 2.3   What is Opinion Mining and Sentiment Analysis

Opinion mining can be seen as a way of extracting sentiment from raw text. It almost always involves textual analysis of some sort. 'Sentiment analysis, also called opinion mining, is the field of study that analyses people's opinions, sentiments, evaluations, appraisals, attitudes, and emotions towards entities such as products, services, organizations, individuals, issues, events, topics, and their attributes' (Liu, 2012, p. 7). The problem of analyzing people's opinions is not an easy one. It involves understanding how people communicate in the first place. Hence, the understanding of language is a must in order to perform sentiment analysis. Only through the understanding of language and its structures is it possible to calculate opinion from text. For example, to analyze reviews on mobile phones, one needs to understand how people review mobile phones in the first place. This involves compiling list of all words, which could be used to describe mobile phones.

Sentiment analysis and opinion mining is a field of research, which is poised to become more significant in the near future. Perhaps, this field of research has

gained more significance since the Internet became a daily part of people's lives. As the number of people with access to the Internet increases, combined with the ever increasing popularity of websites like Facebook, Twitter, Amazon, etc. the amount of raw data that is generated from these website will increase tremendously. Therefore, it is interesting to know what people are sharing on these social media websites. A study found that the majority of the people use Twitter, a micro-blogging website for sharing what they may like to do during the day or what they may be doing at the moment (Java, Finin, Tseng, & Song, 2007). This indicates that the information that people share on the Internet can be analyzed and meaning can be extracted from it. However, it is a cumbersome task and requires lot of resources and effort from the researcher. Opinion mining is in today's world is quintessential as there now exists a kind of parallel online world, which reflects the real world. Having a significant online presence is essential to every organization. We may be slowly transitioning to a predominantly online world. It may as well be that in the near future the real world would become a reflection of the Internet world. Businesses strive to have a strong positive online image and presence. As people spend large portions of their lives in the online world, businesses do not want to miss out on this opportunity to engage them. Failure to do so puts the business at direct risk of its customers changing brand loyalties.

Opinions are central to peoples' decision-making process. People tend to form opinions on the products before they purchase them. 'Long before awareness of the World Wide Web became widespread, many of us asked our friends to recommend an auto mechanic or to explain who they were planning to vote for in local elections, requested reference letters regarding job applicants from colleagues, or consulted Consumer Reports to decide what dishwasher to buy' (Pang & Lee, Opinion Mining and Sentiment Analysis, 2008, p. 1). It is beneficial to know what other customers think about the product before the purchase decision is made. This has been the case for as long as humanity has existed and will be for the foreseeable future. Knowing others opinion on products or services helps people feel safer about their purchase decisions. People understand the fact that people who have already used the product or service can help them by letting them know how they felt about the product. This helps curb the customers' anxiety before making the purchase and helps them make the purchasing decision more comfortably.

Opinion mining has uses in almost every field and it can be easily seen why this may be the case. Every company would benefit from knowing what its customers think about its products and services. There now exist specialized review based websites, which provide not just reviews from other users, but also expert reviews from professionals. Hence for the average consumer who has accessibility to the Internet, it is not difficult to find out a website which would provide reviews and opinions from other users about the product they are looking to purchase. Feldman & Sanger in 2007 defined three main areas of application for opinion mining. The first being corporate finance where banks are interested in recognising trends to gather business intelligence to provide better services to their customers. Second is the field of patent research where the focus is on identifying development strategies for patents. Finally there is the field of life sciences. However, in reality the application of opinion mining and sentiment analysis are practically limitless. It may also interest not just business but also not for profit organisations like the government. Practically any entity, which is interested in finding out the sentiment associated with something will find opinion mining and sentiment analysis a very useful tool.

It is also interesting to talk about how people used to form opinion during the pre-internet era. When people needed to form opinion they used to talk to friends and family, whereas nowadays there exists a trove of information on the Internet. Similarly, organisations may not need to conduct extensive research studies to gather public opinion. Opinion mining has also been used extensively in recent years to predict public opinion before elections.  Hence, it becomes clear that the availability of such large amounts of raw data on social media websites has completely changed peoples' way of gathering opinion and purchasing products. They now find themselves with a large amount of options, which did not exist a few decades ago.

There has also been a significant amount of research conducted on this topic in the recent years. For example, Pak & Paroubek in 2010 used Twitter to extract relevant information for analysis and assign value to the sentiment attached to the data using linguistic analysis. Tetlock in 2007 found that correlation exists between stock markets and high media pessimism, where exceptionally high or low  pessimism influenced high market trading. Tumasjan, Sprenger, Sandner, & Welpe in 2010

analysed raw data on twitter which had something to do with politicians or political parties, and found that only with the amount of Twitter posts the election results could be predicted. It is clear that the applications of opinion mining are limitless.

## 2.4   Sentiment calculation and Aspect-based analysis

Reviews are made up of sentences and each sentence contains information on the aspects and the sentiment associated with it. It can be seen from this definition of reviews that there arise two different tasks associated with sentiment analysis namely, overall sentiment calculation and aspect-based sentiment analysis. Opinion extraction simply focuses on calculating the overall sentiment present in a piece of text.

Document level sentiment analysis involves calculating overall sentiment score from a piece of text, this is done by assigning sentiment to each sentence in the document and then calculating an overall score based on individual sentence scores. 'One of the main challenges for document-level sentiment categorization is that not every part of the document is equally informative for inferring the sentiment of the whole document. Objective statements interleaved with the subjective statements can be confusing for learning methods, and subjective statements with conflicting sentiment further complicate the document categorization task' (Yessenalina, Yue, & Cardie, 2010, p. 1046). Hence when it comes to document level sentiment analysis a key aspect is the process of eliminating objective sentences, which do not contain any opinion (Pang & Lee, A sentimental education: Sentiment analysis using subjectivity summarization based on minimum cuts, 2004). The text could be an entire document, a small review or a tweet.

The second task that comes into view is the task of aspect-based sentiment analysis. This is different from the first task as it does not focus on a piece of text as a whole but rather looks at each sentence and extracts entity-aspect based characteristics and the sentiment associated with them. 'Aspect, also called feature, is usually the property or function of the product, for example, price, degerm, and moisturizing are the aspects of the body wash products (Zhang, Xu, & Wan, Beijing, p. 10283). Jo & Oh in 2011 defined an aspect as 'a multinomial distribution over words that

represents a more specific topic in reviews' (p. 816). Aspect-based sentiment analysis has been discussed in detail later in this study.

## 2.5 What Defines an Opinion

Before moving on to discuss the tasks involved in aspect-based opinion mining in detail, what constitutes an opinion must be defined. A few random sentences are used in order to be able to clearly define an opinion. Consider the following sentence, 'I purchased this phone a month ago', in this sentence, there exists no opinion, or in other words, it is an objective sentence. Looking at another sentence, 'I really like it' it can be seen that the user truly likes the purchased phone. This sentence is a subjective sentence as it contains sentiment or opinion of the user. Moreover, it tells us that the user has enjoyed his experience with the phone as a whole. A sentence can contain both positive and negative emotions. This sentence expresses positive emotion of the user associated with the phone. Considering another sentence, 'The battery life is amazing'. This sentence shows that the user holds positive emotion towards the battery life of the purchased phone. 'The build quality is also not bad', this sentence shows that the user is talking about the build quality of the purchased phone and he expresses positive emotion about it. Lastly, looking at the sentence 'But I don't like the weight of the device', it can be conferred that the user does not appreciate the heavy weight of the purchased phone and expresses negative emotion in their sentence. It is important to note that in this particular sentence the user refers to the purchased phone as 'device'. This is a common feature among many online reviews, where users use different words for referring to the same purchased object. Another key finding from looking at these sentences is that the user first talks about the phone as a whole and in latter sentences mentions particular aspects of the phone like 'battery life' and 'weight'. Two key observations can be made from these sentences, that an online review may contain sentiment about the object as a whole, as well as about individual aspects about the object, which they may or may not have liked. The phone as well as its battery-life and weight are discussed in the above sentences.

Furthermore, most online reviews also contain information about the reviewer such as their name. Opinion holder is either a person or an organization that holds the opinion (Chu, 2014). Moreover, most online reviews also contain information about

the date on which the review was posted. This is relevant with regard to online reviews as opinions may change over time about a product or service. Using the above set of sentences, it is possible to define opinion as something that contains 4 key components or a quadruple. The four things are namely, topic, sentiment, the opinion holder and the time of observation. It is difficult to fully describe the main topic as, it sometimes may not appear in the sentence itself. In the sentence 'The battery life is amazing', the opinion holder describes the battery life of the phone, which is the topic but the sentence only mentions battery life. Therefore, it would not be possible to use this sentence alone, if it had not been known that the user is actually talking about the battery life of the purchased phone and the sentiment in this sentence would be of no use. Hence, the topic must be defined in a structural way at more than just one level. This helps improve the task of mining raw data for opinion and also further using it for analysis. For example, the battery life of the purchased phone can be represented as an entity and an attribute of the entity represented pairwise, (phone, battery-life). An entity contains a component and an attribute (Chu, 2014). Furthermore, each aspect can also have its own sub-aspects for example battery weight, size. Following this line of reasoning the original definition of opinion mining is slightly modified. It is now contains 5 key aspects namely, entity, entity-aspect, sentiment, opinion owner and date.

In another similar piece of work by Aggarwal & Zhai, Mining Text Data in 2012, opinion is defined as a quintuple (p. 416). They further go on to state that their definition serves as a basis of converting structured data into unstructured data. However, they clearly state that this definition should not be considered final and other items could be added to this definition as per the need of the study. For example, the age and gender may be interesting to look at for each reviewer.

## 2.6   Different Types of Opinions

It is important to mention that opinions can be of different types. A regular opinion is the simplest form, where the opinion is expressed as how the author feels. Aggarwal & Zhai, Mining Text Data in 2012 defined regular opinion as a positive or negative sentiment about an entity. Furthermore, an opinion can be expressed directly or in-directly. It is easier to assign sentiment value to directly expressed opinions than indirectly expressed ones (Liu, 2012). Direct opinions are easier to

handle and therefore majority of the research has focused on that fact. An opinion can also be expressed comparatively by saying entity A is better than B. 'A comparative opinion is usually expressed with the use of the comparative or superlative form of an adjective or adverb' (Pozzi, Fersini, Messina, & Liu, 2017, p. 7). Opinions can also be explicit and implicit. Explicit opinions deal with subjective statements that give regular or comparative opinion. For example, 'IPhone works great' and 'The IPhone is better than the Pixel'. On the other hand, implicit opinions are objective statements that imply regular or comparative opinion. For example, 'I bought the IPhone a month ago and already had to go to the service store thrice' and 'The battery of an IPhone is better that the Pixel'.

## 2.7    Different Levels of Opinion Mining

Opinion mining has been investigated at three major levels in research as already mentioned by (Liu, 2012) and (Aggarwal & Zhai, Mining Text Data, 2012).

i.    Document level opinion mining involves classifying opinion based on the entire document. The focus is not on individual sentences, rather on what the document expresses as a whole. An important pre-requisite for this type of opinion mining is that each document can only express opinion about a single entity. For example, a single product or a campaign.

ii.   Sentence level opinion mining involves classifying opinion contained in individual sentences. This kind of opinion mining caters to the fact that opinion mining is closely related to subjectivity. As already mentioned sentences are usually of two types objective and subjective. Objective sentences merely state a fact, whereas subjective sentences contain opinion of the reviewer. This kind of opinion mining is mostly interested in subjective sentences. However, we cannot say with certainty that objective sentences cannot imply opinion. For example, 'I purchases an IPhone last month and it's already broken merely states facts, but indirectly implies that the IPhone may not be highly durable.

iii.  Entity and aspect level opinion mining tend to divert us away from what the users really like or did not like. Aspect level opinion mining is a way of looking at the fine details. Instead of focusing on the language side of the story, aspect level focuses directly on the opinion itself. It focuses on the

opinion as well as the target of the opinion. An opinion without a target serves no purpose. A sentence may also refer to two aspects at the same time. For example, 'The IPhone has a good battery life, but it has a poor image quality. In this case, the sentence tells us about the positive sentiment associated with the battery life and the negative sentiment associated with the image quality.

## 2.7.1 Document level opinion mining

The topic of document level opinion mining is discussed in detail in this section. The goal of document level opinion mining is to classify an entire document and determine the sentiment associated with it. The important thing to note here is that document level opinion mining focuses on the entire document to calculate emotion. The emotion associated with the document can be expresses as positive or negative, on the other hand the emotion can also be expressed on a scale (ex. 1 to 5). In the first case we focus on classifying if the emotion is positive or negative. In the second case of emotion being expressed as a scale we can also perform regression to find causality. A document may contain opinion on just one entity or on more than one entity. In order to be able to extract sentiment from the document we need to assume that the entire document focuses on just one entity and holds the opinion of just one opinion holder. This however may not hold true in reality and a document may not just express the opinion of just one opinion holder. Moreover, the opinion of one opinion holder may be different from another. This plays a role in the case of online reviews on websites like Amazon, Tripadvisor, Yelp, etc. These reviews express opinions only on one entity like a product on Amazon, a hotel on Tripadvisor or a restaurant on Yelp. Each review is usually about a single product like the IPhone and is written by just one person. However online opinion does not just exist in the form of online reviews, other kinds like blogs and forum posts also need to be considered. The problem is that blogs and forum posts contain a lot of information from different sources.

Another key issue arises when it comes to document level opinion mining. It is understandable that any sentiment classification technique employed will only work well for the domain, which it is being applied to. The reason for this is clear, words and phrases can imply different things in different domains. One word in a domain

may not even have the same meaning in another domain. Researchers have explored the topic of using knowledge from one domain and employing it to gain knowledge about a different domain. Yang, Si, & Callan in 2006 used knowledge from the domains of movie reviews and product reviews to generate a dataset for a blog. Their approach was straightforward, using the most commonly occurring words in both movie and product reviews they were able to generate a domain-independent list of words and phrases. In order to apply knowledge from one domain to another, the features ranked highly in both domains were considered to be domain independent features. They then applied this to the draw sentiment on blogs.

Researchers have also discussed the topic of cross language sentiment analysis in detail. Most of the research done on sentiment analysis so far has been in English. Hence, researchers from other regions of the world want to develop sentiment analysis systems in their own languages. Businesses now have to deal with customers from all over the world, which do not necessarily speak the same language. Such businesses would also like to know their customers opinion. Hence, there is clear need for such systems. Wan in 2008 devised a system to use the vast resources available for sentiment analysis in English and used it to analyze Chinese reviews. Chinese reviews were first translated to English language and then existing English language analysis tools were used. This technique proved to be highly effective.

### 2.7.2 Sentence level opinion mining

This field of opinion mining is not very different from document level opinion mining as here one sentence is simply treated as an entire document. This process involves classifying each sentiment based on whether it expresses positive, negative or neutral emotion. The problem of sentiment analysis can be seen as a two-step problem. This is because some sentences may not even contain sentiment at all. Such sentences are called objective sentences. Objective sentences merely state a fact. The focus here is on extracting sentiment from subjective sentences, which are ripe in sentiment. Therefore the first step is to classify whether a sentence is subjective or not. The next step then involves calculating sentiment.

### 2.7.3 Aspect-based opinion mining

In aspect based sentiment extraction the focus is on first defining entities such as a product and then the opinion about different aspects of the product is extracted and finally positive or negative values for the sentiment associated with the aspect are assigned. Another final step is summarizing the opinions, as the goal is to make sense of the opinions of many and not just the few.  Pontiki, Galanis, Pavlopoulos, Papageorgiou, Androutsopoulos, & Manandhar in 2014 used aspect based sentiment analysis to identify the aspects of the entities being reviewed and to determine the sentiment the reviewers express for each aspect. They devised a stepwise procedure to first identify the aspects from the sentences for example, food, service and ambience in relation to restaurants. Finally, they assigned polarities to each aspect.

Hence, the task of aspect-based opinion extraction is clear. First, the entity, which is being investigated, needs to be defined. An entity could be a product, a campaign for a political party or an organization. Then the aspects and the sentiment associated with them needs to be extracted. Third, the words need to be assigned strength in terms of the positive or negative sentiment they posses. Lastly, the polarities for different aspects need to be summarized. There are two kinds of sentences that express sentiment, objective and subjective sentiments. Objective sentences tend to focus on facts about the world we live in whereas subjective sentences express feelings and beliefs. It is important to note that although sentiment can be extracted from any kind of sentence, it is better to use only subjective sentences as they express sentiment explicitly which is what the goal is in the first place.

Having talked about the steps involved in aspect-based opinion mining, lets move on to discussing the tasks involved in opinion mining in detail. The tasks stem from the definition of opinion as a containing 5 key aspects. The first task is entity extraction. The extracted entities are then categorized. One clear problem that emerges here is the fact that people use different words to represent an entity or an emotion. For example, the IPhone may be referred to as Apple's latest flagship device, apple phone, ifone and so on. A mechanism to make sure that they all target towards the same thing needs to be devised.

Next is the problem of entity aspect extraction. Same as the case with entities, aspects also need to be categorized and each aspect category can only represent one aspect. Aspects can be represented in real text through nouns and noun phrases, as well as verbs, verb phrases, adjectives, and adverbs. Aspects that are represented through noun and noun phrases are called explicit aspects. For example, battery life in the sentence 'The battery life is amazing' is an explicit aspect expression. Aspects that are represented in real text that are not nouns and noun phrases are called implicit aspects (Hu & Liu, 2004). Lets try to explain this with the help of an example. In the sentence 'The IPhone is beautiful' the word beautiful is an implicit aspect. It implies to the aspect of looks. A majority of implicit aspects are adjectives and adverbs being used to describe or recognize certain aspects, e.g., beautiful (looks). In a number of cases we find that they can also be verb and verb phrases, e.g., 'I can use the IPhone easily' use indicated the aspect usability. Implicit aspects can sometimes also be expressed in a complicated way in real text. For example, 'The IPhone will crack easily when it falls'. Here 'crack easily' points to the aspect of durability.

The third task involves determining whether the sentiment is positive, negative or non-existent (Liu, 2012). Furthermore, the strength of the sentiment also needs to be determined, which helps later during the analysis. The fourth task is extracting the identity of the opinion holder from the text. Last but not the least, the date when the opinion was expressed also needs to be extracted. In the case of product reviews and blogs, opinion holders are the authors of the respective posts. Opinion holders have high significance when it comes to news articles as they directly point out to the holder of the opinion. The opinion can also be of the organization as a whole. After having discussed the tasks involved in opinion mining in detail, a brief summary of the tasks involved in opinion mining is presented. It becomes clear that it involves 6 significant tasks. These tasks are coherent with the research done by (Liu, 2012); (Pang & Lee, Opinion Mining and Sentiment Analysis, 2008);

i. Extract and categorize all the entities present in the text. Each entity category represents a unique entity.

ii. Extract and categorize all aspects present in the text. Each aspect category represents a unique aspect.

iii.    Extract and categorize the opinion holder from the text. This task is similar to the above two tasks.

iv.    Extract the time when the opinion was made. Here it is important to record the findings in a standardized time format for easy analysis in the later stage.

v.    Extract the opinion present in the text. This can be as positive, negative or neutral or as an integer value.

vi.    Generate the final quintuple for the opinion as defined in the definition based on the results of the above 5 tasks.

Opinion summarization is the final step, which results in a kind of report, which helps clearly infer the opinion of a large group of people on different aspects or features of the entity in question. 'One opinion from a single person is usually not sufficient for action. This indicates that some form of summary of opinions is desirable' (Aggarwal & Zhai, Mining Text Data, 2012, p. 421). For example aspects of an IPhone could be battery life, camera performance, general impression, the build quality, weight, durability and so on. Summarization allows us to assign a score to the different aspects and thereby determine the opinion associated with it.

## 2.8   Opinion Spam

Business owners understand that having a high online rating leads directly to more business for them and hence it is understandable that they may want to add more positive reviews on the Internet about their product and services. They may do this by asking their employees or friends to post fake reviews. The goal is to spam the Internet with false positive reviews and try to sway the general opinion in their favor. Although this is ethically wrong and may cause potential harm to the credibility of online review websites, it is still practiced online. The people who engage in such activities are called opinion spammers. Pozzi, Fersini, Messina, & Liu in 2017 discussed in detail various techniques available to detect opinion spam. The problem of opinion spamming is increasingly evident in the world of online reviews. Fake reviews are posted in order to benefit a particular business like a restaurant or a movie hall. These reviews can help sway public opinion in order to help the business being reviewed. Fake reviews are very difficult to spot and opinion spam poses a problem that has not been seen with other forms of spam on the Internet like email spam (Aggarwal & Zhai, Mining Text Data, 2012): (Liu, 2012).

Malbon in 2013 proposed five steps that can be taken to eliminate the menace of fake online reviews. The first step involves identifying clear objectives, which in the case of online reviews may be to detect fake online reviews and in-turn increase customer trust and market fairness. Step two involves getting together all the businesses, individuals, organization and law enforcement agencies that would be at an advantage if, tough laws and regulations are implemented against fake reviews. Some businesses may feel under pressure to post fake online reviews in order to stay competitive against other businesses. The third step involves, selecting matters of concern to all the parties involved such as customers and businesses. For example, websites containing online reviews may be concerned about their legitimacy which in-turn decreases customer confidence for the website, customers may feel they are being forced to make wrong purchasing decisions through fake online reviews. The fourth step involves clearly deciding the powers given to the group of individuals affected by online fake reviews. For example, businesses may be allowed to privately investigate and identify fake online reviews, customers may choose not to purchase products from an entity engaged in the practice of publishing fake online reviews. Furthermore, businesses may go one step further and report the organizations engaged in the practice of posting fake online reviews to the media and publically shame them. The final step involves assigning powers to individuals to come down strongly against fake online reviews. For example, businesses may be penalized legally for engaging in such activities moreover, a regulator may be setup, which can provide ratings to website which contain online reviews to the extent that what steps are being taken by the website to detect and deter fake online reviews, informational seminars could be conduct to inform customers and businesses about the risks and damages possible by engaging in the activity of posting fake online reviews, small financial incentives to customers and businesses that identify and report fake online reviews could also be provided.

'Even though the existence of online reviews fraud is acknowledged by online vendors, these online vendors rarely discussed publicly how they should fight online reviews fraud. There was no commonly agreed conceptual definition of online reviews fraud based on which vendors could mandate some appropriate legal action. Similar to the case of digital rights management, vendors believed that one way to filter online reviews fraud was to never disclose exactly how they identified

such fraudulent reviews. They had the apprehension that unethical users would take advantage of such disclosures. Due to the above challenges, a method for the determination of existence of manipulation in online reviews is crucial' (Hu, Bose, Koh, & Liu, 2012, p. 676).

## 2.9 Sentiment Lexicon

### 2.9.1 What is it?

Sentiment words also called opinion words are most commonly used to express positive or negative opinion. 'Opinion words are words that are commonly used to express positive or negative opinions (or sentiments)' (Ding, Liu, & Yu, 2008, p. 1). For example, awesome, great, perfect, wonderful, mind-blowing are words associated with positive opinion whereas, bad, horrible, awful are words associated with negative opinion. Not only words but also idioms and phrases can be used to express opinion e.g., 'Buying an IPhone made a deep hole in my pocket' expresses that the IPhone is an expensive product according to the reviewer. It becomes obvious that such sentiment words or opinion words are critical to opinion analysis. When put together in a list, such a list is called an opinion lexicon.

### 2.9.2 Problems with Sentiment Lexicon Analysis

Sentiment Lexicon analysis comes with its own set of problems. Several of them are pointed out below as already discussed by (Liu, 2012):

  i. Sentiment words can have different meaning associated to them in different scenarios.
 ii. Sometimes sentences contain sentiment words but fail to express an opinion at all. This occurs mostly in the case of questions and conditional sentences.
iii. Sarcastic sentences mostly found in political scenarios are difficult to analyze for opinion.
iv. Objective sentences, which do not contain any sentiment words, can also contain opinion in certain cases.

### 2.9.3 Sentiment Lexicon generation

There are few major approaches to generating sentiment lexicon lists. These are discussed them briefly below as already mentioned by (Liu, 2012):

i.   The manual approach is the most time and labor intensive. It is normally conducted by a group of people. It involves a group of researchers looking into reviews and trying to gather a list of words used to describe products and services. Secondly, each word must be assigned with a positive or negative polarity. Finally, researchers can then calculate the overall sentiment using the generated list of lexicons.

ii.  The dictionary-based approach focuses on compiling word lists based on synonyms and antonyms for each word. This process starts by compiling a small list of words with known positive or negative emotion and an algorithm is used to expand this list using online dictionaries. This is a step-wise process, the algorithm is fed a list of words and each step includes the words found by the algorithm in the previous step. 'The approach generally uses a dictionary of opinion words to identify and determine sentiment orientation (positive, negative or neutral). The dictionary is called the opinion lexicon' (Zhang, Ghosh, Dekhil, Hsu, & Liu, 2011, p. 2). There are various online dictionaries available for this like Wordnet and Sentiwordnet.

iii. The corpus-based approach involves making the use of linguistic rules to predict sentiment. The corpus-based is different from the dictionary-based approach as it builds a new corpus of words with their respective negative or positive sentiment through projection (Denecke, 2008). Corpus based approaches involve compiling various online reviews from a domain and then looking for words which relate to sentiment and opinion. This list of words is then assigned polarities of negative or positive sentiment (Prekopcsak, Makrai, Henk, & Ǵasṕar-Papanek, 2011). A corpus-based approach also allows researchers to use the corpus from a specific domain over and over for various sets of online reviews.

# 3 Tripadvisor, Rapidminer and 25hours hotel

## 3.1 Tripadvisor as a website

Tripadvisor claims to be the world's largest travel site. 'Tripadvisor offers advice from millions of travelers and a wide variety of travel choices and planning features with seamless links to booking tools that check hundreds of websites to find the best hotel prices' (Tripadvisor, 2016, p. about). 'For many, TripAdvisor has become a first stop for travel planning. Thanks in part to its prominence in Google searches, some 24 million visitors a month check out what other users have to say about where to stay, eat and play around the world.' (The Wall Street Journal, 2007, p. 1). Tripadvisor contains advice from other users based on their personal experiences. It encompasses hotels, restaurants and the so-called things-to-do activities. People can visit this website to see what other people have said about the product or service which they are interested in. Tripadvisor claims to make travel easier by providing people with advice even before they embark on their journeys.

The reviews published on the Tripadvisor website contain many aspects, for example, the title of the review. This is a short title, which the reviewer must provide. It can be seen as a sort of headline to the actual review text. The review text also contains detailed information about the user's experiences and some advice for others. Moreover, the review also contains a star rating which can be provided by the reviewer measured on a scale of 1-5 with 1 corresponding to terrible, 2 representing poor, 3 corresponding to average, 4 pointing to very good and 5 being excellent. Another key feature in these reviews is that they also mention the kind of traveller. Tripadvisor provides this option to its users by allowing them to mention whether they are travelling solo, as a couple, as a family, as friends, or for business purposes.

According to the Content Integrity Policy of Tripadvisor (2016), Tripadvisor aims to provide a clean playing field to all business involved with the website irrespective of their size. As per Tripadvisor when somebody spams a review means that they themselves or through other people have reviewed their own business positively, or negatively reviewed a competitor's business. They aim to block and penalize such reviews by detecting them with their propriety fraud detection system. Moreover,

when a fake review is detected, this negatively impacts the business' ranking on the website.

A number of studies have been conducted on the topic of Tripadvisor reviews. A study by Vásquez in 2010 analyzed negative Tripadvisor reviews in order to compare them to complaints made through other conventional mediums. Tripadvisor reviews were different from the fact that they contain complaints with a sense of advice and recommendations whereas conventional complaints are associated with threats and warnings. Wu, Greene, Cunningham, & Smyth in 2010 analyzed the effect of the Tripadvisor website on the hotel industry in Las Vegas. Tuominen in 2011 used reviews on Tripadvisor to check if they affected hotel performance in any way. Their aim was to see if the advent of Tripadvisor had any difference in how the hotels changed their strategies or way of operating.

Data from Tripadvisor can be used to analyze reviews in a variety of ways, for example, Schuckert, Liu, & Law in 2016 devised a method to look into the aspect of suspicious fake reviews on the Tripadvisor website. Tripadvisor allows customers to review a hotel on an overall basis as well as individual basis. They found that in the case of fake online reviews, there was a significant difference in the overall rating and the specific rating. Moreover, service emerged as the most discussed topic in reviews, which implies that hotels need to pay more attention on improving the service quality of their hotels in order to gain higher online ratings. It was also found that about 20% of the reviews were suspicious and possibly fake. This is alarming as nowadays travellers rely heavily on online reviews to make travel decisions. Another key finding in their research was that the problem of fake reviews is more significant for lower-class hotels. This shows that travellers need to pay more attention when going through online reviews to make their travel decisions, moreover better mechanisms need to be designed by websites such as Tripadvisor and Amazon in order to detect and remove suspicious fake reviews.

In a study by Banerjee & Chua in 2016 reviews from the Tripadvisor website were analyzed based on travelers' profiles and regions which hotels belong to. The difference was also analyzed based on independent and chain hotels. It was found that in general chain hotels were rated relatively higher than independent hotels. Moreover, chain hotels also tended to receive a higher volume of ratings as

compared to independent ones. However, in the Asia-Pacific region independent hotels were found to get a higher volume of ratings than chain hotels. Differences were also found across regions for example, independent hotels were rated higher in the Europe region whereas they lagged behind in the Asia-Pacific region. Also, couples tended to be more lenient when it comes to rating hotels online. Business travellers on the other hand tended to be stricter when it comes to handing out ratings.

Another topic of significance is the trustworthiness of online reviews posted on travel websites. It is important to discuss their reliability as travelers rely heavily on them to make their travel decisions. In a study by Chua & Banerjee in 2013 this particular topic was discussed in detail. Their methodology comprised of looking into multiple hotels reviewed by the same individual. It was found that users' ratings were consistent across various hotels and hence online reviews should be considered largely reliable. However, there were some reviews, which were not consistent, as they seemed to have high overall ratings, however their content mentioned many negative aspects about the hotel.

In order to understand what motivates people to post online reviews, (Ögut & Cezar, 2012) conducted a study. They found that a higher rating and lower price motivates people to write reviews. However, a higher star rating is not a significant factor to motivate people to write reviews. A key finding was that people are more motivated to write reviews when the number of reviews already present for the hotel is higher. This implies that hotels with a higher number of reviews are more likely to attract online reviews. This has implications for small independent hotels, which do not get a high number of reviews when compared to larger chain hotels. Thus, small hotels must try to create strategies, which attract more reviews from customers online in order to compete with the larger chain hotels.

## 3.2    Rapidminer as a software

Rapidminer is an open source, data science software, which allows users to perform data analysis tasks. It aims to make the process easier by offering an intuitive GUI interface (Rapidminer Inc., 2016). An interesting feature of this software is that it allows for connection to third party Application Program Interfaces (APIs). In this

study, Rapidminer was used to analyze the data for sentiment as well as to perform aspect-based sentiment analysis. The third party API used for text mining is called 'Text Analysis by Aylien'. 'AYLIEN Text Analysis API is a package of Natural Language Processing, Information Retrieval and Machine Learning tools for extracting meaning and insight from textual and visual content with ease.' (Aylien Text Analysis, 2014). 'Rapidminer enables one to design data mining processes by simple drag and drop of boxes representing functional modules called operators into the process, to define data flows by simply connecting these boxes, to define even complex and nested control flows, and all without programming. Rapidminer stores the data mining processes in machine-readable XML format, which is directly executable in Rapidminer with the click of a button, and which along with the graphical visualization of the data mining process and the data flow serves as an automatically generated documentation of the data mining process, makes it easy to execute, to validate, to automatically optimize, to reproduce, and to automate (Hoffman & Klinkenberg, 2014, p. 25).

Rapidminer has been used in many recent studies. For example, Prekopcsak, Makrai, Henk, & Gaspar-Papanek in 2011 used the Hadoop extension of Rapidminer for data transformation tasks. They found out that Rapidminer was a good tool for analyzing big data. The Apache Hadoop project is an open-source software for processing large amounts of data. The advantage of Apache Hadoop is that it allows for processing of data to run on several computers simultaneously. The applications of the Rapidminer software are practically limitless for the field of data mining for example, text processing, web mining and python scripting. Another very good feature of this software is that it allows for integration of various extensions, which further enhance the functionality of the software. Third party APIs allow for even further integration of Rapidminer with other available technologies. In another research by Jungermann in 2009 the information extraction extension of Rapidminer was used, which converts documents containing natural language to machine-readable form, which can then be analyzed.

## 3.3   25hours hotel Vienna

The 25hours hotel located in Vienna, Austria is a 4-star hotel. The hotel is located very close to the city center. It has 217 rooms including 34 suites. Furthermore, the

hotel provides meeting space for about 200 people. A very famous feature of the hotel is its rooftop terrace with a bar. It also offers a variety of very slick restaurants. Underground parking is also available for guests with cars. The hotel distinguishes itself from other hotels by offering a similar feeling as home to its guests. The hotel has an unorthodox style. Moreover, the hotel aims to be a place where guests come not just to sleep, but also rather to socialize with other guests. The hotel offers a diverse range of designs, in order to be able to surprise their guests on every visit. This helps the hotel to be the perfect place to click pictures and post them on various social media platforms. The rooms at the hotel aim to be functional yet unique, which makes them stand out from the competition. The hotel also lays great focus on its public areas like the lobby. The hotel also aims to hire staff, which has strong interests outside of their work lives. By doing this, their aim to not let their guests' interaction with their staff become standard and mundane. Overall, the hotel focuses on aiming to make the experience at the hotel a unique one.

# 4 Methodology

## 4.1 Hypothesis setup

The hypothesis setup for this study stems from the initial objectives. The interest lies in finding correlation between the star rating given by the reviewer and the positive, negative or neutral emotion calculated from the review text. It is expected that these two will be positively correlated, as positive emotion is likely to go hand in hand with higher star rating. Hence, the first hypothesis is setup as follows:

H1 – The star rating and the calculated emotion are positively correlated.

Moving on to the next hypothesis. Both the review text and the title text were analyzed for sentiment. The review text contains more words than the title text and therefore it is interesting to find out if there is any correlation between them. It is expected that there will be some correlation between the review text and the title text. Hence the second hypothesis is setup as follows:

H2 – The calculated emotion of the review text and the title text are positively correlated.
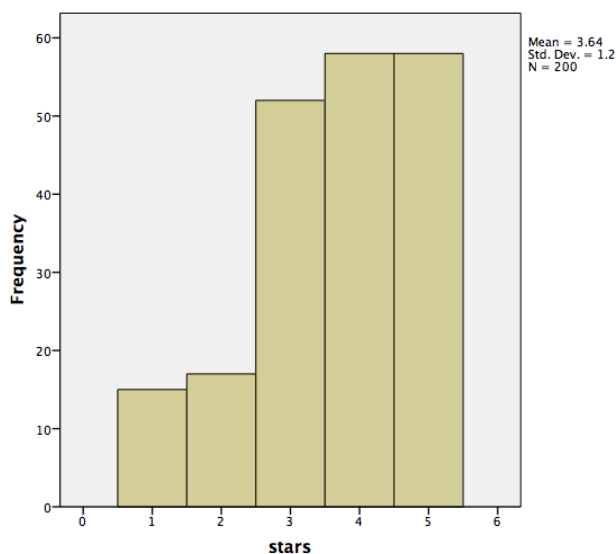
Furthermore, using the data generated from aspect-based sentiment analysis, the aim was to find out which aspects are useful in predicting the star rating. Hence, the hypothesis is setup as follows:

H3 – Aspects are helpful in predicting star rating.
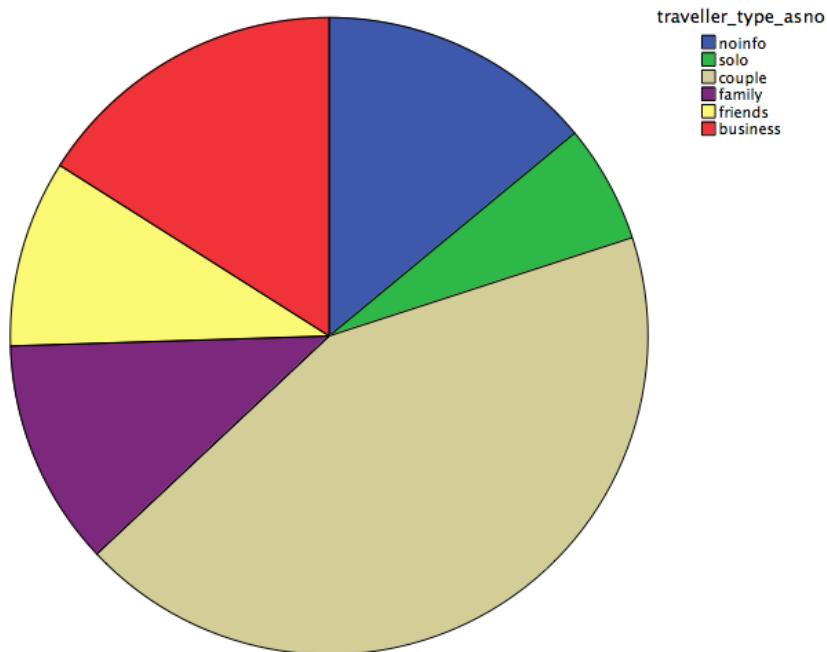
## 4.2 Data collection and analysis

Data in the form of reviews was collected from the Tripadvisor website. As already mentioned the hotel used for collecting reviews was the 25hours hotel located in Vienna, Austria. The number of reviews posted on the website were approximately 1,700 at the time of data collection. In total, 200 reviews were collected from the website. The website additionally asks users to provide a review rating from 1 stars until 5 stars. The collection of 200 reviews contained 58 5-star reviews, 58 4-star reviews, 52-3 star reviews, 17 2-star reviews and 15 1-star reviews.

Figure-1: Star rating



Furthermore, the website allows users to provide information on whether they travelled alone, as a couple, as a family, with friends, or for business purposes. The data set contains comparable number of reviews from all traveller types with couple travellers being more than double than any other traveller type.

Figure-2: Traveller type



The analysis of the reviews through Tripadvisor also generated a variable, which relates to how much confidence can be placed on the emotion that was calculated during the analysis. It was found that the mean confidence in the reviews text analysis was 74.678%.

Table-1, Polarity confidence of the review text

**Descriptives**

| | | | Statistic | Std. Error |
|---|---|---|---|---|
| polarity_confidence | Mean | | .74678 | .015008 |
| | 95% Confidence Interval for Mean | Lower Bound | .71719 | |
| | | Upper Bound | .77638 | |
| | 5% Trimmed Mean | | .75204 | |
| | Median | | .79091 | |
| | Variance | | .045 | |
| | Std. Deviation | | .212251 | |
| | Minimum | | .363 | |
| | Maximum | | 1.000 | |
| | Range | | .637 | |
| | Interquartile Range | | .407 | |
| | Skewness | | −.256 | .172 |
| | Kurtosis | | −1.478 | .342 |

The mean confidence for the analysis of the review title was 60.12%.

Table-2: Polarity confidence of the review title

**Descriptives**

| | | | Statistic | Std. Error |
|---|---|---|---|---|
| polarity_confidence_title | Mean | | .60120 | .013493 |
| | 95% Confidence Interval for Mean | Lower Bound | .57459 | |
| | | Upper Bound | .62780 | |
| | 5% Trimmed Mean | | .59032 | |
| | Median | | .53046 | |
| | Variance | | .036 | |
| | Std. Deviation | | .190826 | |
| | Minimum | | .388 | |
| | Maximum | | 1.000 | |
| | Range | | .611 | |
| | Interquartile Range | | .284 | |
| | Skewness | | .846 | .172 |
| | Kurtosis | | −.585 | .342 |

Further through aspect-based sentiment analysis it was calculated whether the review mentions an aspect and what is the sentiment associated with it. Table summarizing the aspects and the number of times they were mentioned in the reviews is presented below:

Table-3: Frequency of the aspects

| Aspect | Number of times mentioned |
|---|---|
| Cleanliness | 54 |
| Comfort | 55 |
| Room amenities | 163 |
| Quietness | 30 |
| WiFi | 20 |
| View | 54 |
| Beds | 55 |
| Value | 72 |
| Facilities | 136 |
| Design | 93 |
| Customer support | 35 |
| Location | 125 |
| Payment | 29 |
| Food & Drinks | 161 |
| Staff | 160 |

It can be seen that the features room amenities, facilities, location, food & drinks and staff are the ones that are mentioned the most. Hence, it can be conferred that these aspects are the most important when it comes to influencing the opinion of the customer about the hotel.

Lastly, taking the aspect-based sentiment analysis one step further, the various aspects mentioned by different categories of travellers in the reviews were looked into. This allowed for a more fine grain analysis by identifying the most frequently mentioned aspects by different types of travellers. This was done through SPSS software and the frequencies of mentioned aspects based on traveller type were analyzed. A brief summary of top 3 most frequently mentioned aspects by different categories of reviewers is presented below:

Figure-3: Aspects and their occurrence in reviews

| Traveller type | Aspects and their frequencies | | |
|---|---|---|---|
| Businesses | Food and drinks (24) | Room amenities (23) | Facilities (18) |
| Couple | Room amenities (75) | Food and drinks (75) | Staff (69) |
| Family | Staff (17) | Room amenities (17) | Food and drinks (16) |
| Friends | Food and drinks (17) | Staff (15) | Facilities (14) |

Some of the most frequently mentioned aspects match with the findings of Levy, Duan and Boo's research in 2013.

# 5   Results

Analysis on the data using SPSS software was performed. Starting with the first hypothesis (H1), the Spearman correlation was calculated.

Figure-4: Spearman's correlation for star rating and sentiment

**Correlations**

| | | | stars | polarity_asno |
|---|---|---|---|---|
| Spearman's rho | stars | Correlation Coefficient | 1.000 | .273[**] |
| | | Sig. (1–tailed) | . | .000 |
| | | N | 200 | 200 |
| | polarity_asno | Correlation Coefficient | .273[**] | 1.000 |
| | | Sig. (1–tailed) | .000 | . |
| | | N | 200 | 200 |

**. Correlation is significant at the 0.01 level (1–tailed).

It is important to note that the strength of both the coefficients is not high. Hence it can be said that there exists a weak correlation. However, since a significant result was found, H0 is rejected and H1 is accepted. That is, the star rating and the emotion calculated are positively correlated. Hence, a positive emotion is likely to get a higher star rating.

Moving on to the second hypothesis (H2). Spearman's coefficient was calculated.

Figure-5: Spearman's correlation for sentiment of the review text and of the title text

**Correlations**

| | | | polarity_asno | polarity_title_asno |
|---|---|---|---|---|
| Spearman's rho | polarity_asno | Correlation Coefficient | 1.000 | .214[**] |
| | | Sig. (1–tailed) | . | .001 |
| | | N | 200 | 200 |
| | polarity_title_asno | Correlation Coefficient | .214[**] | 1.000 |
| | | Sig. (1–tailed) | .001 | . |
| | | N | 200 | 200 |

**. Correlation is significant at the 0.01 level (1–tailed).

It can be seen that the above results are significant hence H0 is rejected and H2 is accepted. Therefore, there exists a positive correlation between the sentiment associated with the review text and the title text. However, this correlation is weak.

For analyzing hypothesis 3 (H3), linear regression analysis using SPSS was conducted. As already mentioned previously, the most frequently mentioned aspects room amenities, facilities, location, food & drinks and staff were used for the linear regression analysis. After the first run, significant results were found for variables room amenities, food and drinks and staff. Facilities and location were non-

significant. Hence the linear regression was re-run however, this time excluding location. However, the facilities variable was still found non-significant. Lastly the facilities variable was also excluded and the final results are present below:

Figure-6: Regression analysis

**Model Summary**

| Model | R | R Square | Adjusted R Square | Std. Error of the Estimate |
|---|---|---|---|---|
| 1 | .738[a] | .545 | .533 | .769 |

a. Predictors: (Constant), Staff, Room_amenities, Food_Drinks

**ANOVA[a]**

| Model | | Sum of Squares | df | Mean Square | F | Sig. |
|---|---|---|---|---|---|---|
| 1 | Regression | 80.740 | 3 | 26.913 | 45.569 | .000[b] |
| | Residual | 67.328 | 114 | .591 | | |
| | Total | 148.068 | 117 | | | |

a. Dependent Variable: stars
b. Predictors: (Constant), Staff, Room_amenities, Food_Drinks

**Coefficients[a]**

| Model | | Unstandardized Coefficients | | Standardized Coefficients | t | Sig. |
|---|---|---|---|---|---|---|
| | | B | Std. Error | Beta | | |
| 1 | (Constant) | 3.271 | .077 | | 42.260 | .000 |
| | Room_amenities | .307 | .094 | .255 | 3.277 | .001 |
| | Food_Drinks | .263 | .101 | .207 | 2.598 | .011 |
| | Staff | .505 | .101 | .415 | 4.983 | .000 |

a. Dependent Variable: stars

The above results show that the overall model is of moderate strength with the adjusted-R square value of 53.3%. The overall model is significant as can be seen from the significance value (ANOVA p-value<0.001). Hence, H0 is rejected and H3 is accepted that is the variables are helpful in predicting the star rating. From the above results it can be said that room amenities, food & drinks and staff are very helpful in predicting the star rating whereby, staff has the highest influence (0,415).

Since, staff and service go hand in hand it can be said that this result is in coherence with the findings of Tsaoa, Hsieh, Shih & Lin's findings in 2015 and that of Schuckert, Liu & Law's in 2016.

# 6 Discussion and Conclusion

The outcome of this study shows that the Rapidminer software can be considered useful for sentiment analysis. There exist also other data mining software available but Rapidminer is one of the few ones, which are available for free. The first hypothesis shows that when the emotions of the customer are impacted in a positive way, the customer is likely to give a higher positive rating. Moving on to the second hypothesis it can be conferred that the title of the review and the text of the review are somewhat correlated. However, using only the title may not always be the correct method to predict the sentiment associated with review text. It is important to consider the full text of the review for sentiment analysis and aspect extraction. These results emphasize that businesses must pay heed to the entire online review.

Finally, from the third hypothesis it can be conferred that aspect-based sentiment analysis is useful in predicting the causes for a good or a bad review. Companies can use this information to fine-tune their services for their customers. Given that a good review is a direct reflection of customer satisfaction, businesses can perform aspect-based sentiment analysis to see which aspects of their products and services are truly effective in increasing customer satisfaction. Businesses, which are looking to find out what their customers think about their products, can not only analyze the overall satisfaction level of the public but also go into the minute details using aspect-based sentiment analysis to find out features in their products and services which customers find exceptionally appealing or unappealing. They can then build strategies to improve the aspects, which customers find particularly annoying.

Using descriptive statistics it was possible to find out the aspects, which were of the highest concern to different categories of travellers. This allows hotels to further fine-tune their offerings based on the type of traveller in order to increase their satisfaction level even further. For example, hotel facilities and food and drinks

emerged as the number one concern for business travellers, which is understandable as they mostly stay in the hotel for shorter periods of time and due to their nature of travel require good facilities in like wireless internet, printers, etc. Couples mentioned room amenities, food and drinks and staff most frequently in the reviews. People travelling as a group of friends mentioned food and drinks, staff and facilities most frequently. Lastly, travellers, which travelled as a family, mentioned staff, room amenities and food and drinks most often.

This has huge implications for businesses particularly hotels, as they can delve deeper into which factors customers talk about most in online reviews, which allows them to work with precision on those particular aspects. Such kind of in-depth analysis was not possible in the past, moreover this analysis can be conducted in short time frame allowing hotels to quickly adjust their strategies accordingly to customer complaints, which eventually allows them to avoid getting negative reviews in the first place. This research also shows that large amounts of resources need not be invested into opinion mining and sentiment analysis and can be done using minimal resources. It would be beneficial for hotels to have specific departments whose only task is to gather intelligence from online reviews.

# 7  Limitations

The study was conducted on a sample of 200 reviews from a corpus of about 1700 reviews. Hence, the sample representativeness may not allow for results to be generalized for the entire set of reviews. Moreover, there were more reviews taken in the 3 to 5-star rating than the 1 to 2-star rating. This was because of the low availability of the reviews in the 1 to 2-star rating. This may cause less negatively associated aspects to be captured by the software. Sentiment calculation is a difficult task and therefore there may be cases where the sentiment may not have been analyzed correctly by the software due to limited linguistic capabilities.

# 8  Implications for future research

This study demonstrates that using limited resources one is able to calculate sentiment through textual analysis. Opinion mining emerges as a very useful

technique for businesses looking to gather information on customer satisfaction about their products and services. Furthermore, aspect-based sentiment analysis is highly useful for businesses looking to fine-tune their products and services. Companies can analyze customer perception of their products and services based on individual aspects, which the customers may or may not like. Opinion mining through textual analysis can be applied to any product, service, political campaign, public awareness campaigns, etc. An important managerial implication for Hotels emerges that they must focus more on increasing the quality of their service, as this will allow them to have greater online ratings, which is equivalent to high sales figures.

# 9   Bibliography

Prekopcsak, Z., Makrai, G., Henk, T., & Gaspar-Papanek, C. (2011). Radoop: Analyzing Big Data with RapidMiner and Hadoop. *In Proceedings of the 2nd RapidMiner community meeting and conference* (RCOMM 2011) (pp. 865-874)..

Ullah, R., Zeb, A., & Kim, W. (2015). The impact of emotions on the helpfulness of movie reviews. *Journal of applied research and technology , 13* (3), 359-363.

Zhu, F., & Zhang, X. M. (2010). Impact of Online Consumer Reviews on Sales: The Moderating Role of Product and Consumer Characteristics. *Journal of marketing , 74* (2), 133-148.

Zhang, W., Xu, H., & Wan, W. (Beijing). Weakness Finder: Find product weakness from Chinese reviews by using aspects based sentiment analysis . *Expert Systems with Applications , 39* (11), 10283-10291.

Zhang, L., Ghosh, R., Dekhil, M., Hsu, M., & Liu, B. (2011). Combining Lexicon-based and Learning-based Methods for Twitter Sentiment Analysis. *International Journal of Electronics, Communication and Soft Computing Science & Engineering (IJECSCSE) , 89*, 1-8.

Wu, G., Greene, D., Cunningham, P., & Smyth, B. (2010, December 8). Does TripAdvisor Makes Hotels Better? Belfield, Dublin, Ireland.

www.internetlivestats.com. (2016, December 5). *www.internetlivestats.com*. Retrieved December 5, 2016, from www.internetlivestats.com

Wan, X. (2008). Using Bilingual Knowledge and Ensemble Techniques for Unsupervised Chinese Sentiment Analysis. *In Proceedings of the conference on empirical methods in natural language processing* (pp. 553-561). Association for Computational Linguistics.

Yang, H., Si, L., & Callan, J. (2006). Knowledge Transfer and Opinion Detection in the TREC2006 Blog Track. *In TREC* .

Ye, Q., Law, R., & Gu, B. (2009). The impact of online user reviews on hotel room sales. *International Journal of Hospitality Management , 28* (1), 180-182.

Yessenalina, A., Yue, Y., & Cardie, C. (2010). Multi-level structured models for document-level sentiment classification. *In Proceedings of the 2010 Conference on Empirical Methods in Natural Language Processing* (pp. 1046–1056). Association for Computational Linguistics.

Yoo, K.-H., & Gretzel, U. (2009). Comparison of Deceptive and Truthful Travel Reviews. *Information and communication technologies in tourism 2009* , 37-47.

Vásquez, C. (2010). Complaints online: The case of TripAdvisor. *Journal of Pragmatics , 43* (6), 1707-1717.

Vermeulen, I. E., & Seegers, D. (2008). Tried and tested: The impact of online hotel reviews on consumer consideration. *Tourism management , 30* (1), 123-127.

Aylien Text Analysis. (2014). *Aylien Intelligence*. Retrieved December 7, 2016, from http://aylien.com

Aggarwal, C. C., & Zhai, C. (2012). *Mining Text Data.* New York, USA: Springer Science+Business Media.

Aggarwal, C. C., & Zhai, C. (2012). *Mining Text Data.* New York, USA: Springer.

Burgess, R. (2013, March 20). *One minute on the Internet: 640TB data transferred, 100k tweets, 204 million e-mails sent*. Retrieved Dec 8, 2016, from http://www.techspot.com/news/52011-one-minute-on-the-internet-640tb-data-transferred-100k-tweets-204-million-e-mails-sent.html

Banerjee, S., & Chua, A. Y. (2016). In search of patterns among travellers' hotel ratings in TripAdvisor. *Tourism Management , 53*, 125-131.

Browning, D., & Sparks, D. A. (2011). The impact of online reviews on hotel booking intentions and perception of trust. *Tourism Management , 32* (6), 1310-1323.

Chu, W. W. (2014). *Data Mining and Knowledge Discovery for Big Data.* Los Angeles: Springer.

Chua, A. Y., & Banerjee, S. (2013). Reliability of reviews on the Internet: The case of Tripadvisor. *In Proceedings of the World Congress on Engineering and Computer Science. 1*, pp. 23-25.

Chatterjee, P. (2001). Online Reviews – Do Consumers Use Them? In M. C. Gilly, J. Myers-Levy, Provo, & UT (Ed.), *In ACR 2001 Proceedings* (pp. 129-134). Association for Consumer Research.

Duan, W., Gu, B., & Whinston, A. B. (2008). Do online reviews matter? — An empirical investigation of panel data. *Decision support systems , 45* (4), 1007-1016.

Dellarocas, C., Gao, G. (., & Narayan, R. (2010). Are Consumers More Likely to Contribute Online Reviews for Hit or Niche Products? *Journal of Management Information Systems , 27* (2), 127–157.

Denecke, K. (2008). Using SentiWordNet for Multilingual Sentiment Analysis. *In Data Engineering Workshop, 2008. ICDEW 2008. IEEE 24th International Conference* (pp. 1-7). Hannover: IEEE.

Ding, X., Liu, B., & Yu, P. S. (2008). A Holistic Lexicon-Based Approach to Opinion Mining. *In Proceedings of the 2008 international conference on web search and data mining* (pp. 1-9). ACM.

Feldman, R., & Sanger, J. (2007). *The Text Mining Handbook: Advanced Approaches in Analyzing Unstructured Data.* NewYork: Cambridge University Press.

Filieri, R., & McLeay, F. (2013). E-WOM and Accommodation: An Analysis of the Factors That Influence Travelers' Adoption of Information from Online Reviews. *Journal of Travel Research , 53* (1), 44– 57.

Gretzel, D., & Yoo, K. H. (2008). Use and impact of online travel reviews. *Information and communication technologies in tourism 2008* , 35-46.

Hu, N., Bose, I., Koh, N. S., & Liu, L. (2012). Manipulation of online reviews: An analysis of ratings, readability, and sentiments. *Decision Support Systems , 52* (3), 674-684.

Hu, N., Liu, L., & Zhang, J. (2008). Do Online Reviews Affect Product Sales? The Role of Reviewer Characteristics and Temporal Effects. *Information Technology and Management , 9* (3), 201-214.

Hu, M., & Liu, B. (2004). Mining and Summarizing Customer Reviews. *In Proceedings of the tenth ACM SIGKDD international conference on Knowledge discovery and data mining* (pp. 168-177). ACM.

Hoffman, M., & Klinkenberg, R. (2014). *RapidMiner: Data Mining Use Cases and Business Analytics Applications.* Boca Raton, Florida: Taylor & Fancis Group, LLC.

Judith, C., & Mayzlin, D. (2006). The effect of word of mouth on sales: Online book reviews. *Journal of marketing research , 43* (3), 345-354.

Jungermann, F. (2009). Information Extraction with RapidMiner. *In Proceedings of the GSCL Symposium'Sprachtechnologie und eHumanities* (pp. 50-61). TU Dortmund.

Java, A., Finin, T., Tseng, B., & Song, X. (2007). Why We Twitter: Understanding Microblogging Usage and Communities. *In Joint 9th WEBKDD and 1st SNA-KDD Workshop '07* (pp. 1-10). ACM.

Jo, Y., & Oh, A. (2011). Aspect and sentiment unification model for online review analysis. *In Proceedings of the fourth ACM international conference on Web search and data mining* (pp. 815-824). ACM.

Levy, S. E., Duan, W., & Boo, S. (2013). An Analysis of One-Star Online Reviews and Responses in the Washington, D.C., Lodging Market. *Cornell Hospitality Quarterly , 54* (1), 49-63.

Liu, B. (2012). *Sentiment Analysis and Opinion Mining.* (G. Hirst, Ed.) Chicago: Morgan & Claypool Publishers.

Malbon, J. (2013). Taking fake online consumer reviews seriously. *Journal of Consumer Policy , 36* (2), 139-157.

Ögut, H., & Cezar, A. (2012). The factors affecting writing reviews in hotel websites. *Procedia-Social and Behavioral Sciences. 58*, pp. 980-986. Elsevier Ltd.

Pak, A., & Paroubek, P. (2010). Twitter as a Corpus for Sentiment Analysis and Opinion Mining. *In LREc, 10,* 1320-1326.

Pang, B., & Lee, L. (2004). A sentimental education: Sentiment analysis using subjectivity summarization based on minimum cuts. *In Proceedings of the 42nd annual meeting on Association for Computational Linguistics* (p. 271). Association for Computational Linguistics.

Pang, B., & Lee, L. (2008). Opinion Mining and Sentiment Analysis. *Foundations and Trends in Information Retrieval , 2* (1-2), 1–135.

Papathanassis, A., & Knolle, F. (2009). Exploring the adoption and processing of online holiday reviews: A grounded theory approach. *Tourism Management , 32* (2), 215-224.

Pozzi, F. A., Fersini, E., Messina, E., & Liu, B. (2017). *Sentiment Analysis in Social Networks.* Cambridge, Massachusetts, USA: Elsevier Inc.

Pontiki, M., Galanis, D., Pavlopoulos, J., Papageorgiou, H., Androutsopoulos, I., & Manandhar, S. (2014). Aspect Based Sentiment Analysis. *In Proceedings of the 8th International Workshop on Semantic Evaluation* (pp. 27–35). SemEval 2014.

Prichystal, J. (2016). Mobile application for customers' reviews opinion mining. *In 19th International conference enterprise and competitive environment 2016* (pp. 373-381). Elsevier Ltd.

Schuckert, M., Liu, X., & Law, R. (2016). Insights into suspicious online ratings: direct evidence from TripAdvisor. *Asia Pacific Journal of Tourism Research , 21* (3), 259-272.

Sohail, S. S., Siddiqui, J., & Ali, R. (2016). Feature extraction and analysis of online reviews for the recommendation of books using opinion mining technique. *Perspectives in Science , 8*, 754-756.

Rapidminer Inc. (2016, December 1). *RapidMiner, the Industry's #1 Open Source Data Science Platform*. Retrieved December 5, 2016, from Rapidminer: https://rapidminer.com

Tumasjan, A., Sprenger, T. O., Sandner, P. G., & Welpe, I. M. (2010). Predicting Elections with Twitter: What 140 Characters Reveal about Political Sentiment. *In Fourth International AAAI Conference on Weblogs and Social Media* (pp. 178-185). Lehrstuhl für Betriebswirtschaftslehre Strategie und Organisation.

Tuominen, P. (2011). *The Influence of TripAdvisor Consumer-Generated Travel Reviews on Hotel Performance.* University of Hertfordshire Business School.

Tetlock, P. C. (2007). Giving Content to Investor Sentiment: The Role of Media in the Stock Market. *The Journal of Finance, 62* (3), 1139- 1168.

The Wall Street Journal. (2007, June 1). *Deconstructing TripAdvisor*. Retrieved Dec 8, 2016, from http://www.wsj.com/europe

Tsaoa, W.-C., Hsieh, M.-T., Shih, L.-W., & Lin, T. M. (2015). The influence of hotel reviews on booking intention from the perspective of consumer conformity. *International Journal of Hospitality Management , 46*, 99-111.

Tripadvisor. (2016, Dec 1). *About TripAdvisor*. Retrieved Dec 8, 2016, from www.tripadvisor.com