# Sentiment Analysis of Customer Data

**Olivera Grljević**
University of Novi Sad, Faculty of Economics in Subotica, Subotica, Serbia
**Zita Bošnjak**
University of Novi Sad, Faculty of Economics in Subotica, Subotica, Serbia

## Abstract

The value of a customer for the company is not measured only by the monetary effect, but also by the degree of its satisfaction. Satisfied customers spread positive word of mouth, while dissatisfied customers spread negative one (in online or offline environment). Their voice shapes company reputation and the reputation is one of the key parameters in choosing a product, service, or a company.

In order to look into customer satisfaction, companies need to get certain feedback from consumers. The traditional way of collecting feedback from consumers is questionnaires. However, the way in which consumers express their opinions on social media opens up new opportunities and makes it easier for companies to cover a larger number of consumers, to collect data of interest and to continuously monitor the brand.

Social media sites have significantly changed the nature of human activities, interactions and ways of disseminating information. Consumer behavior has changed accordingly. In search for information about products and services consumers are planning to buy or use, offline sources of information are increasingly replaced with online sources and e-WoM. Before making a purchase decision, consumers visit many sites and read the content and comments that other users have generated. Within these contents – texts that users voluntarily post on the Internet and make it publicly available – users freely express their views, opinions, describe their consumer experience, the problems they have encountered and the way in which the problems were addressed, they point to the aspects of products or services they are satisfied or dissatisfied with. These contents shape the opinions of future consumers in great extent and affect their consumer actions. Numerous studies have confirmed the impact these sources have on consumer behavior:

• Consumers appreciate opinions of other individuals and trust them more than the company's promotional campaigns,
• Consumers have equal trust in online comments and reviews as in personal recommendations from friends,
• Online comments on products and services are on the third place according to the influence on buying decisions (after coupons and discounts).

Each business can benefit from analysis of social media content since it comprises useful feedback from consumers in form of expressed opinions and attitudes. Opinions and attitudes are extremely subjective. Due to subjectivity it is necessary to analyze a collection of opinions of different people instead of a single opinion which expresses a subjective view of an individual. In addition to this large number of available sources, the amount of data makes it impossible to manually process them and identify the general pattern, problem or source of (dis)satisfaction. Hence, automated analysis of unstructured content from social media sites is required, i.e. the application of sentiment analysis or opinion mining.

Sentiment analysis is a young research area which has rapidly developed during the past 10 years and it has achieved significant commercialization. The surrounding industries also experienced significant expansion. Sentiment analysis allows companies to analyze public opinion, attitudes and emotions directed at a particular entity (e.g. a particular person, a political candidate, a party, a law, a company, a specific product, or a product feature) which expresses particular sentiment or points to (usually positive, negative or neutral), as well as all variations and gradations of the sentiment. In addition to the significance of sentiment analysis for business, this paper deals with sentiment analysis process and underlying data mining techniques.

## Keywords

Sentiment analysis, online reviews, user-generated contents, data mining techniques.

## Introduction

Consumer behavior when searching for information about products and services has changed significantly. Offline resources have been replaced by means of word-of-mouth e-marketing, such as social media sites or sites for online reviews of products/services – online comments (Gruen et al., 2006). When making a purchase decision, consumers consult publicly available content on the Internet, or online comments, to get insight into other consumers' attitudes and experiences regarding the product/service. The subject of online comments varies from technical devices (computers, cameras, mobile phones, etc.), books, hotels, restaurants, cars to specific services such as education. After coupons and discounts, online comments have the most influence on purchasing decisions (Yang et al., 2015/2) and consumers value and trust the opinions of other individuals more than the company's promotional campaigns (Berthon et al., 2012; Pitt et al., 2002). Furthermore, there is no difference in consumers' level of trust in online comments and their trust in personal recommendations of a friend (Park et al., 2007; Gligorijevic & Luck, 2012). In addition to the impact on consumers' attitudes and actions, online comments increase the level of consumers' trust in the company, directly affect which companies customers choose among competitors and influence the success rate of attracting new consumers. It is of utmost importance for each company to respond promptly and adequately to complaints and praises of its consumers and to manage carefully the company's online reputation, as negative online comments could have unfavorable business implications – even a sole unhappy consumer can ruin the reputation of a company (Tripp & Grégoire, 2011). Consumers are inclined to choose companies with a positive reputation and are willing to pay more for their products. The degree to which consumers are engaged in activities on a company's social media reflects their perception of the reputation of a company. Per definition engagement contains a cognitive aspect (e.g. the consumer is interested in the company's activities), a behavioral aspect (e.g. participation in activities) and / or the emotional aspect (e.g. positive attitude towards the company's activities) (Dijkmans et al., 2015). Achieving a high degree of consumer engagement promotes the company's reputation and brand loyalty and influences purchasing decisions.

The authors (Hudson et al., 2015) emphasize the relevance of emotions in the process of building relationships between consumers and a brand, as well as their impact on the readiness of consumers to recommend the product. Research on how emotions are processed confirms that emotions are formed before any processing of information and that consumer behavior is strongly influenced by them (Hudson et al., 2015). Factors with the most powerful overall impact on the individual's purchasing intentions are emotional factors that result in love that is based on sensory pleasures (Pawle & Cooper, 2006). In other words, it is suggested that the key emotional trigger for strengthening the relationship between consumers and a brand is creating a strong intimate relationship with the brand. Due to the fact that consumers today react more positively to content distributed on social media than on classic advertising campaigns, their emotional response to the first strengthen the brand's connection to customers. The same result is achieved when consumers are given the opportunity to interact with the company on social media through sharing brand knowledge, consumer experience or suggestions. Consumers' expectations lead to an emotional response, which results in a positive or negative feeling, thus shaping their satisfaction, trust and loyalty.

With an adequate analysis of the online content that users generate and especially by analyzing the emotions hidden in it, companies can acquire critical information for improving their business. Sentiment analysis is an approach to data analysis that allows companies to analyze a variety of publicly available opinions, sentiments, attitudes and emotions directed to a specific entity. An entity may be, for example, a particular person, a political candidate, a party, a law, a company, a specific product or some of its characteristics. User-generated content expresses or points to a certain sentiment (usually positive, negative or neutral, as well as all variations and gradations of the sentiment). Through sentiment analysis companies familiarize themselves with their customer base at a higher level, which includes familiarizing with the customers emotions, which further enables optimization of marketing messages, trends forecasting, product/service development and monitoring of the companies' online reputation. By introducing a customer perspective into services and business transactions, companies can deliver a product/service that is most suited to consumers' needs.

In addition to the importance of sentiment analysis for business operations, we talk about the

flow of sentiment analysis and data mining techniques in the background of this process. The paper is structured as follows. The next section describes the process of sentiment analysis. The second section gives an overview of the basic methods and techniques used in sentiment analysis of a content, while the third section describes diversified applications of sentiment analysis of online comments. Additionally, this section illustrates the sentiment analysis conducted on collections of comments from three leading sites in the domain of e-commerce, entertainment and ranking of local restaurants (Amazon, IMDB, and Yelp). The last section sums up the author's conclusions.

## 1. The process of sentiment analysis

The sentiment analysis consists of three tasks that are described below together with the challenges of these tasks.

1.  ***Creating resources for sentiment analysis.*** Supervised learning techniques for data mining are at the heart of sentiment analysis. They require pre-processed training sets that include clearly defined and distinguished examples of each class. Text data collected from social media sites have to be adequately prepared in order to be presented to the selected data mining methods and techniques. The task of creating the resources for sentiment analysis usually requires including additional metadata to a set, i.e. by text annotation. Metadata is defined as information attached to statements expressing sentiment, whereas annotation means each metadata tag used to mark the element of a data set. The content in which the sentiment is expressed can be annotated by its polarity (positive, negative, neutral), the goal of the sentiment – the so-called aspect, or by other semantic, syntactic or lexical information. In order to achieve an effective training of the learning algorithm, the annotations have to be accurate and relevant to the task. For this reason, the annotation of texts is a critical step in the development of applications for natural language processing (Pustejovsky & Stubbs, 2012). Sets of texts are called corpora, and a set of texts annotated based on the same specification is called an annotated corpus. The annotation of a corpus is valuable primarily because it creates a basic knowledge base for training machine learning algorithms (MLA) for the implementation of sentiment analysis, as well as for validating the theory of phenomena in texts (structural, logical, semantic, and syntax) (Hovy & Lavid, 2010). The main challenges in the process of creating resources for sentiment analysis are the ambiguity of words, the granularity (the opinion can be expressed by words, sentences or whole phrases), differences in the way of expressing opinions in different types of texts (blogs, newspaper articles, content from the forum, etc.) and so on. These challenges impose the need to annotate each type of text that is subject to analysis separately, as well as the need for the existence of necessary annotation resources for different languages (Montoyo et al., 2012).

2.  ***Classification of text according to the polarity of the expressed opinion.*** The early classification attempts were made on online comments in the film industry (Pang et al., 2002), online comments on technical products, books, alternative music and films (Dave et al., 2003) and on online comments within the automotive and banking industry and on tourist destinations (Turney, 2002). The classification was done at a rough level of granulation – at document level (comment) – based on associated sentiment indicators, such as the number of stars or the presence of positive and negative adjectives and adverbs. In these first attempts the central concept was based on the automatic learning of the sentiment classification model using a set of previously categorized documents as training data. Additional annotation of collected documents was not required. The classification considered the whole document and the starting premise was that one document discussed one topic. The developed models were used to classify or predict the sentiment category (positive, negative, and, rarely, neutral) of novel documents. The demand of many applications to provide enough detail about the prevailing customer opinion on the different aspects of the observed entity is often not satisfied by this type of classification (Medhat et al., 2014). In addition, the rapid increase in the number of online user comments makes it difficult for companies and consumers to

obtain a comprehensive picture of the prevailing customer opinion. Hence, the traditional sentiment analysis at the document level has been abandoned and the number of studies that examine the possibility of a more detailed analysis of thoughts and feelings, at the level of the word/phrase (Breck et al., 2007; Zirn et al., 2011), sentence/paragraph level (Gamon et al., 2005; Kim & Hovy, 2006; Wiegand & Klakow, 2009) or part of the sentence level (Wilson et al., 2005; Choi & Cardie, 2008) is increasing.

Classification at the sentence level treats each sentence separately and starts from the assumption that each sentence contains exactly one opinion (positive, negative or neutral) (Jagtap et al., 2013). This approach requires separating objective from subjective sentences prior to the sentiment analysis, which then determines if attitudes expressed in subjective sentences are positive or negative. Classification at the word/phrase level is accomplished by determining the polarity of the word/phrase used for classification. In most online reviews, more aspects or more product characteristics are mentioned, different entities are compared, and various sentiments about them are expressed. The early sentiment analyses mentioned above are not adequate in such cases, because they focus on the discovery of the prevailing sentiment, not the aspect and therefore rather reveal the customers' general attitudes than the exact customers' preferences (Liu, 2012; Broß, 2013).

For more detailed analysis, an aspect-oriented approach has to be applied, which analyzes the customers' sentiment according to an individual aspect /characteristic of the product. The starting point of the classification at the aspect level is an opinion based on a positive or negative sentiment referring to a person, object or aspect (the goal of the sentiment). The basic task is to determine the sentiment about the entity or aspect (Liu, 2012). This is the reason why the classification at the aspect level enables more detailed analysis of the sentiment according to the aspects and identifies the customers' preferences more precisely. In order to conduct an aspect-oriented analysis of the sentiment, it is necessary to create

representative corpora by the annotation procedure.

The main challenges of classifying the text are manifested in the need for contextual analysis of the expressed opinions. Such classification requires analysis at different levels (lexical, syntactical, semantic and pragmatic), the utilization of robust, supervised, semi-supervised or unsupervised methods, the analysis of the influence of negation and its scope, detection of irony and detection of implicitly expressed sentiment in objective sentences (Montoyo et al., 2012). Concerning the aspect-oriented analysis, the key challenge is to identify the parts of the text in which the opinion is expressed and the span of the sentiment expression, as well as the target of the expressed opinion (on who/what is the opinion on).

3. *Application of sentiment analysis.* This task can be a stand-alone task, combined with other tasks of natural language processing (summing up opinions, finding opinions, etc.), or a part of a complex application (for example, ranking the product according to the expressed opinion of users in intelligent recommendations or trend detection systems). The main challenge that this task needs to tackle is the adequate combination of various sentiment analysis techniques and methods from the area in which it is applied (Montoyo et al., 2012).

## 2. Sentiment analysis methods and techniques

The task of sentiment analysis in its basis is the task of classifying expressed feelings, usually to positive or negative ones. Neutral opinions are expressed less frequently than positive/negative ones and have significantly less or no influence on the other customers' opinions. Therefore, companies are often not interested in them. The sentiment classification techniques can be divided into machine learning techniques that use linguistic attributes, a lexicon-based approach (vocabulary or corpus), which accepts a list of known and previously prepared sentiment expressions as an input, and a hybrid approach that combines the previous two. Figure 1 shows the most commonly used techniques of the sentiment classification.
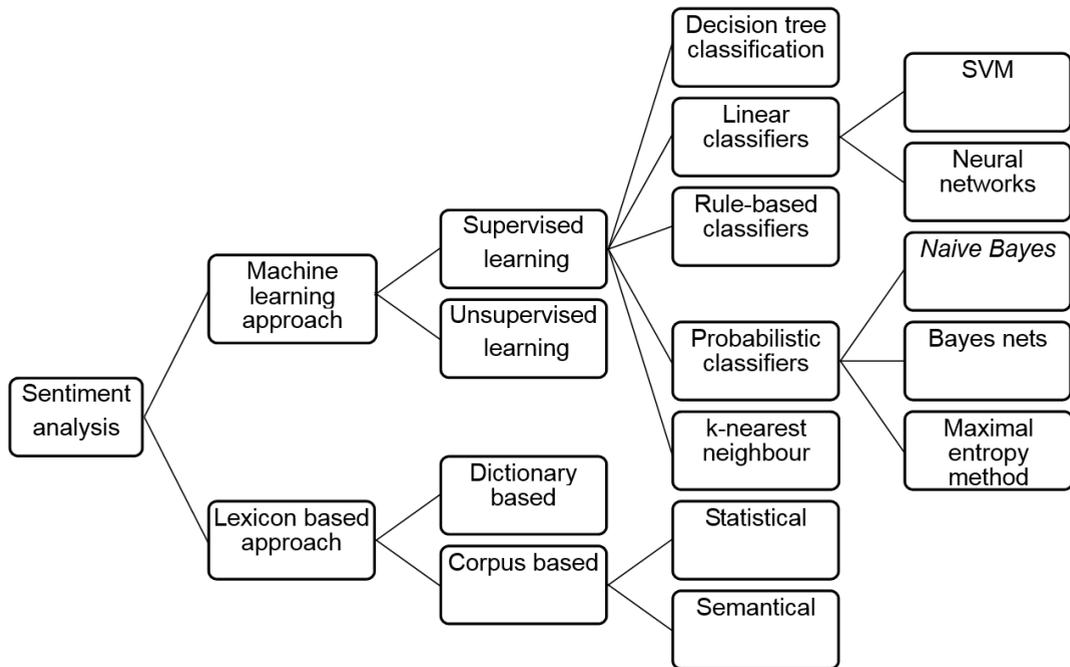
**Figure 1** Sentiment classification techniques
**Source:** Medhat et al., 2014

Machine learning techniques are divided into supervised and unsupervised learning. Supervised learning methods require a large amount of pre-processed (tagged) training data, while unsupervised learning methods do not impose such a requirement.

Classifiers generated by supervised learning methods are created automatically based on the characteristics of the categories from the collection of previously classified and labeled training documents. In the development of classification models with supervised learning, the following issues should be considered (Feldman & Sanger, 2013):

- Which categories should be used in the classification of instances?
- Provide appropriate training sets for each category.
- Decide on the attributes that describe each instance. Since most of the supervised learning algorithms can focus on relevant attributes, it is advisable to generate as many attributes as possible.
- Which algorithms should be applied? The most commonly used algorithms of supervised learning, as shown in Figure 1, are decision trees, linear classifiers, probabilistic classifiers, rule-based classifiers and k-nearest neighbors.

Usually it is easy to collect a large amount of data from the Internet, but it is difficult to prepare a sufficient number of clearly separable and labeled documents (for example, documents containing exclusively positive or negative comments). For this reason, unsupervised learning methods are often used. Unsupervised sentiment analysis can be based on sentiment words, patterns of syntax identified in the sentiment (Turney, 2002) or sentiment lexicons.

Lexicon-based approaches use sentiment dictionaries to recognize sentiment in texts. There are three basic ways of collecting and making a sentiment lexicon: 1) Manual access, which requires intensive work and is time-consuming. It is rarely applied alone, but rather as the final test of an automated approach to avoid errors that automated approaches sometimes make (Medhat et al., 2014); 2) A dictionary-based approach starts from manually collected basic sentiment words that are associated with the sentiment orientation. Afterwards, adequate synonyms and antonyms are searched in the dictionary (WordNet or other online dictionary) 3) Corpus-based approaches start from the list of sentiment words specific to the domain that is the subject of the analysis and search other sentiment words together with their contextual orientation in a large corpus (Medhat, et al., 2014).

## 3. Diversified objectives and approaches to sentiment analysis

Social media sites offer their users a free exchange of thoughts and experiences about the purchasing process, product and services' characteristics and quality. Because of their importance and widespread distribution, online comments become a key type of user-generated content and take a central place of research in the sentiment analysis domain. The emphasis of this research is on connecting different business objectives and desired benefits of the analysis with diversified approaches, to obtain best results. For companies, online users' comments represent an excellent source of information for competitive intelligence used for developing and improving marketing activities. Through customer comments, companies receive sincere, impartial feedback on their products, as well as on their competition's products. In this way, the market segments for which the product is most suitable can be identified along with the mismatches between offered products and customers' preferences and differences in the capabilities of the company and its competition. By interpreting a large amount of data and associated feelings, companies can improve business decision-making by developing new ideas and solutions for their technological or economic problems (Balahur, 2011).

From the customers' point of view, online comments provide useful information that shorten the search time and make purchase decisions more efficient (Yang et al., 2015/2). More than 80% of customers expect the option to consult other customers or professionals before the act of purchase to be implemented on the seller's website (D'Avanzo & Kuflik, 2013). Hence the development of the sentiment analysis applications for the needs of e-commerce is a wide research field. E-commerce sites exploit the possibilities customers' online reviews offer; they analyze, classify and summarize opinions from product reviews that customers visiting the site can use to compare or recommend a product, or consult each other.

The authors D'Avanzo and Pilato (2015) introduce a cognitive-based procedure that analyzes customer opinions from specific types of marketing and sums them up visually to help resolve the problem of customer being overwhelmed by large amounts of online comments and speed up their purchasing activity. The approach mimics the

"zone of proximal development" of Vigotsky[1], which is widely exploited in collaborative learning communities.

Recommendation systems based on online comments are typical applications incorporated into e-commerce sites. Thanks to these systems, users of the site do not search the comments, but get recommendations of a product that best suits their preferences. However, these systems heavily rely on structured metadata when providing recommendations, such as star rating estimates, and ignore text content despite the text being the most important source of information in online comments.

The authors Ganu et al. (2009) believe that the customer experience would be significantly improved if the way, in which the content of online comments is structured would be considered. This means also considering the parts of the comments on certain characteristics of products and the sentiment of the comments on each characteristic.

The authors Hu and Liu (2004) as well as Kim et al. (2014) and Broß, (2013) describe different approaches to product reviews analysis. The approach described by Hu and Liu (2004) constitutes of three steps of efficiently summarizing online product reviews. Their approach entails summarizing comments, which allows potential customers to easily see the opinions of existing customers about the product of interest. At the same time this step allows manufacturers to combine reports from various retail sites to see the general public's attitude towards their own or competitors' products. The authors have developed a system that allows users to summarize comments on the product as a whole, as well as on individual product characteristics, thus gaining a better insight into positively and negatively assessed aspects. In the first step, only those characteristics that consumers expressed their opinions on are analyzed, through the application of data mining and natural language processing (NLP) methods. The next step identifies all sentences that contain expressed opinions and determine their sentiment. NLP techniques are used to identify a set of adjectives that usually express an opinion, and for each identified sentiment word its orientation (positive or negative) is determined

---

[1]   The distance between the most difficult task that one can solve independently and the most difficult task that one can solve with someone's assistance.

using the bootstrapping[2] technique and WordNet[3]. In the last step, the results are summarized.

The authors Yang et al. (2015/1) also dealt with the problem of summarizing reviews. Their approach is based on the ontology tree which is the starting point for extracting the most representative expressions and opinions of customers about the product by forming attribute-sentiment pairs. Their system then calculates the polarity of these pairs based on entropy and gives a comprehensive report to producers or customers.

The authors Kim et al. (2014) dealt with the classification of the sentiment according to the aspects of the product. Their methodology enables the identification of positive and negative opinions about products. Unlike their predecessors, the authors have achieved a detailed measurement of expressed opinions, which allows determining the information relevant to the particular mobile phone device, as well as the part of the specification of the mobile phone that is its competitive advantage in relation to other products. The authors used a domain dictionary composed of 500 words that were associated with values ranging from -5 to +5 for each sentiment word. This method proved difficult for implementing in other domains.

In his paper, Broß (2013) dealt with an aspect-oriented sentiment analysis, which involves identifying the aspect of the product written online comments refer to and classifying the content as positive or negative. In the sentiment classification, a lexicon-based approach was applied together with supervised learning techniques. The author also examined the possibility of automatically generating sentiment vocabulary according to the predefined aspects of the product, as well as dictionaries of all sentiment concepts for the given domain.

## 3.1. The art of modeling for better performance – an illustration

This chapter illustrates the various modes of sentiment analysis application on online product reviews and how these approaches affect the performance of the classification model. Under various modalities we understand diversified sources of data, different classification models and different performance measures available to analysts, which they select exclusively based on their ex-

perience and a few heuristic rules in the domain of sentiment analysis. Data for the analysis is taken from the research conducted by Kotzias et al (2015) and refer to reviews on mobile phones, movies and restaurants. Each set of comments is made up of randomly selected 500 positive and 500 negative comments downloaded from Amazon, IMDB and Yelp sites, that the authors have prepared for research purposes - labeled and preprocessed for the classification. Each set was divided into a training set (80% of reviews) and a test set (20% of reviews). The preprocessing of data from the source files in the RapidMiner[4] tool is shown in Figure 2.

In the first step (the Transform Cases process), the conversion of uppercase letters to lowercase ones is conducted. The second step (the Tokenize process) transforms the content of a document into a list of words. The third step (the Filter Stopwords process) filters frequent words that do not bring new knowledge. The fourth step (the Generate n-grams-terms process) allocates bigrams (two words in a row) and the last step (the Filter tokens by length process) extracts words that are composed of one or two characters and those that have more than 25 characters. The excerpt from the data set after the pre-processing is shown in Figure 3.

The term vector is created using the TF-IDF measure (Term Frequency – Inverse Document Frequency) that shows the relative significance of the terms for a particular document compared to the term's importance for all documents in the data set (Weiss et al., 2010). The result of the preprocessing and application of the TF-IDF measure is the set of data in which one line represents a document (one comment or one sentence), while the columns represent words that appear in the whole corpus. Allocated values indicate the importance that a word has for a specific document. If the observed term appears relatively often in one document, TF-IDF measures will be higher and will indicate that the given word is characteristic for a specific document. If a large number of documents in a corpus contains the observed term, the TF-IDF measure will be lower and will indicate that the given word is not important. Hence, the values closer to one indicate greater significance, while zero indicates that the word does not appear in the analyzed document. This prepared set is an input for the construction of a classifier.

---

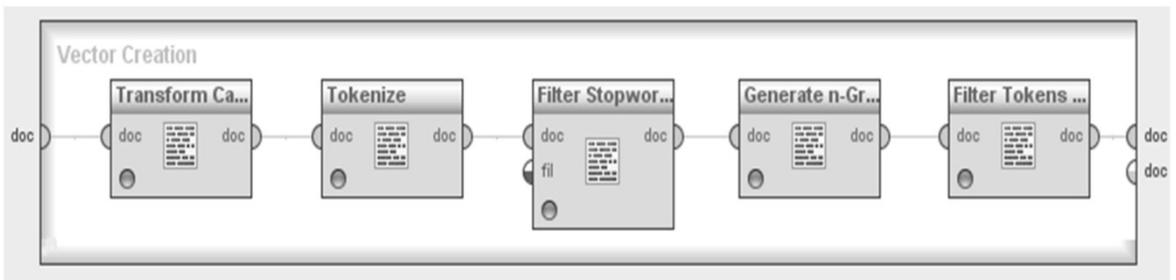[2] A more detailed explanation of the technique is available at http://www.stat.rutgers.edu/home/mxie/rcpapers/bootstrap.pdf
[3] Available at https://wordnet.princeton.edu/

[4] https://rapidminer.com/

**Figure 2** Comments pre-processing
**Source:** The authors (from RapidMiner tool)

| label | text | able ↓ | able_use | able_voice | ac_charger | accessable |
|---|---|---|---|---|---|---|
| negativno | I would have given no star if I was able. | 0.408 | 0 | 0 | 0 | 0 |
| pozitivno | "The range is very decent, I've been able to roam around my house with the ph... | 0.350 | 0 | 0 | 0 | 0 |
| pozitivno | I was able to do voice dialing in the car with no problem. | 0.294 | 0 | 0.360 | 0 | 0 |
| pozitivno | I love being able to use one headset for both by land-line and cell. | 0.257 | 0.315 | 0 | 0 | 0 |

**Figure 3** The dataset after pre-processing
**Source:** The authors (from RapidMiner tool)

Three representative models have been developed, using Naive Bayes, SVM and kNN binary classifiers, which identify the sentiment expressed in novel comments. Classifier performance is measured using Precision, Recall and F measures (Baeza-Yates & Ribeiro-Neto, 1999; Wiebe et al, 2005). Precision shows the relation of correctly classified entities to the majority (positive) class against the total number of entities classified in this class by the model. Recall shows which percentage of positive entities is correctly classified (as positive). Since these two measures cannot be directly compared, the F-measure is calculated as their harmonic mean.

Figure 4 gives a parallel view of the classifiers' performance when online comments from different sources are treated as individual sets. Based on the results, it can be noticed that the classification algorithms behave differently depending on data. In the case of comments on products from the Amazon site, regardless of whether it is a positive class (marked with a "+") or a negative class (marked with a "-"), the SVM classifier has the best performance. In the case of comments on IMDB movies, the kNN algorithm gives the best results in predicting sentiment orientation, while in the case of restaurant reviews from Yelp, the performance of the classifiers is very similar.

When comments are combined into a single dataset to provide a broader word base for classifier training, the SVM algorithm is superior to others, regardless of whether the performance is measured on a positive or a negative class. Namely, the model successfully identifies negative comments in 79.70% of cases, while positive comments are recognized in 75.67% of cases, as shown in Figure 5. All three classifiers are more successful when identifying comments with a negative connotation compared to positive comments, which may be the consequence of a wider commonly used vocabulary for expressing negative sentiment than a positive one.

The obtained results indicate the possibility of successful implementation of classification algorithms on online comments, but at the same time, they confirm that data analysis, although scientifically founded, requires a great deal of analytical skills (which could partially be acquired through experience) and a large number of experiments with different parameters in order to achieve the best result possible. By more thorough insight into the content structure and by comprehensive analysis of classification errors, a more detailed insight into the possibilities of improving the performance of the classifier can be obtained.
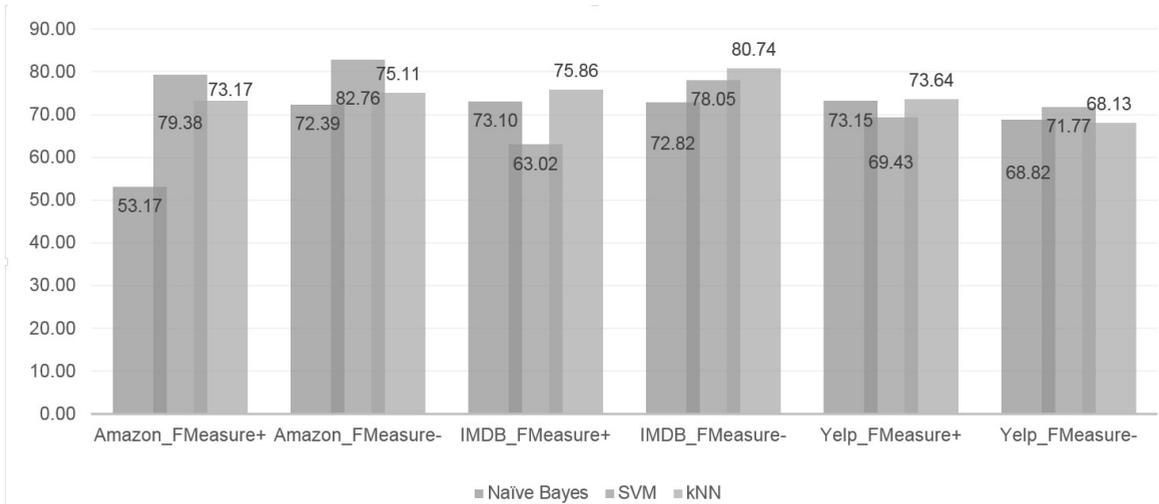
**Figure 4** Parallel view of the classification results on individual sets
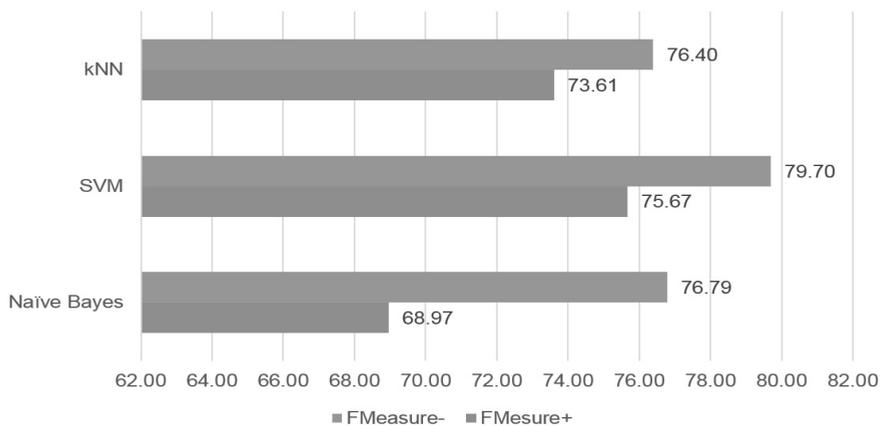**Source:** The authors



**Figure 5** The result of the classification of a comprehensive corpus
**Source:** The authors

The successfully classified content obtained after the utilization of the developed classifier can be used for:

1. Identifying global customer (dis)satisfaction by summarizing positive/negative online comments;
2. Comparison of the online reputation of a company, brand or product with the competition;
3. Analysis of positive/negative comments on product aspects in order to identify those characteristics that consumers prefer. For this purpose, additional modeling of the results is necessary to identify the prevailing topics within the positive and negative comments;

4. Monitoring the sentiment over time, with the aim of identifying variations in customers' (dis)satisfaction.

## Conclusion

Social media sites and sites with online products/services reviews have a great impact on customers' attitudes and actions: they increase customers' confidence in the company, directly affect which company customers choose among competitors and influence the process of attracting new customers. As the content posted by users contain feedback on products/services or the company itself in the form of expressed opinions and attitudes, their adequate analysis, and primarily the analysis of the emotions hidden in them,

can give companies critical information to improve their business. Sentiment analysis allows companies to analyze unstructured content from social media sites, such as expressions of the prevailing opinion, attitudes and emotions directed at a particular entity (e.g. a specific product or a product characteristic). This unstructured data expresses or points to a certain sentiment (usually positive, negative and neutral, as well as all variations and gradations of the sentiment). However, the opinions expressed in this data are very subjective. For this reason, it is necessary to analyze the set of opinions of a large number of people instead of a single opinion expressed as a subjective view of an individual. Considering the number of data sources and the amount of available data, it is impossible to manually process and identify a general pattern or causes of customers' (dis)satisfaction. Therefore, automated approaches and procedures from machine learning domain are used, especially sentiment analysis or opinion mining.

The paper firstly describes the process of sentiment analysis on three groups of tasks and the specific challenges that each task tackles. It was pointed out that the data resource has to be adequately prepared, which is usually accomplished by adding metadata to a set, that is, annotating the texts according to polarity, the goal of the sentiment – aspect, but also by attaching other semantic, syntactic or lexical information to the source data. Annotations have to be accurate and relevant to the task in order to achieve effective training by the data mining algorithm, so adequate annotation is critical for the sentiment analysis. It is pointed out that text classification (to positive, negative or neutral polarity of the expressed opinion in the comments) could be conducted at different levels of analysis (at the level of the document, sentence, word or phrase). The paper particularly highlights the importance of properly selected and robust supervised and unsupervised methods. It also underlines the specific difficulties in analysis that arise from the potential existence of negation with varying range of influence, irony that is hard to detect and implicitly expressed sentiment in objective sentences. It is also pointed out that the implementation of sentiment analysis could be even more complex because it requires combining with other tasks (summarizing opinions, finding opinions, ranking products according to some expressed opinion, etc.) and methods from the specific domain of interest.

An illustrative example of sentiment analysis over three different sources of data (Amazon, IMDB and Yelp sites) with three different classifiers (Naive Bayes, SVM and kNN) showed how significant the influence of analysts' choice among diversified input parameters could be on the performance of the built classification model. More precisely, it has been shown that classification algorithms perform differently depending on the specific data, and also that the selection of the classifier directly affects the results achieved. In the case of comments on Amazon's products, the SVM classifier had the best performance. In the case of comments on IMDB movies, the kNN classifier gave the best results. In anticipation of the sentiment orientation of restaurant reviews from the Yelp site, the performance levels of all three classifiers was very similar. In case of comments being combined into a single set of data, the SVM classifier was superior to others. All three classifiers identified comments with negative connotations more successfully than positive comments.

Social media sites allow users a free exchange of opinions on products'/services' quality, purchasing process and buying experience. For companies this is a freely available source of data for the development and improvement of marketing activities. Honest, impartial feedback can identify the market segments for which a particular product/service is most fitting, data can point to mismatches between products offered by a company and customers' preferences and differences in capabilities of a company and those of its competitors. By introducing a customers' perspective into services and business transactions, organizations can offer the customers a product/service that is tailored according to their needs. Furthermore, companies can improve business decision-making by identifying new ideas and solutions for technological or economic issues. Considering everything mentioned, we can say that companies that successfully master the sentiment classification techniques and the sentiment analysis of their client base have a significant competitive advantage over companies that have not yet mastered intelligent technologies and still rely on classical marketing approaches. **SM**

# References

Baeza-Yates, R. & Ribeiro-Neto, B. (1999). Modern Information Retrieval. Addison-Wesley Professional.

Balahur, A. D. (2011). Methods and Resources for Sentiment Analysis in Multilingual Documents of Different Text Types. Universidad de Alicante.

Berthon, P. R., Pitt, L. F., Plangger, K. & Shapiro, D. (2012). Marketing meets Web 2.0, social media, and creative consumers: Implications for international marketing strategy. Business Horizons, 55 (3), 261-271.

Breck, E., Choi, Y. & Cardie, C. (2007). Identifying expressions of opinion in context. Proceedings of the 20th International Joint Conference on Artificial Intelligence, IJCAI-2007. Hyderabad, India.

Broß, J. (2013). Aspect-Oriented Sentiment Analysis of Customer Reviews Using Distant Supervision Techniques. Berlin, Germany: Freie Universität Berlin, Fachbereich für Mathematik und Informatik.

Choi, Y. & Cardie, C. (2008). Learning with compositional semantics as structural inference for subsentential sentiment analysis. Proceedings of the Conference on Empirical Methods in Natural Language Processing, EMNLP '08. Stroudsburg, PA, USA: Association for Computational Linguistics.

D'Avanzo, E. & Kuflik, T. (2013). E-commerce websites services versus buyers expectations: An empirical analysis of the online marketplace. International Journal of Information Technology & Decision Making, 12 (4), 651-677.

D'Avanzo, E. & Pilato, G. (2015). Mining social network users opinions' to aid buyers' shopping decisions. Computers in Human Behavior, 51, 1284-1294.

Dave, K., Lawrence, S. & Pennock, D. M. (2003). Mining the Peanut Gallery: Opinion Extraction and Semantic Classification of Product Reviews. Proceedings of the 12th International Conference on World Wide Web, (pp. 519-528). New York, NY, USA

Dijkmans, C., Kerkhof, P. & Beukeboom, C. J. (2015). A stage to engage: Social media use and corporate reputation. Tourism Management, 47, 58-67.

Feldman, R. & Sanger, J. (2013). The text mining handbook: Advanced approaches in analyzing unstructured data. Cambridge University Press.

Gamon, M., Aue, A. & Corston-Oliver, S. (2005). Pulse: Mining customer opinions from free text. In Advances in Intelligent Data Analysis VI (pp. 121-132). Berlin, Heidelberg: Springer Berlin Heidelberg.

Ganu, G., Elhadad, N. & Marian, A. (2009). Beyond the stars: Improving rating predictions using Review Text Content. Twelfth International Workshop on the Web and Databases (WebDB 2009), 9, (pp. 1-6). Providence, Rhode Island, USA.

Gligorijevic, B. & Luck, E. (2012). Engaging Social Customers – Influencing New Marketing Strategies for Social Media Information Sources. In Contemporary Research on E-business Technology and Strategy (pp. 25-40). Springer Berlin Heidelberg.

Gruen, T. W., Osmonbekov, T. & Czaplewski, A. J. (2006). eWOM: The impact of customer-to-customer online know-how exchange on customer value and loyalty. Journal of Business Research, 59 (4), 449-456.

Hovy, E. & Lavid, J. (2010). Towards a Science' of Corpus Annotation: A New Methodological Challenge for Corpus Linguistics. International Journal of Translation, 22 (1), 13-36.

Hu, M. & Liu, B. (2004). Mining and summarizing customer reviews. Proceedings of the tenth ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, KDD '04 (pp. 168-177). New York, NY, USA: ACM.

Hudson, S., Roth, M. S., Madden, T. J. & Hudson, R. (2015). The effects of social media on emotions, brand relationship quality, and word of mouth: An empirical study of music festival attendees. Tourism Management, 47, 68-76.

Jagtap, V. & Pawar, K. (2013). Analysis of different approaches to Sentence-Level Sentiment Classification. International Journal of Scientific Engineering and Technology, 2 (3), 164-170.

Kotzias, D., Denil, M., De Freitas, N. & Smyth, P. (2015). From Group to Individual Labels using Deep Features. KDD '15 Proceedings of the 21th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, (pp. 597-606). New York, NY, USA: ACM.

Kim, S. M. & Hovy, E. (2006). Automatic identification of pro and con reasons in online reviews. In ACL (pp. 483-490). Morristown, NJ, USA: Association for Computational Linguistics.

Kim, J., Choi, D., Hwang, M. & Kim, P. (2014). Analysis on Smartphone Related Twitter Reviews by Using Opinion Mining Techniques. Advanced Approaches to Intelligent Information and Database Systems, Studies in Computational Intelligence, 551, 205-212.

Liu, B. (2012). Sentiment Analysis and Opinion Mining. Morgan & Claypool Publishers.

Medhat, W., Hassan, A. & Korashy, H. (2014). Sentiment analysis algorithms and applications: A survey. Ain Shams Engineering Journal, 5, 1093-1113.

Montoyo, A., Martinez-Barco, P. & Blahur, A. (2012). Subjectivity and Sentiment Analysis: An Overview of the current state of the area and envisaged developments. Decision Support Systems, 53 (4), 675-679

Pang, B., Lee, L. & Vaithyanathan, S. (2002). Thumbs up?: Sentiment classification using machine learning techniques. Proceedings of the ACL-02 Conference on Empirical Methods in Natural Language Processing. 10, (pp. 79-86). Stroudsburg, PA, USA: Association for Computational Linguistics.

Park, D.-H., Lee, J. & Han, I. (2007). The Effect of On-Line Consumer Reviews on Consumer Purchasing Intention: The Moderating Role of Involvement. International Journal of Electronic Commerce, 11 (4), 125-148.

Pawle, J. & Cooper, P. (2006). Measuring emotion e lovemarks, the future beyond brands. Journal of Advertising Research, 46 (1), 38-48.

Pitt, L. F., Berthon, P. R., Watson, R. T. & Zinkhan, G. M. (2002). The Internet and the birth of real consumer power. Business Horizons, 45 (4), 7-14.

Pustejovsky, J. & Stubbs, A. (2012). Natural Language Annotation for Machine Learning. O'Reilly Media, Inc.

Tripp, T. M. & Grégoire, Y. (2011). When Unhapy Customers Strike Back on the Internet. MIT Sloan Management Review, 52 (3), 37-44.

Turney, P. D. (2002). Thumbs up or thumbs down?: semantic orientation applied to unsupervised classification of reviews. Proceedings of the 40th Annual Meeting on Association for Computational Linguistics (pp. 417–424). Stroudsburg, PA, USA: Association for Computational Linguistics.

Wiebe, J., Wilson, T. & Cardie, C. (2005). Annotating expressions of opinions and emotions in language. Language Resources and Evaluation 39 (2-3) 165-210.

Weiss, S. M., Indurkhya, N. & Zhang, T. (2010). Fundamentals of Predictive Text Mining. London: Springer-Verlag.

Wiegand, M. & Klakow, D. (2009). The role of knowledge-based features in polarity classification at sentence level. Proceedings of the 22nd International Florida Artificial Intelligence. AAAI Press.

Wilson, T., Wiebe, J. & Hoffmann, P. (2005). Recognizing contextual polarity in phrase-level sentiment analysis. Proceedings of the Conference on Human Language Technology and Empirical Methods in Natural Language Processing, HLT '05 (pp. 347-354). Stroudsburg, PA, USA: Association for Computational Linguistics.

Yang, C., Chen, Z., Wang, T. & Sun, P. (2015/1). Research on the Sentiment Analysis of Customer Reviews Based on the Ontology of Phone. International Conference on Education, Management and Computing Technology, (pp. 273-280).

Yang, C.-S., Chen, C.-H. & Chang, P.-C. (2015/2). Harnessing consumer reviews for marketing intelligence: a domain-adapted sentiment classification approach. Information Systems and e-Business Management, 13 (3), 403-419

Zirn, C., Niepert, M., Stuckenschmidt, H. & Strube, M. (2011). Fine-grained sentiment analysis with structural features. Proceedings of 5th International Joint Conference on Natural Language Processing. Chiang Mai, Thailand: Asian Federation of Natural Language Processing.

✉ **Correspondence**

**Olivera Grljević**

Faculty of Economics in Subotica
Segedinski put 9-11, 24000, Subotica, Serbia

E-mail: oliverag @ef.uns.ac.rs