

Review-Level Aspect-Based Sentiment Analysis Using an Ontology

Sophie de Kok
408695sk@student.eur.nl

Linda Punt
409135lp@student.eur.nl

Rosita van den Puttelaar
414662rp@student.eur.nl

Karoliina Ranta
413163kr@student.eur.nl

Kim Schouten
schouten@ese.eur.nl
Erasmus University Rotterdam
P.O. Box 1738, NL-3000 DR
Rotterdam, the Netherlands

Flavius Frasincar
frasincar@ese.eur.nl

ABSTRACT

The rapid growth of the World Wide Web has led to an explosion of information that is available on this platform. This has resulted in an increased interest in sentiment analysis, where the goal is to determine the opinion regarding a topic. Aspect-based sentiment analysis aims to capture the sentiment within a segment of text for mentioned aspects, rather than for the text as a whole. The task we consider is aspect-based sentiment analysis at the review-level for restaurant reviews. We focus on ontology-enhanced methods that complement a standard machine learning algorithm. For this task we use two different algorithms, a review-based and a sentence aggregation algorithm. By using an ontology as a knowledge base, the classification performance of our models improves significantly. Furthermore, the review-based algorithm gives more accurate predictions than the sentence aggregation algorithm.

CCS Concepts

•Information systems → Sentiment analysis; Information extraction; Web mining;

1. INTRODUCTION

The rapid growth of the World Wide Web has led to an explosion in the amount of information that is available on this platform [7]. As a result, the recognition of information retrieval as a value-added field has increased correspondingly. The Web has made it possible for consumers to share their opinions and experiences about products and services and they love to do so, according to a survey among more than 2000 American individuals [5]. The demand for user information is one of the major driving forces behind the

interest in opinion mining [9], where the goal is to determine the opinion of a group of people regarding a topic [17]. However, research states that information about goods and services is often missing or confusing, and the amount of it can be overwhelming [5]. An improved way of accessing consumer opinions is thus needed to aid businesses and consumers alike.

A commonly used way to extract information from review texts is to perform sentiment analysis, also referred to as opinion mining. This approach is defined as finding the quadruple (g, s, h, t) , where g represents the sentiment target, s the sentiment, h the opinion holder, and t the time at which the opinion was expressed [9]. A general approach is to take a whole sentence or review for g . A downside to this approach is that it only assigns a single polarity value to a sentence or a review. Consequently, it can not capture different sentiments within one segment of text. Rather than finding only the general sentiment of a document or a sentence, aspect-based sentiment analysis captures different aspects of the discussed entity and the sentiments expressed about these. For example, when dealing with restaurant reviews, a consumer can be positive about the service and be negative about the food quality. Therefore, aspect-based sentiment analysis allows for a more detailed analysis that utilizes more of the information provided by the available text [17].

Most aspect-based sentiment analysis approaches are concerned with two tasks, namely, aspect detection and aspect sentiment classification [17]. Aspect detection is defined as determining the different aspects of an entity in a particular part of the text, like a sentence or a review. For example, in the sentence “*Service is not what one would expect from a joint in this price category*”, ‘service’ is an aspect represented by this review. Sentiment classification assigns a sentiment to the aspects found in the text. In this example, the sentiment expressed about the aspect ‘service’ is negative.

A method that has been shown to perform well for aspect detection and sentiment classification is the Support Vector Machine (SVM) model [12]. When an SVM model is used, a feature vector of values is created for every instance to be classified and the model is taught using training data

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee.

SAC 2018, April 9–13, 2018, Pau, France

Copyright 2018 ACM 978-1-4503-5191-1/18/04...\$15.00

<https://doi.org/10.1145/3167132.3167163>

to interpret these values. Machine learning methods such as SVM models need a substantial amount of training data to obtain acceptable accuracy, as individual elements have little predictive value [2]. This reveals the need for large training data, but large amounts of annotated data are often not available for a product or service. To address this problem, we hypothesize that using external information in the form of an ontology could lessen the need for training data. Ontologies use a shared vocabulary for a certain domain and include axioms to define the relationships between different domain concepts [4]. These axioms can help derive information that is only implicitly stated. For this reason, the employment of an ontology is expected to improve the performance of sentiment analysis [13] and lessen the need for training data.

In our paper, we focus on the second sub-problem of aspect-based sentiment analysis, namely sentiment classification. Many solutions to the first task have already been provided, for example [6, 18]. We focus on sentiment classification for entity aspects that appear across a review (review-level aspect-based sentiment analysis) [14]. By using an ontology as a knowledge base, we can define concepts and relationships which could help in performing our task. For example, by knowing that someone liked the pasta we can infer that the food was liked, as pasta is a type of food. We further choose to use a linear SVM as our machine learning algorithm because it deals well with large amounts of features [3]. In this paper we hypothesize that review-level aspect-based sentiment analysis using an ontology gives better results than methods which do not include the use of an ontology. Furthermore, we investigate two approaches for review-level aspect-based sentiment analysis, one at the review level and one that aggregates sentence level sentiment label predictions. We expect that the review-level approach gives better results, because reviewers tend to write in interconnected sentences. Lastly, we consider the proposition that less training data is needed when we include an ontology in our model.

This paper is structured as follows. First, we discuss the related work in Sect. 2, followed by information about the used data set in Sect. 3. Then, in Sect. 4 and 5 we explain the proposed methodology and analyze the performance of our algorithms. Last, Sect. 6 contains our concluding remarks and possible directions for future work.

2. RELATED WORK

In this section we discuss work related to the field of sentiment analysis. Firstly, [17] provides a survey on aspect-level sentiment analysis. For this task [17] considers three different types of methods: dictionary-based, supervised machine learning, and unsupervised machine learning. In this paper we use a supervised machine learning method. We do so because we have supervised data available and supervised methods work better than unsupervised ones.

Various studies investigate whether the inclusion of an ontology improves results. In [18], the focus is on a knowledge-based approach that complements standard machine learning algorithms. The authors of [18] enhance the sentiment analysis using domain ontology information. By incorpo-

rating common domain knowledge into an ontology, classification performance for both aspect detection and aspect sentiment classification can be improved. The authors found words within sentences that appear in the ontology and are related to the aspect under consideration. They then provided all the superclasses of the ontology concept to the employed machine learning algorithm for the classification tasks. For both classification tasks, [18] works with an existing classifier, the linear Support Vector Machine. Contrary to [18], which focuses on aspect-based sentiment analysis at the sentence-level, this paper considers aspect-based sentiment analysis at the review-level. Furthermore, we enhance the ontology application using additional ontology related features such as synonyms.

[19] proposes an approach that the authors call HL-SOT. HL-SOT is a hierarchical learning (HL) process in combination with a sentiment ontology tree (SOT). A sentiment ontology tree has a tree-like appearance and the complete SOT consists of numerous sub-SOTs. A SOT has an attribute root node that has two leaf children which each represent a sentiment (positive/negative) associated with the attribute. Each sub-SOT represents a sub-attribute and is given as a child of the root node of the parent attribute SOT. Furthermore, each sub-SOT is assigned its own classifier with its own threshold value. The ontology is used by the authors to ensure that a text is labeled to contain an attribute only if all its parent attributes have also been mentioned within the same text segment. The proposed approach, however, can be disadvantageous when considering short reviews as these may express opinions on certain attributes without the mention of parent attributes. Unlike [19], we do not use different classifiers for different attributes. However, we do take into account that the sentiment polarity of some words is dependent on the product attribute being described.

Last, in [8] a system is considered in which an individual can search for information on a specific product. The authors suggest using an ontology to improve this system. They recommend a procedure in which the ontology is used as an alternative representation of a domain specific sentiment lexicon. The described ontology contains products, product features, sentiment words that are specific to a product feature, and the associated polarity. This sentiment ontology is then used in combination with manually crafted sentiment lexicons and NLP rules to determine the polarity of a product feature and sentiment word pair. This results in an accuracy comparable with previous works utilising machine-learning techniques. The authors of [8] compare their ontology-based approach to machine learning techniques, however, in this paper we combine both these approaches.

3. DATA

In this paper, we use a dataset of restaurant reviews from SemEval 2016 [14]. The dataset consists of training data and test data. The training data is used to develop and train our machine learning algorithm and the test data is used to evaluate the performance of our algorithm. We define a **notion** as an aspect category paired with a review (or sentence) in which it is mentioned. Each **notion** has a textual unit which contains the text of the review (or sentence). Our training data contains 335 reviews with 1435

review-level `notion` instances and 2455 sentence-level `notion` instances. The test data contains 90 reviews with 404 review-level `notion` instances and 859 sentence-level `notion` instances.

Our main task is to determine aspect sentiment polarities at the review level. To compare the review-based and sentence-aggregation approaches, we use data that is annotated for both reviews and sentences with respect to aspect sentiment. Each review (sentence) in the dataset is annotated with its occurring aspect categories, we do not differentiate between aspects and their categories, and the corresponding sentiment polarities.

Each review in the dataset can contain multiple aspect categories and each of these is labeled as positive, neutral, negative, or conflict. An aspect is labeled as ‘conflict’ in the case of conflicting opinions. Each aspect mentioned in the review has a unique sentiment, however an aspect can be mentioned multiple times in a review with different sentiments. In this case, all different sentiments are taken into account to assign an appropriate label. The sentences in the dataset can also contain multiple aspect categories, however, contrary to the review level, the aspect categories are labeled as positive, neutral, or negative. At the sentence level, the dataset is not annotated with the conflict label. This is presumably due to the small size of a sentence which makes the appearance of conflicting opinions less likely.

In Fig. 1 and 2 we show some statistics related to the aspects and sentiments in our dataset annotated for reviews. Fig. 1 shows the relative frequency of each aspect category in the data. Each review is assigned an overall sentiment label about the restaurant, therefore the aspect category `RESTAURANT#GENERAL` has a frequency of 100%. Fig. 2 shows the distribution of the sentiment values. We see that in both the training and test data the distribution of the sentiment values is unbalanced, as the positive label appears more frequently than the other sentiment labels.

4. METHOD

For the sentiment classification we use a linear Support Vector Machine (SVM). For the review level we train a multiclass SVM model with the classes: positive, negative, neutral, and conflict. For the sentence level we train a multiclass SVM model with the classes: positive, negative, and neutral. We use an SVM model because it has shown good results for sentiment analysis in the past [12].

Before the available data can be used for aspect sentiment classification it has to be pre-processed first. For this we use the Stanford CoreNLP Natural Language Processing Toolkit [10]. The first step in the processing of the data is tokenization. By tokenizing the data, we can separate meaningful terms and characters in the text from each other and remove white spaces. After tokenization, Part-of-Speech tags, such as ‘noun’ and ‘adjective’, are attached to the words of the sentences. In order to be able to recognize different forms of a word as the same, we lemmatize the words. This means that we find the dictionary form of a word. The last step is to parse the data, which determines the grammatical structure of sentences. This information can later be used to determine related words. Our proposed

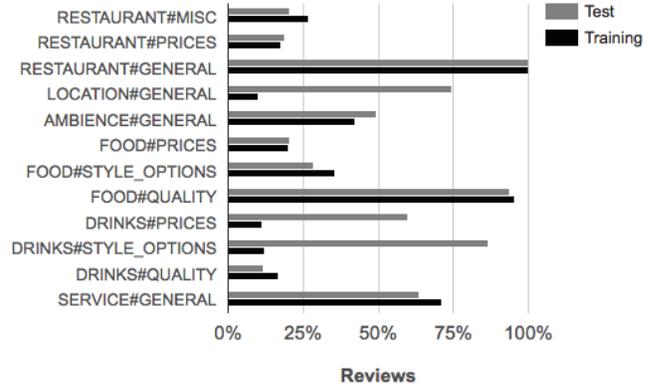


Figure 1: Relative frequencies of each aspect category label

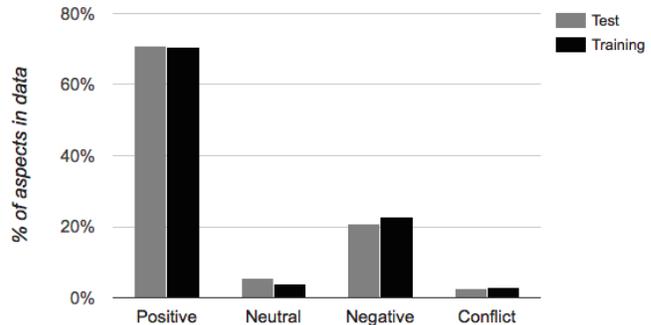


Figure 2: Relative frequencies of each sentiment label

algorithms then use `notion` instances, an aspect category paired with a review (or sentence) in which it is mentioned, from this pre-processed data.

4.1 Ontology

The restaurant ontology¹ consists of three main classes: *Entity*, *Property*, and *Sentiment*. Our first main class is the *Entity* class. This class contains terms pertaining to the domain of restaurant reviews. Its subclasses are *Ambience*, *Experience*, *Location*, *Person*, *Price*, *Restaurant*, *Service*, *StyleOptions*, and *Sustenance*. Most of these classes represent one or more aspect categories. If relevant, the property *aspect* connects the class to the corresponding aspect category as annotated in the data set (e.g., `FOOD#QUALITY`). The mentioned subclasses of *Entity* further have their own subclasses dividing them into more specific aspect categories.

Property is our second main class and it is divided into numerous subclasses containing descriptive terms that can be encountered in the restaurant domain. We created subclasses of terms that describe general negative and positive properties that can be related to several *Entity* subclasses (e.g., *GenericPositiveProperty*), and subclasses that represent adjectives describing characteristics of only one *Entity* subclass. For example, the class *AmbienceNegativeProperty*

¹The used ontology can be downloaded at www.kimschouten.com/publications/#sac2018

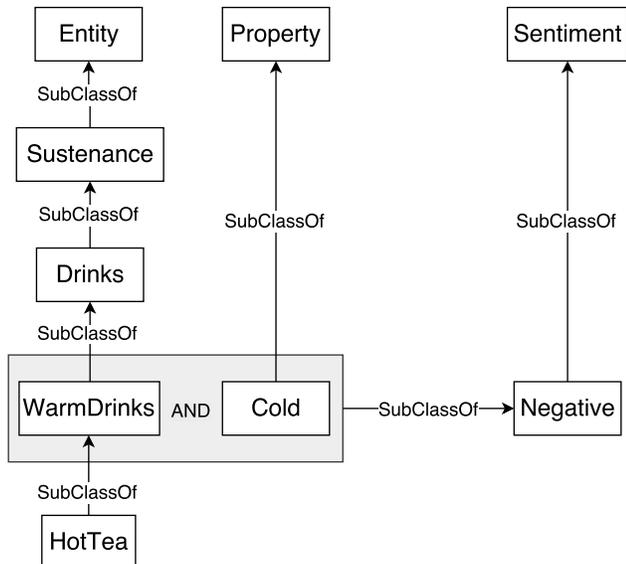


Figure 3: Ontology snippet

is a subclass of *Property* and a subclass of *Ambience*.

Our last main class is *Sentiment*. *Sentiment* is the superclass of the various polarity values, excluding conflict. This class has positive, negative, and neutral as subclasses. Subclasses of the *Property* class that represent positive (negative) properties are also a subclass of the *Positive* (*Negative*) class. Yet subclasses of the *Property* class representing adjectives that are more context specific are not subclasses of the different sentiment classes (e.g., *Cold*). Such classes, however, in combination with an entity (e.g., *WarmDrinks*), may be a subclass of one of the sentiment classes. This can be seen in Fig. 3. Here *HotTea* is a subclass of *WarmDrinks* and *Cold* is a subclass of *Property*. The intersection of *WarmDrinks* and *Cold* is a subclass of *Negative*, as warm drinks should not be cold. Thus, because *HotTea* is a subclass of *WarmDrinks*, the intersection of *HotTea* and *Cold* is a subclass of *Negative*.

The majority of our classes have the *lex* property attached. This property links the class to the associated lexicalizations. These lexicalizations can later be used to search for the presence of a concept in the ontology.

To extend our ontology we use refined lists of terms² pertaining to our domain and available on the Web and add these terms as subclasses of the corresponding classes in our ontology. For example, we expand our *AmbienceNegativeProperty* class with a list of negative adjectives describing ambience. We further augment our ontology by adding words that frequently occur within our training data and are relevant to the restaurant domain.

Let us consider an example to illustrate the workings of our ontology: the word ‘cramped’ is related via the *lex* property to *Cramped* and *Cramped* is a subclass of *AmbienceNegativeProperty*. *AmbienceNegativeProperty* is a subclass of *Ambience* and *Ambience* is related to the aspect AMBIENCE#GENERAL via the *aspect* property. There-

fore, we know that the word ‘cramped’ refers to AMBIENCE#GENERAL. Furthermore, because *AmbienceNegativeProperty* is also a subclass of *Negative*, we know that ‘cramped’ expresses a negative sentiment about AMBIENCE#GENERAL. Therefore, if the word ‘cramped’ is found in the review, the aspect AMBIENCE#GENERAL is labeled negative. This can be formally represented as follows:

$$\begin{aligned}
 Cramped &\equiv \exists lex.\{“cramped”\} \\
 Cramped &\sqsubseteq AmbienceNegativeProperty \\
 AmbienceNegativeProperty &\sqsubseteq Ambience \\
 AmbienceNegativeProperty &\sqsubseteq Negative \\
 Ambience &\sqsubseteq \exists aspect.\{“AMBIENCE#GENERAL”\} \\
 \text{Thus ‘cramped’ is negative about } &AMBIENCE\#GENERAL.
 \end{aligned}$$

4.2 Algorithms

We distinguish between two different algorithms. The first algorithm is review-based and uses a linear multiclass SVM model. The SVM determines the sentiment of the aspect categories in the review based on a feature vector. The aspect categories can be labeled as positive, neutral, negative, or conflict. Per **notion** we create a new feature vector instance. Our second algorithm is a sentence aggregation algorithm and a more refined method for the prediction of the aspect sentiments in reviews. We once again use a linear multiclass SVM, though now with the classes positive, negative, and neutral. Contrary to the review-based algorithm, we predict the sentiment of aspects in a single sentence instead of a review. Using these predictions, we use an aggregation algorithm that sums up the predicted polarities of each aspect per review. Thus, if a review has for example five sentences, where in three of them the FOOD#QUALITY aspect appears, we sum up the predicted polarity of these three sentences. The pseudocode of the summation algorithm is given in Alg. 1. The different features of both algorithms are discussed in the next section.

4.3 Model Features

Our SVM models use a variety of features to determine the sentiment classification. We use several procedures to construct these features, which we can split into two groups: the *feature generators*, which each create one or more features, and the *feature adaptors*, which adjust existing features. Some of these are independent from the employed ontology while others stem from the use of the ontology. We determine which features to include in our final models by comparing the average F_1 score for different model feature combinations using the training data with 10-fold cross-validation. Furthermore, for each model we optimize the SVM complexity parameter c over the range 10^{-6} to 10^3 with steps of 10^1 .

4.3.1 Feature Generators

The following feature generators are **independent** from the ontology:

Aspect. In the data, we have the aspects that are mentioned within each review (or sentence). Thus, for each no-

²<http://quizlet.com>
<http://www.macmillandictionary.com>

Algorithm 1: Sentence Aggregation Algorithm

Let r be a review and a an aspect. Then we define $S_{r,a}$ as the set of sentences for review r and aspect a .

Furthermore, we denote the polarity of a review for a certain aspect as $p_{r,a}$.

Data: $S_{r,a}$ set of sentences for review r and aspect a

Result: $p_{r,a}$ the polarity for aspect a of review r

begin

```

  score  $\leftarrow$  0
  posOrNeg  $\leftarrow$  false
  for  $s \in S_{r,a}$  do
    if computeSentiment( $s, a$ ) = positive then
      score  $\leftarrow$  score + 1
      posOrNeg  $\leftarrow$  true
    else if computeSentiment( $s, a$ ) = negative then
      score  $\leftarrow$  score - 1
      posOrNeg  $\leftarrow$  true
  if score > 0 then
     $p_{r,a} \leftarrow$  positive
  else if score < 0 then
     $p_{r,a} \leftarrow$  negative
  else if score = 0 and posOrNeg = false then
     $p_{r,a} \leftarrow$  neutral
  else
     $p_{r,a} \leftarrow$  conflict
  return  $p_{r,a}$ 

```

tion we use its corresponding aspect category as a feature in the SVM model.

Sentence count. This feature generator counts the number of sentences in a review-level **notion** and adds this value to the feature vector.

Lemma. This feature generator keeps track of the occurrence of words within the textual unit of a **notion**. For this item, all words within the data set are added to the SVM feature vector and the instance value is set equal to one if the word appears in the textual unit of the current **notion**, and zero otherwise.

The following feature generators are **dependent** on the ontology:

Ontology concepts. We inspect for each word in the textual unit of a **notion** whether it is a lexicalization of a concept in our ontology. If this is the case, we then find all superclasses of this class. If at least one of these superclasses has the property *aspect* linking it to the current aspect category (e.g., SERVICE#GENERAL), we add all the superclasses of the ontology concept related to this word as attributes to the feature vector with a value of one. By adding all the superclasses, we can make use of implicitly stated information.

Sentiment count. This feature generator counts the number of positive and negative text hits within the ontology. Thus, whenever one of the superclasses of a class associated to a word in the considered textual unit is the *Positive* or *Negative* class we update the respective counter.

4.3.2 Feature adaptors

The following feature adaptors are **dependent** on the ontology:

Ontology concept score. This feature adaptor influences the *ontology concepts* feature generator. We set the value for all superclasses to one, like in *ontology concepts*, however, superclasses that have the current category (e.g., SERVICE#GENERAL) as a value for the *aspect* property get a larger importance score. We denote this importance score with the parameter m and we determine the value of this parameter using optimization. By assigning a larger value to the superclasses directly related to the current category, the SVM can take into account that these superclass features are more important than superclasses that are not directly related to the current category.

Negation handling. It is common that reviewers make use of expressions such as “not bad”. In order to account for this we adapt the feature generator *ontology concepts*. If a word, that is a hit in our ontology, has a negating word directly preceding it, we replace all positive (negative) superclasses associated with the word with their negative (positive) counterpart. In this way when expressions such as “not bad” are used, we correctly specify it as a positive expression rather than a negative one.

Synonyms. Since our ontology adds useful information to our SVM feature vectors, we want to increase the number of relevant words that occur as lexical representations of concepts in our ontology. For this we use synonyms from WordNet [11] to complement the feature generator *ontology concepts*. If a word in the textual unit does not appear as a lexicalization in our ontology, we check if one of its synonyms is included. If this is the case, we add the superclasses associated with the synonym that does occur in the ontology to the feature vector. Word-sense disambiguation, identifying the meaning of a word within its context, is included in the design of the *synonyms* feature adaptor is implemented. Since only words that correspond to the restaurant domain specific meaning are included in our ontology, synonyms that do not relate to the restaurant domain are ignored. For example, the word ‘starter’ may appear in the textual unit of a **notion**, however, ‘starter’ does not appear as a lexicalization in the ontology. We thus consider the synonyms of ‘starter’. The word ‘starter’ has, among others, the synonyms ‘newcomer’ and ‘appetizer’, yet only ‘appetizer’ appears in our ontology. Therefore, the restaurant domain specific meanings are automatically selected. We assume that a word is used with only one meaning (the domain related one) in our domain text.

Weight. In order to take into account that some words are less important than others, we adjust *ontology concepts* generator by determining weight scores for every word that appears in the data. In the calculation of the weight scores we take into account that words that frequently appear in a review, but also frequently appear in all reviews, are less important than words that do not frequently appear in all reviews. This is called term frequency-inverse document frequency (TF-IDF). We use the following formula to determine the weight score of a word [16]:

$$weightScore = tf_{t,r} \cdot \log \frac{N}{df_t}, \quad (1)$$

where $tf_{t,r}$ is the frequency of the term t in the current review r , N is the total number of reviews, and df_t is the number of reviews in which the term t appears. We take the natural logarithm because words that appear ten times more often are not necessarily ten times more important. Thus, the logarithm helps to scale down the importance of the term. When the *weight* property is applied, the instance value of each superclass in *ontology concepts* is replaced by:

$$\max_i w_{s,i}, \quad (2)$$

where $w_{s,i}$ is the weight of word i , which is a lexicalization of a class in the ontology that has s as one of its superclasses. We take the maximum as we do not want a superclass to count more heavily when it appears more frequently. When *weight* is applied in combination with the *Ontology concept score* feature adaptor, the instance value of each superclass is set equal to:

$$\sum_i w_{s,i}, \quad (3)$$

where $w_{s,i}$ is as described above. However, when a superclass has the current aspect category as the value for the *aspect* property, the instance value described in Eq. 3 is multiplied by the importance score m . In this case we take the summation of the weight scores because if multiple ontology concepts are related to the current aspect category, we want the ontology concept that appears more often to have a larger score.

Word window. Rather than using the whole textual unit of a *notion* to create features, we determine a set of word windows. We initially iterate over all the words in the textual unit and when a word (or a synonym of it, when applied in combination with the feature adaptor *Synonyms*) appears as a lexicalization in our ontology, we determine a window of related words that are at most $k+1$ grammatical steps away from the original word. We determine these grammatical steps using the dependencies found during the pre-processing of the data. The value of k is optimized. We then use the word windows to create the features related to that *notion*. To illustrate the concept of a word window, consider the word ‘prices’ which appears in the sentence “Prices too high for this cramped and unappealing restaurant.”. The word window surrounding ‘prices’ is [Prices, too, high, restaurant,], where we, for this example, assume that k is equal to one.

4.4 Model Evaluation

To evaluate our models we calculate the accuracy, which is equal to the F_1 score. When an instance is correctly predicted, we define this as true positive. False positives and false negatives are both found when the predicted sentiment value is incorrect. Last, to compare models with each other we use a two-sample, two-tailed paired t-test.

5. EVALUATION

In this section we present and discuss our results. We hypothesize that the use of an ontology results in an improved

Table 1: Model performance for the aspect sentiment classification at the review level. P-value is given for a paired two-sided t-test.

	10-fold cross-valid.		p-value	in-sample	out-of-sample
	avg. F_1	st.dev.	vs. base	F_1	F_1
base	0.7852	0.0524	-	0.8718	0.8020
final	0.8001	0.0506	<0.0001	0.8753	0.8119

sentiment prediction accuracy. For this we describe and compare the predictive capabilities of our models. We also expect that the review-based algorithm outperforms the sentence aggregation algorithm. Furthermore, we hypothesize that the addition of an ontology to aspect-based sentiment analysis lessens the need for training data, and thus decreases the data size sensitivity. To analyze this we run experiments with differing training data sizes. Last, we take a look at which features are most important to our models.

5.1 Performance

In Table 1, we compare our final review-level aspect-based sentiment analysis model, *final*, to a baseline model, *base*. Our *base* model is a combination of the feature generators *aspect*, *sentence count*, and *lemma*. While the *final* model uses the feature generators *aspect*, *sentence count*, *lemma*, *ontology concepts*, and *sentiment count*, and the feature adaptors *negation handling*, *synonyms*, and *weight*. For both the *base* and the *final* models we find the optimized complexity parameter to be $c = 0.1$.

For each model we present the average value and standard deviation of the F_1 -measure over five runs of 10-fold cross-validation using only the training data. Furthermore, we report the p-values of the two-sided paired t-tests comparing the predictive abilities of the different models. Finally, the right half of the table contains the F_1 scores for a single run of the models using the test data.

Table 1 shows that the model that includes the use of an ontology has significantly better performance than the base case for the task of aspect sentiment classification at the review-level. This result can be seen in both the 10-fold cross-validation using the training data, and in the single run using both the training and test data. The model with the ontology displays an increase in accuracy of approximately 1.0 percentage points in comparison to the baseline.

In Table 2 we compare our sentence-level aspect-based sentiment analysis model, *ontSL*, to a baseline model, *baseSL*. The *baseSL* model is a combination of the feature generators *aspect* and *lemma*. The *ontSL* model is a combination of the feature generators *aspect*, *lemma*, *ontology concepts*, and *sentiment count*, and the feature adaptors *ontology concept score*, *negation handling*, *synonyms*, *weight*, and *word window*, all applied to the sentence level. The optimal value of the importance score for the *ontology concept score* feature adaptor is equal to $m = 5.0$ and the optimal value of the parameter k for the *word window* feature adaptor is $k = 1$. The *baseSL* model and the *ontSL* model have an optimal complexity parameter of $c = 1$ and $c = 0.1$ respectively.

Similarly to Table 1, Table 2 shows the model performance for aspect sentiment classification, but now at the sentence level. This table shows that the model that includes the

Table 2: Model performance for the aspect sentiment classification at the sentence level

	10-fold cross-valid.		p-value	in-sample	out-of-sample
	avg. F_1	st.dev.		F_1	F_1
baseSL	0.7008	0.0513	-	0.8436	0.7229
ontSL	0.8217	0.0500	<0.0001	0.8811	0.7963

Table 3: Model performance for the aspect sentiment classification using the sentence aggregation algorithm

	10-fold cross-valid.		p-value	out-of-sample
	avg. F_1	st.dev.		F_1
baseSA	0.6897	0.0616	-	0.6824
ontSA	0.8130	0.0512	<0.0001	0.7717
	Gold value (upper bound):			0.9633

use of an ontology has significantly better performance than the base case. The approach that applies the ontology displays an increase in accuracy of approximately 7.3 percentage points.

In Table 3, we show the model performance for review-level aspect sentiment classification using the sentence aggregation algorithm. The *baseSA* model and the *ontSA* model in this table respectively use the predictions of the *baseSL* model and *ontSL* model in Table 2, in combination with the sentence aggregation algorithm. Similarly to the previous results, the model that includes the use of an ontology has significantly better performance than the baseline. The computed gold value is an upper bound on the F_1 score for the sentence aggregation algorithm. Instead of summing up the predicted sentiment polarities with the sentence summation algorithm, we aggregate the real sentiment polarities of the aspects in sentences, as present in the annotated data.

The results show that the review-based algorithm outperforms the sentence aggregation algorithm for the test data. The *final* model has an accuracy that is approximately 4.0 percentage points higher than the accuracy of the *ontSA* model. Table 4 shows the performance of our *final* model compared to SemEval submissions for this task [14]. In the table it is indicated whether the submission was constrained (C), used only the training data, or unconstrained (U), used other resources. The accuracy of our *final* model differs less than one percentage point from the top two submissions.

5.2 Data Size Sensitivity

Since we hypothesize that less training data is needed when we include an ontology in our model, we perform an experiment by training our *base* and *final* SVM models with a decreasing proportion of training data. The test dataset, however, remains the same for each run. In this way we can compare the F_1 scores over all the runs. We obtain the average F_1 scores over 10 single runs with randomly generated seeds. The result, shown in Fig. 4, is a mapping between the prediction performance and the proportion of training data.

The figure shows that the gap between the two lines remains approximately the same. Therefore, the use of an ontology does not influence the data size sensitivity. This is in line

Table 4: Ranking of proposed aspect sentiment classification methods SemEval-2016

Team	(Un)Constrained	Accuracy
UWB	Unconstrained	0.8193
ECNU	Unconstrained	0.8144
final	Unconstrained	0.8119
UWB	Constrained	0.8094
ECNU	Constrained	0.7871
bunji	Unconstrained	0.7055
bunji	Constrained	0.6658
GTI	Unconstrained	0.6411

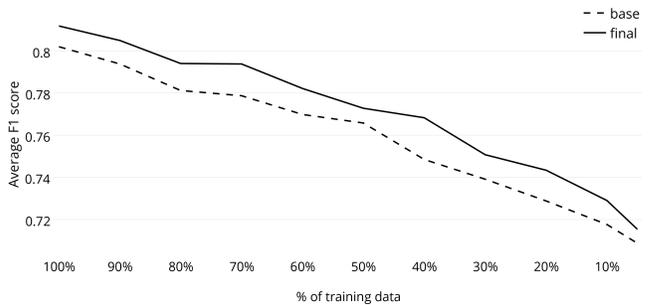


Figure 4: Data size sensitivity analysis of an ontology enhanced (*final*) and a non-ontology enhanced (*base*) approach

with the results reported in [18]. However, the ontology-enhanced method does, on average, perform better at all proportions of the training data when compared to the *base* model.

5.3 Feature Analysis

In Table 5, we list the ten most important features of our *final* model. The values listed represent the Information Gain (IG) [15]. The term reported next to the Information Gain states from which generator this feature originates, and next to that the name of the feature is given. For example, the first feature *numNegative* is generated by *sentiment count*, which counts the number of negative concept hits in the ontology.

As can be seen in Table 5 the majority of the most important SVM features are related to the negative sentiment class. As most of the **notions** within the dataset are labeled as positive, features that expose the negativity of a textual unit are important. Furthermore, we also calculate the internal attribute weights of the *final* SVM model. The 80 features with the largest weight are all ontology related features such as *Negative*, *Boring* and *Cozy*. This emphasizes the added value of an ontology.

6. CONCLUSION

In this paper, we investigated the added value of an ontology to the task of review-level aspect-based sentiment analysis. We proposed two different algorithms, a review-based algorithm and an aggregated sentence-based algorithm, which we both enhanced with the use of an ontology. For both algorithms, the accuracy is significantly higher when the ontology is used. The importance of the ontology is also supported by the results of the feature analysis. When com-

Table 5: Top 10 most important features for the *final* review-based model (*ontology concepts* is abbreviated to *ontology*)

0.2381	Sentiment count: numNegative
0.1159	Ontology: Negative
0.0821	Lemma: ‘not’
0.0638	Ontology: SustenanceNegativeProperty
0.0557	Sentiment count: numPositive
0.0539	Ontology: ServiceNegativeProperty
0.0522	Lemma: ‘do’
0.0517	Sentence count: numSentences
0.0515	Ontology: GenericNegativeProperty
0.0443	Lemma: ‘horrible’

paring the two algorithms, we see that, as expected, the review-based algorithm gives better results than the aggregated sentence-based algorithm. Furthermore, we hypothesized that the inclusion of an ontology would lessen the need for training data. However, contrary to our expectations, this is not the case. The ontology incorporating method is not less sensitive to the size of the training data than the method that does not incorporate an ontology, which is in line with earlier results [18]. A reason for this might be that the ontology features are not robust, meaning that the model needs training data to learn how to interpret them.

Since building an ontology manually is labor-intensive and time consuming, we would suggest to look into automating the process of creating the ontology for future work, for example, by extracting information from text [1]. This could make the use of an ontology more efficient when considering multiple domains. Furthermore, in our paper, we only assign the polarity values positive, neutral, negative, and conflict. However, an opinion can often not be categorized by merely four sentiment polarities. To account for this, sentiment scores could be assigned in order to extract the strength of an opinion.

7. REFERENCES

- [1] P. Buitelaar, P. Cimiano, and B. Magnini. *Ontology Learning from Text: Methods, Evaluation and Applications*. IOS Press, 2005.
- [2] E. Cambria. Affective Computing and Sentiment Analysis. *IEEE Intelligent Systems*, 31(2):102–107, 2016.
- [3] C.-C. Chang and C.-J. Lin. LIBSVM: A Library for Support Vector Machines. *ACM Transactions on Intelligent Systems and Technology*, 2(3):27:1–27:27, 2011.
- [4] S. Decker, M. Erdmann, D. Fensel, and R. Studer. Ontobroker: Ontology based access to distributed and semi-structured information. In R. Meersman, Z. Tari, and S. Stevens, editors, *Database Semantics: Semantic Issues in Multimedia Systems*, pages 351–369. Springer US, Boston, MA, 1999.
- [5] J. A. Horrigan. Online Shopping. *Pew Internet & American Life Project Report*, 36:1–24, 2008.
- [6] M. Hu and B. Liu. Mining Opinion Features in Customer Reviews. In *Proceedings of the 19th National Conference on Artificial Intelligence (AAAI 2004)*, volume 4, pages 755–760. AAAI, 2004.
- [7] Internet World Stats. [WWW Document, accessed: Aug 31, 2017] www.internetworldstats.com/stats.htm.
- [8] R. Y. Lau, C. C. Lai, J. Ma, and Y. Li. Automatic Domain Ontology Extraction for Context-Sensitive Opinion Mining. In *Proceedings of the 30th International Conference on Information Systems (ICIS 2009)*, pages 35–53. AIS Electronic Library, 2009.
- [9] B. Liu. *Sentiment Analysis and Opinion Mining, Synthesis Lectures on Human Language Technologies*, volume 16. Morgan & Claypool, 2012.
- [10] C. D. Manning, M. Surdeanu, J. Bauer, J. Finkel, S. J. Bethard, and D. McClosky. The Stanford CoreNLP Natural Language Processing Toolkit. In *Proceedings of 52nd Annual Meeting of the Association for Computational Linguistics: System Demonstrations*, pages 55–60. ACL, 2014.
- [11] G. Miller, R. Beckwith, C. Felbaum, D. Gross, and K. Miller. Introduction to WordNet: An On-Line Lexical Database. *International Journal of Lexicography*, 3(4):235–312, 1990.
- [12] T. Mullen and N. Collier. Sentiment Analysis using Support Vector Machines with Diverse Information Sources. In *Proceedings of the 9th International Conference in Empirical Methods on Natural Language Processing (EMNLP 2004)*, volume 4, pages 412–418. ACL, 2004.
- [13] J. Polpinij and A. K. Ghose. An Ontology-Based Sentiment Classification Methodology for Online Consumer Reviews. In *Proceedings of the IEEE/WIC/ACM International Conference on Web Intelligence and Intelligent Agent Technology (WI 2008)*, volume 1, pages 518–524. IEEE, 2008.
- [14] M. Pontiki, D. Galanis, J. Pavlopoulos, H. Papageorgiou, I. Androutsopoulos, and S. Manandhar. Semeval-2016 Task 5: Aspect Based Sentiment Analysis. In *Proceedings of the Tenth International Workshop on Semantic Evaluation (SemEval 2016)*, pages 27–35. ACL, 2016.
- [15] D. Roobaert, G. Karakoulas, and N. V. Chawla. Information gain, correlation and support vector machines. In I. Guyon, M. Nikravesh, S. Gunn, and L. A. Zadeh, editors, *Feature Extraction: Foundations and Applications*, pages 463–470. Springer Berlin Heidelberg, Berlin, Heidelberg, 2006.
- [16] G. Salton and M. J. McGill. *Introduction to Modern Information Retrieval*. McGraw-Hill, 1983.
- [17] K. Schouten and F. Frasincar. Survey on Aspect-Level Sentiment Analysis. *IEEE Transactions on Knowledge and Data Engineering*, 28(3):813–830, 2016.
- [18] K. Schouten, F. Frasincar, and F. de Jong. Ontology-Enhanced Aspect-Based Sentiment Analysis. In *Proceedings of the 17th International Conference on Web Engineering (ICWE 2017)*, volume 10360, pages 302–3320. Springer, 2017.
- [19] W. Wei and J. A. Gulla. Sentiment Learning on Product Reviews via Sentiment Ontology Tree. In *Proceedings of the 48th Annual Meeting of the Association for Computational Linguistics (ACL 2010)*, pages 404–413. ACL, 2010.