

2017

Aspect Identification and Sentiment Analysis in Text-Based Reviews

Sean Byrne
Lehigh University

Follow this and additional works at: <http://preserve.lehigh.edu/etd>



Part of the [Industrial Engineering Commons](#)

Recommended Citation

Byrne, Sean, "Aspect Identification and Sentiment Analysis in Text-Based Reviews" (2017). *Theses and Dissertations*. 2535.
<http://preserve.lehigh.edu/etd/2535>

This Thesis is brought to you for free and open access by Lehigh Preserve. It has been accepted for inclusion in Theses and Dissertations by an authorized administrator of Lehigh Preserve. For more information, please contact preserve@lehigh.edu.

Aspect Identification and Sentiment Analysis in Text-Based
Reviews

by

Sean Byrne

A Thesis

Presented to the Graduate and Research Committee

of Lehigh University

in Candidacy for the Degree of

Master of Science

in

Industrial & Systems Engineering

Lehigh University

May 2017

© Copyright 2017

Sean Byrne

This thesis is accepted and approved in partial fulfillment of the requirements for the Master of Science.

Date

Martin Takáč, Thesis Advisor

Tamás Terlaky, Chairperson of Department

Acknowledgements

I would like to express my gratitude to Professor Martin Takáč for his guidance and encouragement throughout my research and my time at Lehigh. I'd also like to thank my parents, Kevin and Jenifer Byrne, and my brothers Matthew and Jake for their constant support in everything I do. In addition, I'd like to thank Dr. Ali Yazdanyar, physician at Reading Hospital, for allowing me to apply what I've learned in a real-world setting.

Contents

Acknowledgements	iv
List of Figures	vii
List of Tables	viii
Abstract	1
1 Introduction	2
1.1 Automatic Aspect-Based Review System	2
1.2 Natural Language Processing	5
1.3 Application: Reading Hospital	6
2 Dataset Structure and Text Features	8
2.1 Datasets	8
2.2 Text Features	13
2.2.1 Token-level Features	13
2.2.2 Sentence-level Features	13
2.2.3 Review-level Features	14
2.2.4 Other Possible Features	14
3 Aspect Identification	16
3.1 Problem Description	16
3.2 Sequential Labeling: Conditional Random Fields	17
3.2.1 Labeling Method	18
3.2.2 Background: Naive Bayes and Maximum Entropy Models	18
3.2.3 Hidden Markov Models	20
3.2.4 CRF Model Description	22
3.2.5 CRF Training	23
3.2.6 CRF Evaluation	24
3.3 Association Mining Method (Hu and Liu)	29
3.3.1 Association Mining Method Description	29
3.3.2 Association Mining Method Evaluation	32
4 Aspect-Based Sentiment Analysis	34
4.1 Problem Description	34

4.2	VADER-based Method	35
4.2.1	Evaluation	37
4.2.2	Ratings-Based Evaluation	38
5	Conclusion	40
	Bibliography	42
	Appendix - Data Processing and Test Functions	45
	Appendix - Class Definitions	49
	Appendix - Aspect Identification	60
	Appendix - Sentiment Analysis	75
	Biography	81

List of Figures

2.1	An example of the dataset format in SemEval 2014.	10
2.2	An example of the dataset format in SemEval 2015.	11
2.3	An example of the dataset format in SemEval 2016.	12

List of Tables

3.1	The results for CRFs using distinct aspect terms.	26
3.2	The results for CRFs using instances of aspect terms.	27
3.3	The results for CRFs using distinct aspect terms across domains.	28
3.4	The results for CRFs using instances of aspect terms across domains.	28
4.1	The results using VADER on aspect terms in the Laptop domain.	38
4.2	The results using VADER on aspect terms in the Restaurant domain.	38
4.3	The results using VADER on aspect categories in the Restaurant domain.	38
4.4	True and predicted ratings for each category in the Restaurant domain.	39

Abstract

Online text-based reviews are often associated with only an aggregate numeric rating that does not account for nuances in the sentiment towards specific aspects of the review's subject. This thesis explores the problem of determining review scores for specific aspects of a review's subject. Specifically, we examine two important subtasks - aspect identification (identifying specific words and phrases that refer to aspects of the review subject) and aspect-based sentiment analysis (determining the sentiment of each aspect). We examine two different models, conditional random fields and an association mining algorithm, for performing aspect identification. We also develop a method for performing aspect-based sentiment analysis based on VADER, a sentence-level sentiment analysis algorithm built for sentiment analysis of social media. We identify key problem considerations, including other important subtasks and ideal training dataset qualities, for future development of a full aspect-based review system.

Chapter 1

Introduction

1.1 Automatic Aspect-Based Review System

Text-based reviews found online have become a common way to evaluate options when making a decision. These reviews span subjects from a variety of topics - products available for purchase online, downloadable applications, movie and music releases, restaurants, hotels, and more. Oftentimes, these reviews are associated with an overall numeric rating (typically on a 5-point or 10-point scale), which can be aggregated to form an average rating for a given subject. However, these ratings oftentimes hide the details present in the text of the reviews. For example, by examining a set of laptop reviews with an average rating of 3.0 out of 5.0, one might find that the screen of the laptop is mostly referred to positively, but the keyboard is mostly referred to negatively. This nuance is not reflected with an overall 5-point numeric rating, despite the fact that users oftentimes have preferences that require a more detailed view of the subject.

In order to more accurately reflect how reviewers feel about different aspects of a subject, it is desirable to develop a system to rate the major features of a subject separately, providing more meaningful information to those who may have specific preferences. A shopper looking to purchase a laptop, for example, may desire a high screen quality while not caring much about the processing power. This shopper would benefit from finding a laptop with a highly-rated score for the aspect "Screen" and may not mind if the laptop's overall score is dragged down by a lower rating for the aspect "Processing Power". It's possible that websites aiming to have a more comprehensive set of ratings could force users to rate specific qualities on a numeric scale, rather than just the overall product. However, this requires more effort on the end user, and ignores the vast amount of text-based review data that already exists.

One way such a system can be developed using existing product reviews is to utilize *sentiment analysis* (also known as opinion mining). Sentiment analysis attempts to derive measures of subjectivity from written text, typically labeling text using either the labels "subjective" and "objective" (ignoring polarity of subjective text), or the labels "positive", "negative", and "neutral" (where "positive" and "negative" are opposite categories of subjectivity, and "neutral" is equivalent to "objective"). Text-based reviews are an important source of data for sentiment analysis because they consist primarily of subjective opinions, making them particularly useful for building models with the ability to determine sentiment polarity.

However, rather than attempting to determine the sentiment of the review as a whole, the sentiment of particular attributes of the product would be measured. If a particular attribute is found to be associated with positive or negative polarity for most instances within a set of reviews, then it is given a high or low rating, respectively, for

that particular attribute. These attributes (or aspects) can be found through *aspect identification* - determining what words and phrases (terms) refer to specific aspects of the subject. For example, in the sentence "The battery life is quite strong and lasts all day long," the phrase "battery life" is an aspect term of the subject.

Once these aspect terms have been identified, sentiment analysis can be used to determine the sentiment polarity of each aspect term. Specifically, *aspect-based sentiment analysis* attempts to determine the sentiment of each aspect term. Accurately determining the polarity of aspect terms is more challenging than the typical sentiment analysis task. Sentiment analysis relies heavily on sentiment lexicons that classify adjectives based on their sentiment polarity, but an adjective that has a positive sentiment when used to describe one aspect may have a negative or neutral sentiment when used to describe another aspect. For example, "long" tends to have a positive sentiment when used to refer to "battery life" in a laptop, but a negative sentiment when used to refer to "wait times" at a restaurant. Another significant issue is when multiple aspect terms are mentioned within the same sentence. If one aspect has a positive sentiment and another has a negative sentiment, determining these sentiments accurately requires understanding which portions of the sentence apply to a given aspect term.

In the remainder of Chapter 1, we examine a brief background of natural language processing and mention an ongoing application of the methods described in this thesis. In Chapter 2, we describe the datasets used and qualities of a useful dataset for the problems of aspect term extraction and aspect-based sentiment analysis, as well as important features that can be derived from the text. In Chapter 3, we examine methods for extracting aspect terms from these datasets, and in Chapter 4, we examine methods for determining the sentiment of aspect terms.

1.2 Natural Language Processing

Natural Language Processing (NLP) is a field of study within computer science and artificial intelligence that focuses on analyzing and deriving meaningful information from human (natural) language. Natural language processing developed as a result of interest in machine translation (MT), the problem of translating sentences automatically from one language to another, in the 1950s. Research was severely limited due to the relatively undeveloped state of computers at the time. Initial research started as dictionary-based, with attempts to translate sentences word-for-word, but issues with determining the correct syntax (the arrangement of words) and semantics (the meaning of words) in translation quickly showed the limitations of such an approach. Despite technological limitations, research of this time period was able to effectively identify the importance of developing an explicit structure and definition for language that could allow methods to be generalized and implemented with computers [16]. The low quality of the methods developed, however, led a committee commissioned by the United States government called ALPAC (Automatic Language Processing Advisory Committee) to express doubts in the merit of continued MT research in a report in 1966 [8]. The committee suggested that significant improvements in computational linguistics was needed before MT could be effectively tackled, leading to a significant shift away from MT in the late 1960s. This shift allowed other problems within NLP to be explored, eventually leading to the broad range of problems studied within the field today.

The massive amount of data and processing power that are accessible today has opened the door to new heights in the world of natural language processing. Modern NLP research examines problems such as converting speech to text [12], answering text-based

questions [13], automatically summarizing large documents, automatic spell-checking, determining grammatical relationships between words, and much more. NLP has been utilized in a large variety of business contexts as well. Lawyers use NLP software to analyze large sets of legal documents to find meaningful information. Spam filters utilize NLP to find patterns within email text that indicate a high likelihood of being spam, and Google uses NLP in their language-translation software. Various social media sites utilize natural language processing so that advertisements can be customized to the interests of each user.

In this thesis, we utilize some commonly-used software for natural language processing. In particular, we make extensive use of the Natural Language Toolkit (NLTK) [6] and Stanford's CoreNLP toolkit [18]. These are packages for Python that provides a large set of functions and datasets useful for natural language processing.

1.3 Application: Reading Hospital

With the rise of electronic medical records, applications have started to appear within the medical field. Taking text-based data from the past (in this case, from physicians' notes) and data related to the eventual treatment of the patient's medical issues (for example, procedures done, diagnoses given, and success/failure rates), patterns can be found within the text of the doctors' notes. In this way, physicians' notes can be analyzed to determine signs of postsurgical complications, or to determine the procedure with the highest likelihood of success for a given diagnosis and set of physical traits.

One potential area for the application of natural language processing techniques is a project recently started at Reading Hospital. Because of the nature of this work, the

specific details cannot be shared in this paper. However, a basic problem outline can be shared. When a scan is done to examine a particular portion of the body, secondary information can be gathered. For example, in a CT scan where the primary objective is to examine a tumor, secondary nodules could be found on the scan that aren't directly related to the tumor. In this case, the doctor often suggests that the patient make a follow-up appointment with another practitioner; however, there is no easy way to connect the patient to the appropriate office for a follow-up appointment, and oftentimes patients end up ignoring these secondary findings until their next appointment months or even years later. Complications that could have been treated easily, if dealt with earlier, can end up becoming much more serious medical issues because of this.

The project's goal is to use NLP techniques to identify keywords in the notes of these scans that suggest a secondary finding should be examined or a follow-up appointment is needed. The details of these patients and scans could then be routed to the appropriate place automatically. In this way, the methods discussed in this paper can have a real impact on the patients at Reading Hospital.

Chapter 2

Dataset Structure and Text

Features

2.1 Datasets

There is a great deal of text-based review data available online - however, the raw data alone isn't enough. In order to perform the three major tasks associated with aspect-based sentiment analysis, the data provided must contain information about which words are aspect terms, which words are a part of which aspect categories, and whether each instance of a term is referred to positively or negatively. This requires human tagging of datasets, along with cross-validation measures to ensure that the tags are consistent across multiple people.

The difficulty of creating adequate data sources causes significant issues when tackling the problem of aspect-based sentiment analysis. It significantly limits the effectiveness of methods that rely heavily on domain-based features, since each subject (spanning

all categories of online products, media, restaurants, and others) may require a different set of training data for these methods to be effective. Thus, the importance of developing a model that is not overly reliant on the domain of the training data is particularly important.

SemEval (also known as the International Workshop on Semantic Evaluation) is "an ongoing series of evaluations of computational semantic analysis systems" hosted annually by SIGLEX (Special Interest Group on the Lexicon of the Association for Computational Linguistics) [1]. Each year, a set of tasks related to semantics within natural language processing are developed, with the goals of developing methods of discerning meaning from language and identifying issues worth exploring further. From 2014 to 2016, one of the tasks was "Aspect-Based Sentiment Analysis" [5] [4] [24]. In this task, participants were given data annotated with aspect terms, aspect categories, and their polarities. The goal of the task was to predict each of these for a set of testing data as accurately as possible.

We utilize datasets from the 2014-2016 SemEval competitions. They have been cross-validated to ensure that inter-annotator agreement is high [23], and there is data available from two different domains: laptop and restaurant reviews. The sentences in each year's data are largely the same, but the format they're stored in (as well as their aspect term and aspect category annotations) vary. In all formats, aspect terms and/or aspect categories are associated with a sentiment polarity from the set {"positive", "negative", "neutral"}, though the 2014 and 2016 formats also allow for a fourth "conflict" value that represents subjective statements without clear overall positive or negative sentiment.

```

<sentence id="2846">
  <text>Not only was the food outstanding, but the little 'perks' were great.</text>
  <aspectTerms>
    <aspectTerm term="food" polarity="positive" from="17" to="21"/>
    <aspectTerm term="perks" polarity="positive" from="51" to="56"/>
  </aspectTerms>
  <aspectCategories>
    <aspectCategory category="food" polarity="positive"/>
    <aspectCategory category="service" polarity="positive"/>
  </aspectCategories>
</sentence>

```

FIGURE 2.1: An example of the dataset format in SemEval 2014.

The 2014 datasets are stored as sentences (without review context) in two different domains: laptop reviews with 3,141 sentences and restaurant reviews with 3,145 sentences. Aspect terms are provided for sentences in the datasets of both domains, and aspect categories are provided for sentences in the dataset of the restaurant domain. The sentiment polarity fields in this dataset support the "conflict" value when the dominant sentiment polarity is not clear. For each aspect term, character offsets are provided (in two fields: "from" and "to", representing the beginning and end of the term, respectively) to identify the location of each aspect term within the sentence. Offsets start at index 0 within a sentence, and the "to" field stores the index of the offset immediately after the last character of the aspect term. Each sentence contains zero or more aspect terms and zero or more aspect categories.

The 2015 datasets are stored as reviews in two different domains: laptop reviews and restaurant reviews. Each review is provided as a list of sentences in order, and each sentence is associated with zero or more aspect categories. Aspect categories are stored as pairs of entities (E) and attributes (A) in the following format: "E#A". Entities are components of the overall topic - for example, entities in the set of Laptop reviews include "CPU", "Software", "Shipping", and "Support". Attributes are specific features or qualities of the entities - for example, attributes in the set of laptop reviews include "Price", "Quality", and "Portability". This dataset does not support the "conflict" value

```

<Review rid="720418">
  <sentences>
    <sentence id="720418:0">
      <text>Great Indian food and the service is incredible.</text>
      <Opinions>
        <Opinion target="Indian food" category="FOOD#QUALITY" polarity="positive" from="6" to="17"/>
        <Opinion target="service" category="SERVICE#GENERAL" polarity="positive" from="26" to="33"/>
      </Opinions>
    </sentence>
    <sentence id="720418:1">
      <text>The owner truly caters to all your needs.</text>
      <Opinions>
        <Opinion target="owner" category="SERVICE#GENERAL" polarity="positive" from="4" to="9"/>
      </Opinions>
    </sentence>
    <sentence id="720418:2">
      <text>When family came in he gave them apps to test their palets, and then ordered for them.</
      text>
      <Opinions>
        <Opinion target="NULL" category="SERVICE#GENERAL" polarity="positive" from="0" to="0"/>
      </Opinions>
    </sentence>
    <sentence id="720418:3">
      <text>Everyone was more then happy with his choices.</text>
      <Opinions>
        <Opinion target="NULL" category="SERVICE#GENERAL" polarity="positive" from="0" to="0"/>
      </Opinions>
    </sentence>
    <sentence id="720418:4">
      <text>Great food and the prices are very reasonable.</text>
      <Opinions>
        <Opinion target="food" category="FOOD#QUALITY" polarity="positive" from="6" to="10"/>
        <Opinion target="NULL" category="RESTAURANT#PRICES" polarity="positive" from="0" to="0"/>
      </Opinions>
    </sentence>
  </sentences>
</Review>

```

FIGURE 2.2: An example of the dataset format in SemEval 2015.

for sentiment polarity. For reviews in the Restaurant dataset, an opinion "target" can be specified - this happens when an entity E is explicitly referenced through a target word or phrase in the sentence. This allows aspect terms to be linked to aspect categories. The keyword "NULL" is used if an aspect category's entity is not explicitly referenced through a target. If an opinion target is specified, "from" and "to" fields are used to specify the location of the target within the sentence. These are set to 0 when the target is "NULL".

The 2016 dataset is provided in two different formats. One is identical to the 2015 dataset format. The other is a review-based format that stores sentences and aspect categories separately. Each review consists of a list of sentences and a separate list of the aspect categories within the review. This means that polarity ratings for each aspect category are review-level rather than sentence-level, and so each aspect category is assigned the sentiment polarity that is dominant within most sentences that contain

```

<Review rid="1032695">
  <sentences>
    <sentence id="1032695:0">
      <text>Every time in New York I make it a point to visit Restaurant Saul on Smith Street.</text>
    </sentence>
    <sentence id="1032695:1">
      <text>Everything is always cooked to perfection, the service is excellent, the decor cool and understated.</text>
    </sentence>
    <sentence id="1032695:2">
      <text>I had the duck breast special on my last visit and it was incredible.</text>
    </sentence>
    <sentence id="1032695:3">
      <text>Can't wait wait for my next visit.</text>
    </sentence>
  </sentences>
  <Opinions>
    <Opinion category="RESTAURANT#GENERAL" polarity="positive"/>
    <Opinion category="FOOD#QUALITY" polarity="positive"/>
    <Opinion category="SERVICE#GENERAL" polarity="positive"/>
    <Opinion category="AMBIENCE#GENERAL" polarity="positive"/>
  </Opinions>
</Review>

```

FIGURE 2.3: An example of the dataset format in SemEval 2016.

the aspect category. Aspect categories are defined in a similar way to the 2015 dataset, using an entity-attribute pair to represent each category. In cases where the dominant sentiment polarity is not clear, the polarity is defined as "conflict". Opinion targets are not provided.

The formats can be summarized as follows. The 2014 dataset identifies specific aspect terms and their associated polarities, as well as aspect categories for the Restaurant dataset that are not explicitly linked to aspect terms. The 2015 dataset is somewhat more specific - it identifies specific entity-attribute combinations that form aspect categories, as well as target aspect terms for the Restaurant dataset that explicitly link aspect terms to aspect categories. It also provides context information by grouping sentences by each review. The 2016 dataset is more general - it identifies specific entity-attribute combinations that form aspect categories that are found within a review as a whole, rather than individual sentences. By comparing methods across dataset formats with very similar data, the value of creating training datasets with a higher or lower level of detail can be found.

2.2 Text Features

2.2.1 Token-level Features

We break each sentence down into tokens consisting of words and punctuation using the Penn Treebank tokenizer within NLTK [20]. This tokenizer splits contractions (for example, “don’t” will become the separate tokens “do” and “n’t”) and stores punctuation as separate tokens.

Some features can be extracted from individual tokens without the need for information from the rest of the sentence or corpus. We store the original token text, as well as a lowercase version of the token. Several binary features are stored - whether or not the token is punctuation, whether or not it is in “titlecase” (the first letter of the token is capitalized, and the following letters are all lowercase), and whether the token is a digit. We use a popular word stemmer, PorterStemmer, to store the stem of a given word, removing all prefixes and suffixes from the token [25].

2.2.2 Sentence-level Features

Some features require sentence-level context. The index of each token within the sentence is stored, with 0 being the first token of the sentence. A part-of-speech (POS) tagger using the Penn Treebank tagset is used to tag the part-of-speech for each token in a sentence [20]. The full POS tag and the first 2 characters of the POS tag are stored as separate features, as the first two characters are indicative of a broader category that the following characters are part of (for example, “NN”, “NNP”, “NNS”, and “NNPS” are all tags to describe nouns). Each token also stores information about the previous and

next tokens in the sentence - the text, lowercase text, stem, and both POS tag features of the previous and next tokens, storing a default value if the previous or next token doesn't exist.

2.2.3 Review-level Features

Oftentimes, text-based reviews are associated with an overall numeric rating. Our datasets do not have contain numerical rating information, but utilizing these review ratings in an aspect-based sentiment analysis model may yield positive results, and is worth future consideration when designing annotated datasets from online reviews.

2.2.4 Other Possible Features

Many other features are commonly used for natural language processing purposes. WordNet is a lexical database designed to store words based on their word sense (the meaning of the word) rather than the word itself [21]. It contains over 155,000 words and 117,000 synonym sets (sets of words with the same meaning), with over 206,000 word-sense pairs in total [2]. Several other semantic relations are stored as well, such as antonyms. Hypernyms, a semantic “parent” of a given word, and hyponyms, semantic “children” of a given word, are stored - for example, the pair “sport” and “baseball” is a hypernym-hyponym pair. Meronyms and holonyms refer to component parts and the collective whole, respectively - for example, the pair “car” and “wheel” is a holonym-meronym pair. Using WordNet in a natural language processing model, particularly the problems of aspect identification and aspect-based sentiment analysis, would give the model a greater understanding of the relationships between words in a sentence. However,

WordNet has been found to not significantly impact the performance of text classification models [19], and the limited tests we performed showed little benefit. Despite this, usage of WordNet in other models for aspect identification and aspect-based sentiment analysis may still be worth exploring.

Word2Vec is a deep learning algorithm that takes sentences as inputs and outputs a vectorization of each distinct word within the training data. This can be used to determine the similarity of one word from another word. Word2Vec also allows for accurate operations among words, meaning syntactic and semantic patterns can be accurately generated. For example, suppose $\text{vec}(\text{word})$ is the Word2Vec vector representation of a word. $\text{vec}(\text{'brother'}) - \text{vec}(\text{'man'}) + \text{vec}(\text{'woman'})$ results in a vector similar to $\text{vec}(\text{'sister'})$. As a result, relationships among words are encoded in the vectors. Word2Vec was designed for massive datasets, ranging from tens of millions to billions of words, and so attempts to train Word2Vec on the datasets available (with only several thousand sentences available) were unsuccessful. Training Word2Vec on larger datasets available, such as the full English Wikipedia, has resulted in positive results in other aspect identification models [23].

Chapter 3

Aspect Identification

3.1 Problem Description

In some texts, particularly text-based reviews, there is an overall subject being discussed throughout the text. Aspect identification (or aspect term extraction) is the process of identifying what words and phrases (terms) refer to specific aspects of a subject in these texts.

Aspect identification typically refers to extracting aspect terms explicitly mentioned within the sentence, rather than implied terms. For example, the sentence "The restaurant was quite expensive" does not explicitly mention price, but "expensive" is an adjective referring to the price of the food, an implicit aspect within the sentence. We consider only explicit aspect terms in this paper.

An ideal system would not rely heavily on the domain of the training data, as otherwise a new set of training data would be required for each new domain examined. Identifying aspect terms requires human identifiers to manually record these aspect terms

and their sentiment, and requires a consistent approach so that these human identifiers mostly agree with each other. When each set of training data requires potentially hundreds of reviews (thousands of sentences), this task becomes infeasible to complete for the many domains available for text-based reviews.

One of the most significant challenges in aspect identification is balancing accuracy with robustness. The most accurate models will likely require more detailed training data - accurate sentence-level datasets identifying aspect terms and their respective polarities (positive, negative, or neutral). But the most domain-neutral models will rely on more general features and potentially unsupervised methods. Thus, we examine both supervised and unsupervised approaches, and test across domains to see how applicable each supervised method is to training data from a different domain.

3.2 Sequential Labeling: Conditional Random Fields

Aspect term extraction can be modeled as a sequence labeling problem, where each sentence is examined as a sequence of tokens, taking the context of an individual token into account. This framework is used for problems such as part-of-speech tagging, named entity recognition, and shallow parsing [26]. We describe and implement a common sequence labeling model called a Conditional Random Field (CRF), a generalization of another model called a Hidden Markov Model. These are sequential labeling models based on generalizations of the single-label models described with the naive Bayes classifier and Maximum Entropy models. The goal of a CRF is to determine the conditional distribution of potential labels (in our case, using the IOB2 tagging format) given the output (each token's text). Using the framework for Maximum Entropy models and

CRFs, feature functions can be defined that can allow a vector of output features to be associated with each word in a sentence.

3.2.1 Labeling Method

We use the IOB2 tagging format, where each token is associated with one of three labels - inside an aspect term ("I"), outside an aspect term ("O"), or the beginning of an aspect term ("B"). All aspect terms start with a "B", so only multi-token aspect terms utilize the "O" label.

3.2.2 Background: Naive Bayes and Maximum Entropy Models

The naive Bayes classifier is used to predict a class label y given a feature vector \mathbf{x} . It is based on the assumption of conditional independence of the individual features given the class label. The model attempts to maximize the joint probability $p(\mathbf{x}, y)$ of the features and the class label, which due to their conditional independence can be described as follows:

$$p(\mathbf{x}, y) = p(y) \prod_{i=1}^n p(x_i | y). \quad (3.1)$$

The Maximum Entropy classifier (also known as multinomial logistic regression) makes the assumption that $\log(p(y | \mathbf{x}))$ can be represented as a linear combination of the features and a constant. This is useful in that the features are not assumed to be independent, and so the relationships among the output features are considered. The Maximum Entropy classifier models the conditional probability $p(y | \mathbf{x})$ as follows:

$$p(y | \mathbf{x}) = \frac{1}{Z} \exp(\beta_{\mathbf{y}} \mathbf{x} + \beta_{y,0}). \quad (3.2)$$

$Z = \sum_y \exp(\beta_{\mathbf{y}} \mathbf{x} + \beta_{y,0})$ is a normalization constant which adjusts to ensure valid probabilities. The parameters $\beta_{\mathbf{y}}$ and $\beta_{y,0}$ can be chosen based on the training data using the expectation-maximization (EM) algorithm [11].

Naive Bayes is a generative model, meaning that the model estimates the joint probability distribution of the state and the feature vector and uses this learned distribution to predict the likelihood of a feature vector \mathbf{x} being assigned a class label y . Maximum Entropy models, on the other hand, are discriminative - they learn the conditional probability $p(\mathbf{y} | \mathbf{x})$ of being in a state \mathbf{x} given an output \mathbf{y} . This is important because unlike generative models, the probability distribution of outputs $p(x)$ does not need to be learned. In the case of natural language processing where the observed outputs are words, there are almost certainly words that don't exist in the training corpus that may occur when using the model, meaning $p(x)$ cannot be accurately estimated without training data that contains every possible word - an unfeasible task.

Because these classifiers only predict a single class label for a set of features, they cannot model the relationships among the hidden states. Graphical models such as Hidden Markov Models and CRFs, on the other hand, are able to account for the dependencies between the nodes' labels.

3.2.3 Hidden Markov Models

One model for labeling sequences of inputs is called a Hidden Markov Model (HMM). The system is assumed to be a Markov process, where the state of the current node is dependent only on the state of the previous node in the sequence. However, instead of observing the state of a given node directly, we observe an output that is dependent on the state, and each state has a probability distribution over the set of outputs. HMMs also assume conditional independence of the output features given each node's state, making them a generalization of the Naive Bayes classifier. Given a sequence of outputs and information about each state's distribution of possible outputs, we can predict a sequence of hidden states.

In our problem, the sequence of words or tokens within the sentence is the sequence of outputs, and the sequence of labels, using the IOB2 standard, are the hidden states. Our goal is to predict the IOB2 labels of each token within a sentence, using the sentence's tokens as the sequence of output features.

Let $X = (x_1, x_2, \dots, x_n)$ be the sequence of observed outputs and $Y = (y_1, y_2, \dots, y_n)$ be the sequence of hidden states. x_i can be any value within a set of possible outputs O and y_i can be any value within a set of possible state labels L . We define the transition probability $p(y_i | y_{i-1})$ of the current state given the previous state. The emission probability $p(x_i | y_i)$ is the probability of observing the current output given the state of the node. The joint probability distribution of a sequence of outputs \mathbf{x} and a sequence of hidden states \mathbf{y} can be defined as follows, denoting $p(y_1)$ as $p(y_1 | y_0)$ for simplicity:

$$p(\mathbf{x}, \mathbf{y}) = \prod_{i=1}^n p(y_i | y_{i-1}) p(x_i | y_i). \quad (3.3)$$

This is a generalization of the joint probability distribution defined in the naive Bayes classifier, and can be rewritten as follows:

$$p(\mathbf{x}, \mathbf{y}) = \exp \left[\sum_{i=1}^n \log(p(y_i | y_{i-1})) + \sum_{i=1}^n \log(p(x_i | y_i)) \right]. \quad (3.4)$$

If we by replace $\log(p(y_i | y_{i-1}))$ with a parameter $\beta_{y_i, y_{i-1}}$, $\log(p(x_i | y_i))$ with a parameter μ_{x_i, y_i} , and adjust by a normalization factor Z , we can rewrite this further as:

$$p(\mathbf{x}, \mathbf{y}) = \frac{1}{Z} \exp \left[\sum_{i=1}^n \beta_{y_i, y_{i-1}} + \sum_{i=1}^n \mu_{x_i, y_i} \right]. \quad (3.5)$$

These parameters can be indexed based on the set of labels by using indicator functions to determine the appropriate parameter:

$$p(\mathbf{x}, \mathbf{y}) = \frac{1}{Z} \exp \left[\sum_{i=1}^n \sum_{j, k \in L} \beta_{j, k} \mathbf{1}_{\{y_i=j\}} \mathbf{1}_{\{y_{i-1}=k\}} + \sum_{i=1}^n \sum_{o \in O} \sum_{l \in L} \mu_{o, l} \mathbf{1}_{\{x_i=o\}} \mathbf{1}_{\{y_i=l\}} \right]. \quad (3.6)$$

Finally, feature functions can be defined to simplify the notation used. Allow $f_{j, k}(y_i, y_{i-1}, x_i) = \mathbf{1}_{\{y_i=j\}} \mathbf{1}_{\{y_{i-1}=k\}}$ and $f_{o, l}(y_i, y_{i-1}, x_i) = \mathbf{1}_{\{x_i=o\}} \mathbf{1}_{\{y_i=l\}}$. Under this notation, each pair of possible labels (j, k) and each observation-label pair (o, l) has a feature function defined. By indexing these feature functions and their corresponding parameters $\beta_{j, k}$ and $\mu_{o, l}$ using q (with F total functions), we can write the joint probability as follows:

$$p(\mathbf{x}, \mathbf{y}) = \frac{1}{Z} \exp \left[\sum_{i=1}^n \sum_{q=1}^F \lambda_q f_q(y_i, y_{i-1}, x_i) \right]. \quad (3.7)$$

This notation will allow the differences between HMMs and CRFs to be highlighted.

3.2.4 CRF Model Description

As with HMMs, we define X as the sequence of hidden states and Y as the sequence of outputs. However, unlike HMMs, feature functions can be defined that can account for other output features. In the case of aspect term extraction, this means that the token features defined in the previous chapter can be used to train the model. [27]

Consider the joint probability distribution for HMMs. The conditional probability $p(\mathbf{y} \mid \mathbf{x})$, derived from the joint distribution, is:

$$p(\mathbf{y} \mid \mathbf{x}) = \frac{\exp \left[\sum_{i=1}^n \sum_{q=1}^F \lambda_q f_q(y_i, y_{i-1}, x_i) \right]}{\sum_{\mathbf{y}'} \exp \left[\sum_{i=1}^n \sum_{q=1}^F \lambda_q f_q(y'_i, y'_{i-1}, x_i) \right]}. \quad (3.8)$$

This is equivalent to a linear-chain Conditional Random Field with feature functions corresponding to each output. This is a specific sub-case of linear-chain CRFs; more generally, we can describe each output x_i as a vector of features. In our case, this means that rather than using only the word itself as a feature, we can use various features related to the word (such as prefixes/suffixes, part-of-speech tags, or whether capitalization is used). A feature function and corresponding parameter can be defined for any function of the current features, the current label, and the previous label. The general model is described below:

$$p(\mathbf{y} \mid \mathbf{x}) = \frac{\exp \left[\sum_{i=1}^n \sum_{q=1}^F \lambda_q f_q(y_i, y_{i-1}, \mathbf{x}_i) \right]}{Z(\mathbf{x})}, \quad (3.9)$$

where $Z(\mathbf{x}) = \sum_{\mathbf{y}'} \exp \left[\sum_{i=1}^n \sum_{q=1}^F \lambda_q f_q(y'_i, y'_{i-1}, \mathbf{x}_i) \right]$ is the normalization constant, computed by summing the feature functions multiplied by their weights over the possible

label combinations. The number of possible label combinations becomes very large, but it will be shown that this problem can be averted during training.

3.2.5 CRF Training

Training requires a set of training data $\{(\mathbf{x}^{(i)}, \mathbf{y}^{(i)})\}_{i=1}^N$ consisting of N 'documents' - in our case, sentences. Each sentence has n_i tokens. For sentence i , $\mathbf{y}^{(i)} = \{y_1^{(i)}, y_2^{(i)}, \dots, y_{n_i}^{(i)}\}$ is a sequence of IOB2 labels for a sentence and $\mathbf{x}^{(i)} = \{\mathbf{x}_1^{(i)}, \mathbf{x}_2^{(i)}, \dots, \mathbf{x}_{n_i}^{(i)}\}$ is a sequence of feature vectors, with one feature vector for each token in the sentence. The goal of training is to maximize the conditional log-likelihood for a set of parameters $\theta = \{\lambda_q\}_{q=1}^F$.

$$l(\theta) = \sum_{i=1}^N \log(p(\mathbf{y}^{(i)} | \mathbf{x}^{(i)})). \quad (3.10)$$

In addition, a technique called *regularization* is often used to smooth the parameters by making a penalty for overfitting:

$$l(\theta) = \sum_{i=1}^N \log(p(\mathbf{y}^{(i)} | \mathbf{x}^{(i)})) - \sum_{q=1}^F \frac{\lambda_q^2}{2\sigma^2}. \quad (3.11)$$

This assumes a Gaussian prior on the parameters θ , each with a mean of 0 and variance σ^2 . The gradient of $l(\theta)$ is:

$$\frac{\partial l}{\partial \lambda_q} = \sum_{i=1}^N \sum_{j=1}^{n_i} f_q(y_j^{(i)}, y_{j-1}^{(i)}, \mathbf{x}_j^{(i)}) - \sum_{i=1}^N \frac{\partial}{\partial \lambda_q} \log(Z(\mathbf{x}^{(i)})) - \frac{\lambda_q}{\sigma^2}. \quad (3.12)$$

where

$$\frac{\partial}{\partial \lambda_q} \log(Z(\mathbf{x}^{(i)})) = \sum_{i=1}^N \sum_{j=1}^{n_i} \sum_{y, y' \in L} f_q(y, y', \mathbf{x}_j^{(i)}) p(y, y' | \mathbf{x}^{(i)}). \quad (3.13)$$

The partial derivative with respect to λ_q can be interpreted as follows: the first component is the number of observed occurrences of the feature f_q , the second component is the expected number of occurrences of the feature f_q , and the third is the regularization adjustment. At the maximum likelihood solution, the expected and observed occurrences should be equal.

The maximum likelihood function $l(\theta)$ with regularization is strictly concave, and so a global optimum can be found [27]. This can be done with nonlinear optimization algorithms such as L-BFGS, stochastic gradient descent, and others. CRFsuite, a software implementation of conditional random fields, allows various optimization algorithms to be used for this purpose [22].

3.2.6 CRF Evaluation

An important consideration is the method with which ATE systems are evaluated. One key question is whether to apply these methods to distinct aspect terms or to each occurrence of an aspect term. If we evaluate based on distinct aspect terms, then we take the set of predicted distinct aspect terms and compare them to the list of actual distinct aspect terms. However, aspect terms with higher frequency are more valuable, given that our eventual goal is to determine polarity scores for a few most common terms/categories. A model that is able to accurately predict high-frequency aspect terms, but is less effective at predicting low-frequency aspect terms, is more valuable than a model that is better at predicting low-frequency terms than high-frequency terms.

On the other hand, evaluation based on instances of each aspect term can lead to overconfidence in models that can identify some of the most common terms with accuracy, but cannot accurately identify most other terms. Aspect terms with the highest frequencies in the dataset aren't always more important to accurately identify than aspect terms with lower frequencies. An individual aspect term may be more frequent than other aspect terms simply because it has few or no synonyms (for example, "Microsoft Office" has no synonyms, while "price" has several different words representing the same concept).

Thus, we evaluate the methods described in the previous sections with respect to both distinct aspect terms and instances of each aspect term. We use 70% of the data available in each domain for training and 30% for testing. As a review, three of the most common methods of evaluating models are precision, recall, and F-measure. Precision describes the fraction of predicted aspect terms that actually exist in the dataset. Recall is the fraction of true aspect terms that are predicted by the model. F-measure is the harmonic mean of precision and recall.

CRFsuite implements several algorithms to solve for the CRF parameters. Two of the most common optimization algorithms for solving CRFs are provided: L-BFGS and stochastic gradient descent. L-BFGS is a common quasi-Newton method that avoids storing a full approximated Hessian, making it useful for problems such as CRFs where there are often a large number of parameters to be found [17]. Stochastic gradient descent (SGD) is an extension of gradient descent that moves in the direction of a random data point at each iteration. In the CRFsuite implementation, SGD is performed with ℓ^2 regularization to prevent overfitting. Both of these algorithms have been shown to be successful when utilized to solve conditional random fields [28].

TABLE 3.1: The results for CRFs using distinct aspect terms.

Algorithm	Dataset	Precision	Recall	F-measure
L-BFGS	Restaurants	0.7003	0.5224	0.5984
SGD	Restaurants	0.6095	0.4187	0.4964
AP	Restaurants	0.6701	0.4004	0.5013
PA	Restaurants	0.6526	0.5346	0.5877
AROW	Restaurants	0.4399	0.5423	0.4859
L-BFGS	Laptops	0.5969	0.3793	0.4639
SGD	Laptops	0.3357	0.3522	0.3438
AP	Laptops	0.5682	0.2463	0.3436
PA	Laptops	0.5935	0.4064	0.4825
AROW	Laptops	0.4349	0.3867	0.4094

Three other algorithms are implemented in CRFsuite as well: Averaged perceptrons (AP), passive aggressive (PA), and Adaptive Regularization of Weight Vectors (AROW). Averaged perceptrons iterates over the training data, updating the feature weights of a perceptron whenever the model cannot make a correct prediction and updating the average feature weights. The final averaged feature weights are returned by the algorithm [7]. Passive-aggressive algorithms define a loss function on predicted instances, aggressively shifting the current parameter estimate when the current training instance has a positive value for the loss function and making no adjustment when the loss function is zero [9]. AROW is a variation of confidence-weighted learning, which maintains a Gaussian distribution to measure the confidence in each parameter estimate. It adjusts the model to prevent overly aggressive shifts that can occur when using passive-aggressive updates [10].

The results for distinct aspect terms for both Restaurant and Laptop datasets (using the 2014 format described in Figure 2.1) can be seen in Table 3.1. The best training algorithms for both datasets evaluated with distinct aspect terms were L-BFGS and PA. Overall, CRFs were more effective on the restaurants domain (with a best

TABLE 3.2: The results for CRFs using instances of aspect terms.

Algorithm	Domain	Precision	Recall	F-measure
L-BFGS	Restaurants	0.8491	0.7231	0.7810
SGD	Restaurants	0.8036	0.5629	0.6621
AP	Restaurants	0.8246	0.6127	0.7030
PA	Restaurants	0.8182	0.7574	0.7867
AROW	Restaurants	0.6664	0.7430	0.7026
L-BFGS	Laptops	0.8025	0.6119	0.6944
SGD	Laptops	0.4668	0.4268	0.4459
AP	Laptops	0.7460	0.3895	0.5118
PA	Laptops	0.7715	0.6436	0.7018
AROW	Laptops	0.6510	0.6312	0.6409

F-measure of 0.5984 when using L-BFGS) than on the laptop domain (with a best F-measure of 0.4825 when using PA).

The results for instances of aspect terms can be seen in Table 3.2. The best training algorithms for both datasets evaluated with aspect term instances were L-BFGS and PA, with F-measures of 0.7810 and 0.7867 respectively for the Restaurant domain and 0.6944 and 0.7018 respectively for the Laptop domain. Again, the CRF seems to be more effective on the Restaurant domain than the Laptop domain.

The model seems to have significantly higher precision than recall regardless of algorithm and across both distinct and instance-based evaluation methods. This suggests that the models may have difficulty identifying some aspect terms; however, the significant increase in both precision and recall when evaluating the instances of each aspect term suggests that much of this may come from failing to identify infrequent aspect terms.

To see how effective the model would be on data outside of the training domain, we attempt to train each model on one domain and evaluate its performance using testing data from the other domain. These results can be found in Table 3.3 and Table 3.4.

TABLE 3.3: The results for CRFs using distinct aspect terms across domains.

Algorithm	Train Domain	Test Domain	Precision	Recall	F-Measure
L-BFGS	Restaurant	Laptop	0.4354	0.1576	0.2315
SGD	Restaurant	Laptop	0.5272	0.0714	0.1258
AP	Restaurant	Laptop	0.3656	0.1675	0.2297
PA	Restaurant	Laptop	0.4084	0.1921	0.2613
AROW	Restaurant	Laptop	0.1635	0.1921	0.1767
L-BFGS	Laptop	Restaurant	0.6176	0.1280	0.2121
SGD	Laptop	Restaurant	0.3918	0.3089	0.3455
AP	Laptop	Restaurant	0.5350	0.1707	0.2589
PA	Laptop	Restaurant	0.5509	0.1870	0.2792
AROW	Laptop	Restaurant	0.3221	0.2134	0.2567

TABLE 3.4: The results for CRFs using instances of aspect terms across domains.

Algorithm	Train Domain	Test Domain	Precision	Recall	F-Measure
L-BFGS	Restaurant	Laptop	0.4900	0.1699	0.2523
SGD	Restaurant	Laptop	0.5795	0.0704	0.1256
AP	Restaurant	Laptop	0.4247	0.1754	0.2483
PA	Restaurant	Laptop	0.4935	0.2113	0.2959
AROW	Restaurant	Laptop	0.1909	0.1851	0.1879
L-BFGS	Laptop	Restaurant	0.8216	0.1792	0.2942
SGD	Laptop	Restaurant	0.5581	0.2824	0.3750
AP	Laptop	Restaurant	0.7289	0.1801	0.2888
PA	Laptop	Restaurant	0.7765	0.2516	0.3800
AROW	Laptop	Restaurant	0.5040	0.2308	0.3166

Overall, the quality of the results suffered significantly, suggesting that the model doesn't perform well on data outside of the domain of the training data. However, some of the algorithms used seem to suffer less reduction in quality than others. Using SGD on the Laptops dataset for training, the F-measures for distinct and instances (0.3438 and 0.4459, respectively) for the Laptop testing set are relatively close to their values when using the Restaurant testing set (0.3455 and 0.3750, respectively). Though the overall results were still poor, this suggests that some methods and training datasets may be more generalizable than others. Discovering than area that may worth exploring in the future.

3.3 Association Mining Method (Hu and Liu)

A method based on association mining to find frequent itemsets was defined in [14]. It is a rule-based method that builds a list of candidate itemsets consisting of nouns and noun phrases in each sentence, then prunes them to identify aspect terms. This is based on the notion that reviewers tend to use similar words when describing aspects of a review topic, and so frequently-occurring sets of words are more likely to be aspect terms.

3.3.1 Association Mining Method Description

First, a set of initial candidate itemsets are generated. A list of nouns and noun phrases N , ordered by their placement within the sentence, are extracted from each sentence i as initial itemsets. Pairs and triples of these nouns and noun phrases within each sentence are also considered candidate terms. This is only done for adjacent nouns and noun phrases. More specifically, the extracted pairs and triples can be described as

$$Pairs = \{N_i \cup N_{i+1} : i \in \{1, 2, \dots, |N| - 1\}\}$$

$$Triples = \{N_i \cup N_{i+1} \cup N_{i+2} : i \in \{1, 2, \dots, |N| - 2\}\}$$

At this point, we have an initial set of candidate terms. We reduce the set of candidate terms down to a set of “frequent” itemsets, as defined by a minimum support level m . This can also be based off a specified percentage of the dataset. All other candidate terms are eliminated.

Two additional pruning measures are taken to reduce the set of candidate terms. An adjusted frequency measure called “p-support” is found that only counts a candidate term in a sentence if it is not a subset of another candidate term within the sentence. For example, if a sentence contained the phrase “ham sandwich” and both “ham” and “ham sandwich” were candidate terms, then this sentence would not count towards the support of “ham” , since it’s a subset of another candidate term “ham sandwich” that exists within the sentence. If the p-support of a candidate term is low, and it appears as part of a larger candidate term, the candidate term is likely a component of the larger term. We define a minimum p-support threshold p - if the p-support of a candidate term is less than p and the candidate term is a subset of some other term, we remove it from the set of candidate terms.

Another pruning measure attempts to correct for issues that can arise from using frequent itemsets. When the initial set of candidate terms is created, pairs and triples of nouns and noun phrases are considered candidate terms. However, these words may be relatively far apart within a sentence, suggesting that they might not be part of the same aspect term. For a term within a given sentence, we find the maximum distance between any two adjacent words in the term. Their distance is measured by how many tokens apart they are in the sentence. If this value exceeds a token distance threshold w , then we consider the term non-compact within the sentence. If a term is found to be non-compact in a greater number of sentences than a maximum non-compact frequency threshold c , then the term is discarded.

Algorithm 1 Association Mining Method (Hu and Liu)

Require: List of sentences S , minimum support threshold m , a token distance threshold w , a maximum non-compact frequency threshold c , and a minimum p-support threshold p .

```

1:  $T = \{\}$ 
2: for sentence  $\in S$  do
3:    $N = \text{getNounsAndNounPhrases}(\text{sentence})$ 
4:    $T = T \cup \{\text{term} \in N, \text{getAdjacentPairs}(N), \text{getAdjacentTriples}(N)\}$ 
5: end for
6: Support = Dictionary() // Default key value is 0
7: for sentence  $\in S$  do
8:   for term  $\in (T \cap \text{sentence})$  do
9:     Support[term] = Support[term] + 1
10:  end for
11: end for
12:  $T.\text{remove}(\{\text{term} : (\text{term} \in T) \text{ and } (\text{Support}[\text{term}] < m)\})$ 
13: P-Support = Dictionary() // Default key value is 0
14: Non-Compact = Dictionary() // Default key value is 0
15: for sentence  $\in S$  do
16:   for term  $\in (T \cap \text{sentence})$  do
17:     if maxTokenDistance(term, sentence)  $> w$  then
18:       Non-Compact[term] = Non-Compact[term] + 1
19:     end if
20:     if (term  $\not\subset$  term2)  $\forall$  term2  $\in ((T - \{\text{term}\}) \cap \text{sentences})$  then
21:       P-Support[term] = P-Support[term] + 1
22:     end if
23:   end for
24: end for
25: for term  $\in T$  do
26:   if (Non-Compact[term]  $> c$ ) then
27:      $T.\text{remove}(\text{term})$ 
28:   end if
29:   if ((P-Support[term]  $< p$ ) then
30:     for term2  $\in T - \{\text{term}\}$  do
31:       if term.contains(term2) then
32:          $T.\text{remove}(\text{term})$ 
33:       end if
34:     end for
35:   end if
36: end for
37: return  $T$ 

```

3.3.2 Association Mining Method Evaluation

An important consideration in aspect term extraction is the idea that aspect terms are likely to be nouns and noun phrases. In the case of one-word aspect terms that are nouns, identification is as simple as finding nouns from each sentence using a part-of-speech tagger, then using other methods to filter out nouns that aren't actually aspect terms. For the case of multi-word aspect terms, noun phrases must be identified. Any unsupervised method for aspect identification must somehow identify these noun phrases without the benefit of training data. The general problem of identifying grammatical structures such as noun phrases is called shallow parsing [3]. Three different methods were explored for identifying noun phrases. We attempted to use NLTK's "Regex" (regular expression) feature, which finds specific patterns in text using a pre-defined search pattern [6]. However, noun phrases take many possible forms, and defining all the possible search patterns that noun phrases may exist in is unfeasible. We also examined bigram and trigram classifiers, trained on a portion of Treebank data available in NLTK [20]. Finally, we examined the default named-entity chunker within NLTK.

In the testing of the Association Mining algorithm, it became clear that noun chunking was a significant issue that hindered the performance of the algorithm as a whole. After tuning the input parameters for the Restaurant domain dataset, the best model using the named-entity chunker had a precision of 0.3777, recall of 0.2480, and F-measure of 0.2994. This is with a minimum support threshold m of 6, a minimum p-support threshold p of 2, a max token distance threshold w of 2, and a maximum non-compact frequency threshold of 1. We examined the full list of candidate terms before pruning (consisting of all nouns and noun phrases in the sentence, as well as all adjacent

pairs and triples of nouns and noun phrases) and found that only 61.18% were detected - this provides an upper bound on the recall of the model. In the future, examining effective ways to identify noun phrases is an important step in improving unsupervised methods, particularly those based on frequent itemsets.

Chapter 4

Aspect-Based Sentiment Analysis

4.1 Problem Description

This chapter is focused on estimating the sentiment of the aspect terms in a sentence, assuming the aspect terms are known. We also examine the case where aspect categories are provided for each sentence rather than individual aspect terms.

Given a set of reviews with aspect terms identified, we would like to accurately estimate the sentiment of each occurrence of an aspect term in a sentence. With one aspect per sentence, an assumption can be made that polarity within the sentence is associated with the polarity of the aspect. When multiple aspect terms are present in one sentence, words associated with one aspect term may incorrectly be associated with another aspect term, causing the polarities of each aspect term within the sentence to affect each other.

An issue with using aspect terms individually is that oftentimes, multiple aspect terms will refer to the same or similar aspects. For example, "price" and "cost" refer

to the same aspect, yet are considered separate aspect terms. This suggests that a way to categorize aspect terms is desirable when designing a system to accurately rate important aspects of a review’s subject. As such, we will focus on accurately identifying the sentiment of instances of aspect terms, rather than the problem of aggregating these terms to provide a more accurate view of a more general “aspect category”.

A secondary formulation of the problem can be given for aspect categories (pre-defined categories that collectively contain the most important or commonly-discussed aspects of a review’s subject). Given a set of reviews with sentence-level aspect categories, we would like to accurately estimate the sentiment of each occurrence of an aspect category in a sentence. There are multiple benefits to using aspect categories rather than specific aspect terms. Typically there will be a much smaller number of aspect categories than aspect terms, and these categories will be present in a larger number of sentences than individual aspect terms. This means a smaller amount of data is needed to have enough instances of an aspect category to provide an accurate rating. However, identifying these aspect categories in the first place can be difficult, and requires a predefined list of categories for each domain. As such, the results provided here are predicated on the availability of a method to identify these aspect categories.

4.2 VADER-based Method

VADER, or the Valence-Aware Dictionary for sEntiment Reasoning, is a rule-based model for performing sentiment analysis on a per-sentence basis. The system was trained on online media text, some of which included movie and product reviews. VADER utilizes a sentiment lexicon constructed with the purpose of being generalizable to multiple

domains. This makes VADER particularly suitable for analyzing online review data. In addition, by classifying on a per-sentence basis and performing unsupervised, VADER can easily be applied and tested on newly-seen data and data across domains.

Their sentiment lexicon was based on several existing sentiment lexicons, as well as common emoticons and acronyms. It includes valence scores (between -4 and 4) that contain information about sentiment intensity (how strongly a word expresses a sentiment) in addition to sentiment polarity.

Given a sentence, VADER calculates a valence score to measure the sentiment intensity and polarity. Five major heuristics are used to determine the valence score of a given sentence:

1. Some types of punctuation, specifically exclamation points, increases the magnitude of the valence score. For example, “The keyboard is great.” is rated with a lower magnitude than “The keyboard is great!”
2. Full-word capitalization, especially when other nearby words aren’t fully capitalized, increases the magnitude of the valence score. For example, “The keyboard is great!” is rated with a lower magnitude than “The keyboard is GREAT!”
3. A set of adverbs called ‘degree modifiers’ is used to increase or decrease the magnitude of the valence score, depending on the word. For example, “The keyboard is great.” is rated with a lower magnitude than “The keyboard is very great.” and with a higher magnitude than “The keyboard is kinda great.”
4. The conjunction “but” signals a shift in sentiment polarity. The sentiment of the portion of the sentence after “but” is considered to be the dominant sentiment, and

contributes a greater amount (two-thirds) to the valence score than the portion of the sentence before “but” (one-third).

5. The trigram before an occurrence of a lexical feature (as determined by the sentiment lexicon) is examined to determine whether a negation is used to express the opposite polarity. For example, “The keyboard is not great” would be given a negative valence score, since “not” is a negation that flips the polarity of “great”.

VADER returns a set of four scores: one each for “positive”, “negative”, and “neutral” (which together sum to 1.0), as well as a “compound” score (ranging from -1.0 to 1.0) reflecting the intensity of the polarity within the sentence. Negative scores are associated with negative polarity within a sentence, and positive scores are associated with positive sentence polarity. Larger magnitudes of the “compound” score are associated with higher intensities.

In order to compare these scores with the available data, for each sentence we return a single label (“positive”, “negative”, or “neutral”) depending on the scores returned by VADER. If a sentence’s “neutral” score is 1.0, we return the label “neutral”. If a sentence’s “negative” score is greater than its “positive” score (or if the “compound” score is less than 0), we return “negative”. Otherwise, we return “positive”. We do not attempt to classify “conflict” values.

4.2.1 Evaluation

We keep track of the predicted and true label values for each occurrence of an aspect term (and additionally for each occurrence of an aspect category, in the case of the restaurant domain dataset). Accuracy is the primary measurement we use to evaluate

TABLE 4.1: The results using VADER on aspect terms in the Laptop domain.

Accuracy:		0.5855	
Label	Precision	Recall	F-Measure
Positive	0.8140	0.6543	0.7254
Negative	0.3362	0.2916	0.3123
Neutral	0.4535	0.7104	0.5536

TABLE 4.2: The results using VADER on aspect terms in the Restaurant domain.

Accuracy:		0.6501	
Label	Precision	Recall	F-Measure
Positive	0.8235	0.7558	0.7882
Negative	0.4028	0.3287	0.3620
Neutral	0.3730	0.6423	0.4719

TABLE 4.3: The results using VADER on aspect categories in the Restaurant domain.

Accuracy:		0.6535	
Label	Precision	Recall	F-Measure
Positive	0.7918	0.7974	0.7946
Negative	0.5226	0.2942	0.3765
Neutral	0.3674	0.6570	0.4713

our model; and Precision, recall, and F-Measure are also calculated with respect to the labels “positive”, “negative”, and “neutral”. The term-based results can be found in Table 4.1 for the laptop domain dataset and in Table 4.2 for the Restaurant domain. The accuracy is reasonably high for an unsupervised model, though it is somewhat lower for the Laptop domain (0.5855) versus the restaurant domain dataset (0.6501). Evaluating based on aspect categories for the restaurant domain dataset provides similar results, without significant variation in any of the evaluation measures. These results can be found in Table 4.3.

4.2.2 Ratings-Based Evaluation

Given the number of positive (p), negative (n), neutral/objective (o), and conflict (c) labels for a given aspect term or aspect category, a rating (r) from 1 to 5 can be

TABLE 4.4: True and predicted ratings for each category in the Restaurant domain.

Category	True Rating	Predicted Rating	Rating Error
food	4.15	4.42	-0.27
ambience	3.81	4.42	-0.61
price	3.44	4.48	-1.04
anecdotes/miscellaneous	3.92	4.15	-0.23
service	3.38	4.05	-0.67

determined as follows:

$$r = 4 \left[\frac{p + 0.5c}{p + n + c} \right] + 1. \quad (4.1)$$

This model assumes that “conflict” labels are associated with an equal split between positive and negative sentiment, and assumes that positive occurrences should be weighted the same as negative occurrences. This assumption is based on the idea that a review with n out of N stars has a fraction of positive to negative sentiment of $\frac{n-1}{N-1}$. However, this may not be true in practice. Given a dataset with quantitative review scores in addition to annotated aspect terms or categories, more accurate proportions of positive to negative sentiment may be developed. Using these proportions, weights can be used to more heavily skew an occurrence of a particular sentiment label versus other sentiment labels.

We use the restaurant domain dataset’s aspect categories to calculate ratings, since there are significantly more occurrences of each aspect category than any one aspect term. The ratings based on the adjusted VADER model’s predictions and based on the true sentiment labels can be found in Table 4.4. Overall, the predicted ratings tended to overestimate the true rating by an average of 0.564; this suggests that VADER is somewhat skewed towards positive ratings, at least on our available dataset.

Chapter 5

Conclusion

In this thesis, we explored some of the key tasks in the development of an aspect-based review system. We outlined the considerations required for developing an annotated dataset for the purpose of training models for aspect identification and aspect-based sentiment analysis. For the task of aspect identification, two algorithms were described and tested: a supervised sequential learning model called a conditional random field and an unsupervised association mining algorithm. The results for conditional random fields suggest that they are an effective classifier for identifying aspect terms, particularly when the parameters are learned using L-BFGS or a passive-aggressive algorithm. The results for the association mining algorithm were relatively poor due to issues with identifying noun phrases, but illuminated a future area for further exploration: accurately identifying noun phrases. For the task of aspect-based sentiment analysis, we describe a modified version of VADER, a rule-based sentiment intensity analyzer, to estimate the sentiment of aspect terms and aspect categories. [15]. The results for this model were

A significant area of future exploration is aspect aggregation - identifying aspect terms that are synonyms of each other (for example, "price" and "cost") and aspect terms that are a part of an overarching category (for example, "water" and "wine" might be part of an overarching category called "beverages"). This can be done with predefined categories, which can allow for a supervised approach to the clustering problem. Review-level and sentence-level training data is difficult to generate for a large number of domains, but having a small number of predefined categories to capture the most common aspect terms for each domain is much more feasible. Unsupervised clustering methods may also be explored, given a fixed number of clusters. In this case, clusters can be identified by their most frequent aspects.

Bibliography

- [1] SIGLEX (ACL Special Interest Group). <http://alt.qcri.org/siglex/>. Accessed: 2017-04-21.
- [2] WNSTATS(7WN) manual page. <http://wordnet.princeton.edu/wordnet/man/wnstats.7WN.html>. Accessed: 2017-04-27.
- [3] S. P. Abney. Parsing by chunks. In *Principle-based parsing*, pages 257–278. Springer, 1991.
- [4] I. Androutsopoulos, D. Galanis, S. Manandhar, H. Papageorgiou, J. Pavlopoulos, and M. Pontiki. SemEval-2015 Task 12: Aspect Based Sentiment Analysis < SemEval-2015 Task 12. <http://alt.qcri.org/semEval2015/task12/>. Accessed: 2017-04-15.
- [5] I. Androutsopoulos, D. Galanis, S. Manandhar, H. Papageorgiou, J. Pavlopoulos, and M. Pontiki. Task Description: Aspect Based Sentiment Analysis (ABSA) < semeval-2014 task 4. <http://alt.qcri.org/semEval2014/task4/>. Accessed: 2017-04-15.
- [6] S. Bird. NLTK: the natural language toolkit. In *Proceedings of the COLING/ACL on Interactive presentation sessions*, pages 69–72. Association for Computational Linguistics, 2006.
- [7] M. Collins. Discriminative training methods for hidden markov models: Theory and experiments with perceptron algorithms. In *Proceedings of the ACL-02 Conference on Empirical Methods in Natural Language Processing-Volume 10*, pages 1–8. Association for Computational Linguistics, 2002.
- [8] N. R. Council, A. L. P. A. Committee, et al. *Language and Machines: Computers in Translation and Linguistics; A Report*. National Academy of Sciences, National Research Council, 1966.
- [9] K. Crammer, O. Dekel, J. Keshet, S. Shalev-Shwartz, and Y. Singer. Online passive-aggressive algorithms. *Journal of Machine Learning Research*, 7(Mar):551–585, 2006.
- [10] K. Crammer, A. Kulesza, and M. Dredze. Adaptive regularization of weight vectors. In *Advances in Neural Information Processing Systems*, pages 414–422, 2009.

-
- [11] A. P. Dempster, N. M. Laird, and D. B. Rubin. Maximum likelihood from incomplete data via the EM algorithm. *Journal of the Royal Statistical Society. Series B (methodological)*, pages 1–38, 1977.
- [12] J. J. Godfrey, E. C. Holliman, and J. McDaniel. SWITCHBOARD: Telephone speech corpus for research and development. In *Acoustics, Speech, and Signal Processing, 1992. ICASSP-92., 1992 IEEE International Conference on Acoustics, Speech, and Signal Processing*, volume 1, pages 517–520. IEEE, 1992.
- [13] L. Hirschman and R. Gaizauskas. Natural language question answering: the view from here. *Natural Language Engineering*, 7(04):275–300, 2001.
- [14] M. Hu and B. Liu. Mining and summarizing customer reviews. In *Proceedings of the tenth ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pages 168–177. ACM, 2004.
- [15] C. J. Hutto and E. Gilbert. Vader: A parsimonious rule-based model for sentiment analysis of social media text. In *Eighth International AAAI Conference on Weblogs and Social Media*, 2014.
- [16] K. S. Jones. Natural language processing: a historical review. In *Current Issues in Computational Linguistics: In Honour of Don Walker*, pages 3–16. Springer, 1994.
- [17] D. C. Liu and J. Nocedal. On the limited memory BFGS method for large scale optimization. *Mathematical Programming*, 45(1):503–528, 1989.
- [18] C. D. Manning, M. Surdeanu, J. Bauer, J. R. Finkel, S. Bethard, and D. McClosky. The Stanford CoreNLP Natural Language Processing Toolkit. In *ACL System Demonstrations*, pages 55–60, 2014.
- [19] T. N. Mansuy and R. J. Hilderman. Evaluating WordNet Features in Text Classification Models. In *FLAIRS Conference*, pages 568–573, 2006.
- [20] M. P. Marcus, M. A. Marcinkiewicz, and B. Santorini. Building a large annotated corpus of English: The Penn Treebank. *Computational Linguistics*, 19(2):313–330, 1993.
- [21] G. A. Miller. WordNet: a lexical database for English. *Communications of the ACM*, 38(11):39–41, 1995.
- [22] N. Okazaki. CRFsuite: a fast implementation of conditional random fields (CRFs). 2007.
- [23] I. Pavlopoulos. Aspect based sentiment analysis. *Athens University of Economics and Business*, 2014.
- [24] M. Pontiki, D. Galanis, H. Papageorgiou, S. Manandhar, and I. Androutsopoulos. Task 5: Aspect-Based Sentiment Analysis < SemEval-2016 Task 5. <http://alt.qcri.org/semeval2016/task5/>. Accessed: 2017-04-15.

-
- [25] M. F. Porter. An algorithm for suffix stripping. *Program*, 14(3):130–137, 1980.
- [26] F. Sha and F. Pereira. Shallow parsing with conditional random fields. In *Proceedings of the 2003 Conference of the North American Chapter of the Association for Computational Linguistics on Human Language Technology-Volume 1*, pages 134–141. Association for Computational Linguistics, 2003.
- [27] C. Sutton and A. McCallum. An introduction to conditional random fields. *arXiv preprint arXiv:1011.4088*, 2010.
- [28] S. Vishwanathan, N. N. Schraudolph, M. W. Schmidt, and K. P. Murphy. Accelerated training of conditional random fields with stochastic gradient methods. In *Proceedings of the 23rd International Conference on Machine Learning*, pages 969–976. ACM, 2006.

Appendix - Data Processing and Test Functions

```
1 import time
2 from collections import defaultdict
3 import xml.etree.ElementTree as ET
4 import libraries.structure as st
5 from libraries.structure import Corpus
6 import aspect_identification as ai
7 import sentiment_analysis as sa
8 from stanford_corenlp_python import jsonrpc
9
10
11 def sentimentAnalysisTest(data):
12     instances = data.corpus
13     trueTermPolsBySent = sa.getTermPolarities(instances)
14     trueCatPolsBySent = sa.getCategoryPolarities(instances)
15
16     predictedTermPolsBySent = sa.vaderTermPolarities(instances)
17     print "Evaluate By Terms:"
18     sa.evaluatePolarities(trueTermPolsBySent, predictedTermPolsBySent)
19
20     if len([i for j in trueCatPolsBySent for i in j]) > 0:
21         predictedCatPolsBySent = sa.vaderCategoryPolarities(instances)
22         print "Evaluate By Categories:"
23         sa.evaluatePolarities(trueCatPolsBySent, predictedCatPolsBySent)
24         print "True Ratings:"
25         print sa.computeRatings(st.fd2([i for j in trueCatPolsBySent for
26     ↪ i in j]))
27         print "Predicted Ratings:"
28         print sa.computeRatings(st.fd2([i for j in
29     ↪ predictedCatPolsBySent for i in j]))
30
31 def aspectIdentificationTest(dataR, dataL, HL = True, CRF = True):
32     # Split into train/test data
33     trainR, testR = dataR.split(threshold=0.7)
34     trainL, testL = dataL.split(threshold=0.7)
35     train = trainR
```

```

34     test = testR
35
36     testFD = st.fd([" ".join(a.tokenized_term) for i in test for a in
    ↪ i.aspect_terms])
37     testBySent = [i.adjustFormat() for i in test]
38
39     numMethods = 0
40     if HL == True:
41         # H&L settings
42         minSupports = [0]
43         minPsupports = [0]
44         maxWordDist = [10.0]
45         maxNonCompact = [10.0]
46         params = [(i,j,k,l) for i in minSupports for j in minPsupports
    ↪ for k in maxWordDist for l in maxNonCompact]
47         numMethods += len(params)
48
49     if CRF == True:
50         # CRF settings
51         algs = ['lbfgs', 'l2sgd', 'ap', 'pa', 'arow']
52         numMethods += len(algs)
53
54     predictedFDs = range(numMethods)
55     predictedTermsBySent = range(numMethods)
56     methodNames = []
57     count = 0
58
59     if HL == True:
60         for p in range(len(params)):
61             # Run Association Mining (Hu & Liu) algorithm
62             i = params[p][0]
63             j = params[p][1]
64             k = params[p][2]
65             l = params[p][3]
66             predictedFDs[count], predictedTermsBySent[count] =
    ↪ ai.HuLiu(dataL.corpus, minSupport = i, minPsupport = j,
    ↪ maxWordDist = k, maxNonCompact = l)
67             methodNames.append("H&L:
    ↪ (minS="+str(i)+" ,minPS="+str(j)+" ,maxWD="+str(k)+" ,maxNC="+str(l)+" ")
68             count += 1
69
70     if CRF == True:
71         for k in algs:
72             # Run Conditional Random Field algorithm
73             crfLabels = ai.crf(train, test, k)
74             predictedFDs[count], predictedTermsBySent[count] =
    ↪ ai.IOB2toAspectTerms(crfLabels, test)
75

```

```

76         methodNames.append(k)
77         count += 1
78     # Evaluate methods
79     ai.evaluateAspectTerms(testFD, testBySent, predictedFDs,
80         ↪ predictedTermsBySent, methodNames, True)
81
82 def process semeval_2015():
83     # the train set is composed by train and trial data set
84     corpora = dict()
85     corpora['laptop'] = dict()
86     train_filename =
87     ↪ 'datasets/ABSA-SemEval2015/ABSA-15_Restaurants_Train_Final.xml'
88     trial_filename =
89     ↪ 'datasets/ABSA-SemEval2015/absa-2015_restaurants_trial.xml'
90
91     reviews = ET.parse(train_filename).getroot().findall('Review') + \
92         ET.parse(trial_filename).getroot().findall('Review')
93
94     sentences = []
95     for r in reviews:
96         sentences += r.find('sentences').getchildren()
97
98     # TODO: parser is not loading aspect words and opiniooss
99     corpus = Corpus(sentences)
100     corpus.size()
101
102 def process semeval_2014(type = "R"):
103     # the train set is composed by train and trial dataset
104     # corpora = dict()
105     # corpora['data'] = dict()
106     if type == "R":
107         train_filename =
108         ↪ 'datasets/ABSA-SemEval2014/Restaurants_Train_v2.xml'
109         trial_filename =
110         ↪ 'datasets/ABSA-SemEval2014/restaurants-trial.xml'
111
112     elif type == "L":
113         train_filename = 'datasets/ABSA-SemEval2014/Laptop_Train_v2.xml'
114         trial_filename = 'datasets/ABSA-SemEval2014/laptops-trial.xml'
115     corpus =
116     ↪ Corpus(ET.parse(train_filename).getroot().findall('sentence') +
117         ↪ ET.parse(trial_filename).getroot().findall('sentence'))
118
119     # corpora['data']['trainset'] = dict()
120     # corpora['data']['trainset']['corpus'] = corpus
121     return corpus

```



```
117
118 def main():
119     # TODO: start corenlp server "python corenlp.py"
120
121     # interface for Stanford-Core-NLP server
122     start = time.time()
123     server = jsonrpc.ServerProxy(jsonrpc.JsonRpc20(),
124
125                                 ↪ jsonrpc.TransportTcpIp(addr=("127.0.0.1",
126                                                         8080)))
127
128     #result = loads(server.parse("Hello world. It is so beautiful"))
129     #print "Result", result
130
131     # Load corpus
132     dataR = process semeval_2014("R")
133     dataL = process semeval_2014("L")
134
135     print 'The restaurant corpus has %d sentences, %d aspect term
136     ↪ occurrences, and %d distinct aspect terms.' % (dataR.size,
137     ↪ sum(dataR.aspect_terms_fd[a] for a in dataR.aspect_terms_fd),
138     ↪ len(dataR.top_aspect_terms))
139
140     print 'The laptop corpus has %d sentences, %d aspect term
141     ↪ occurrences, and %d distinct aspect terms.' % (dataL.size,
142     ↪ sum(dataL.aspect_terms_fd[a] for a in dataL.aspect_terms_fd),
143     ↪ len(dataL.top_aspect_terms))
144
145     end = time.time()
146     print "Load Corpus: " + str(end - start) + " seconds"
147     start = end
148
149     aspectIdentificationTest(dataR, dataL, HL=True, CRF=False)
150     sentimentAnalysisTest(dataR)
151
152 if __name__ == '__main__':
153     main()
```

Appendix - Class Definitions

```
1 import xml.etree.ElementTree as ET, getopt, logging, sys, random, re,
   ↪ copy
2 from xml.sax.saxutils import escape
3 import nltk
4 from nltk.tokenize import WordPunctTokenizer
5 from nltk.tokenize import TreebankWordTokenizer as Tokenizer
6 from nltk.stem.porter import PorterStemmer as Stemmer
7
8 import string
9 from collections import defaultdict
10
11 def fd(counts):
12     '''Given a list of occurrences (e.g., [1,1,1,2]), return a
   ↪ dictionary of frequencies (e.g., {1:3, 2:1}.)'''
13     d = defaultdict(lambda:0)
14     for i in counts: d[i] = d[i] + 1 if i in d else 1
15     return d
16
17 freq_rank = lambda d: sorted(d, key=d.get, reverse=True)
18 '''Given a map, return ranked the keys based on their values.'''
19
20 def fd2(counts):
21     '''Given a list of 2-uplets (e.g., [(a,pos), (a,pos), (a,neg),
   ↪ ...]), form a dict of frequencies of specific items (e.g.,
   ↪ {a:{pos:2, neg:1}, ...}).'''
22     d = {}
23     for i in counts:
24         # If the first element of the 2-uplet is not in the map, add it.
25         if i[0] in d:
26             if i[1] in d[i[0]]:
27                 d[i[0]][i[1]] += 1
28             else:
29                 d[i[0]][i[1]] = 1
30         else:
31             d[i[0]] = defaultdict(lambda: 0)
32             d[i[0]][i[1]] += 1
33     return d
```

```

34
35 def validate(filename):
36     '''Validate an XML file, w.r.t. the format given in the 4th task of
37     ↪ **SemEval '14**.'''
38     elements = ET.parse(filename).getroot().findall('sentence')
39     aspects = []
40     for e in elements:
41         for eterms in e.findall('aspectTerms'):
42             if eterms is not None:
43                 for a in eterms.findall('aspectTerm'):
44                     aspects.append(Aspect('', '').createEl(a).term)
45     return elements, aspects
46
47 fix = lambda text: escape(text.encode('utf8')).replace('\\"', '&quot;')
48 '''Simple fix for writing out text.'''
49
50 # Dice coefficient
51 def dice(t1, t2, stopwords=[]):
52     tokenize = lambda t: set([w for w in t.split() if (w not in
53     ↪ stopwords)])
54     t1, t2 = tokenize(t1), tokenize(t2)
55     return 2. * len(t1.intersection(t2)) / (len(t1) + len(t2))
56
57 # Find the index of the nth occurrence of a word within a tokenized text
58 def findNthOccurrence(tokenized_text, word, n):
59     if n < 1:
60         print "Error: n must be an integer > 1"
61         exit()
62     k = 0 # How many occurrences we've seen so far
63     for index in range(len(tokenized_text)):
64         if word in tokenized_text[index]:
65             k = k + 1
66             if k == n:
67                 return index
68     print "Error: Could not find nth occurrence"
69     return -1
70
71 def generate(sentences):
72     features = [[token.toDict() for token in s.tokens] for s in
73     ↪ sentences]
74     labels = [[token.actualIOB2 for token in s.tokens] for s in
75     ↪ sentences]
76     return features, labels
77
78 class Category:

```

```

76     '''Category objects contain the term and polarity (i.e., pos, neg,
77     ↪ neu, conflict) of the category (e.g., food, price, etc.) of a
78     ↪ sentence.'''
79
80     def __init__(self, term='', polarity=''):
81         self.term = term
82         self.polarity = polarity
83
84     def createEl(self, element):
85         self.term = element.attrib['category']
86         self.polarity = element.attrib['polarity']
87         return self
88
89     def update(self, term='', polarity=''):
90         self.term = term
91         self.polarity = polarity
92
93 class Token:
94     '''Token objects contain information about an individual token -
95     ↪ usually a word or punctuation. '''
96
97     def __init__(self, text='', index=-1):
98         self.text = text # The text of the
99         ↪ token
100        self.index = index # Index of the token
101        ↪ in the tokenized sentence
102        self.isBOS = not index # isBOS (Beginning
103        ↪ of sentence): True if index = 0, False otherwise
104        self.lower_text = text.lower() # The lowercase text
105        ↪ of the token
106        self.isTitle = text.istitle() # True if token is
107        ↪ "titlecased" (first letter is uppercase and other letters
108        ↪ are lowercase)
109        self.isPunct = text in string.punctuation # True if the token
110        ↪ is punctuation rather than a word
111        self.isDigit = text.isdigit() # True if the token
112        ↪ is a digit rather than a word
113        self.stem = Stemmer().stem(text) # Word stem of the
114        ↪ token (Ex: the stem of "running" is "run")
115        self.actualIOB2 = "O" # "O" if token is
116        ↪ outside, "I" if token is inside, "B" if token is the
117        ↪ beginning of an aspect term
118        self.polarity = "" # Positive ("pos"),
119        ↪ negative ("neg"), or neutral ("neu")
120        self.POS = "" # Part of speech of
121        ↪ the token
122        self.POS2 = "" # First 2 characters
123        ↪ of the POS tag

```

```

107
108     def toDict(self):
109         features = dict(self.__dict__)
110         features.pop('actualIOB2')
111         return features
112
113     def setIndex(self, index):
114         self.index = index
115
116     def setPrev(self, prev):
117         self.prev_text = prev.text
118         self.prev_lower_text = prev.lower_text
119         self.prev_POS = prev.POS
120         self.prev_POS2 = prev.POS2
121         self.prev_stem = prev.stem
122
123     def setNext(self, next):
124         self.next_text = next.text
125         self.next_lower_text = next.lower_text
126         self.next_POS = next.POS
127         self.next_POS2 = next.POS2
128         self.next_stem = next.stem
129
130     def setActualIOB2(self, IOB2):
131         self.actualIOB2 = IOB2
132
133     def setPredictedIOB2(self, IOB2):
134         self.predictedIOB2 = IOB2
135
136     def setPOS(self, POS):
137         self.POS = POS
138         self.POS2 = POS[:2]
139
140     def setPolarity(self, polarity):
141         self.polarity = polarity
142
143
144 class Aspect:
145     ''' Aspect objects contain information about each aspect term. '''
146
147     def __init__(self, term='', id='', tokens=''):
148         self.term = term           # The text of the aspect term
149         self.id = id              # The sentence id
150         self.offsets = ''         # The offsets within the sentence
151         ↪ {'from':startIndex, 'to':endIndex}
152         self.polarity = ''        # The polarity (pos, neg, neu,
153         ↪ conflict)

```

```

152     self.lower_term = ''           # The lowercase text of the aspect
    ↪ term
153     self.tokens = ''              # An ordered list of Tokens
    ↪ representing the sentence
154     self.tokenized_term = ''      # An ordered list of Strings
    ↪ representing the sentence
155     self.termSize = ''            # Number of elements in
    ↪ tokenized_term
156     self.headIndex = ''           # The index of the term's first
    ↪ token
157     self.endIndex = ''            # The index after the term's last
    ↪ token
158     if tokens != '':
159         self.createFromTokens(tokens)
160     elif len(term) > 0:
161         self.lower_term = self.term.lower()
162         self.tokenized_term = Tokenizer().tokenize(self.term)
163         self.lower_tokenized_term = [t.lower() for t in
    ↪ self.tokenized_term]
164         self.termSize = len(self.tokenized_term)
165
166     def createFromTokens(self, tokens):
167         ''' Create an Aspect from tokens (used after initial file
    ↪ processing) '''
168         self.tokens = tokens
169         self.tokenized_term = [t.text for t in tokens]
170         self.lower_tokenized_term = [t.lower for t in
    ↪ self.tokenized_term]
171         self.termSize = len(tokens)
172         self.headIndex = tokens[0].index
173         self.endIndex = self.headIndex + self.termSize
174
175     def createEl(self, element):
176         ''' Create an Aspect from an XML element (used when reading from
    ↪ file)
    '''
177
178         self.term = element.attrib['term']
179         self.lower_term = self.term.lower()
180         self.polarity = element.attrib['polarity']
181         self.offsets = {'from': str(element.attrib['from']), 'to':
    ↪ str(element.attrib['to'])}
182         self.lower_term = self.term.lower()
183         self.tokenized_term = Tokenizer().tokenize(self.term)
184         self.lower_tokenized_term = [t.lower() for t in
    ↪ self.tokenized_term]
185         self.termSize = len(self.tokenized_term)
186         return self
187

```

```
188     def compareWithinSentence(self, otherAspect):
189         ''' Comparison based on same sentence - only returns true if the
           ↳ aspects are in the same position within the sentence (ex:
           ↳ the first occurrence of the aspect "keyboard" does not equal
           ↳ the second occurrence of the same aspect term within the
           ↳ sentence)
           '''
190
191         if self.headIndex == otherAspect.headIndex:
192             if self.termSize == otherAspect.termSize:
193                 return True
194         return False
195
196     def compare(self, otherAspect):
197         ''' Comparison based on the words within the aspect - returns
           ↳ true if all Tokens within the aspect are equivalent.
           '''
198
199         result = False
200         if self.termSize == otherAspect.termSize:
201             result = True
202             for i in range(termSize):
203                 if self.tokenized_term[i].text !=
           ↳ otherAspect.tokenized_term[i].text:
204                     result = False
205                     break
206         return result
207
208     def setTokens(self, tokens):
209         self.tokens = tokens
210
211     def getHeadToken():
212         return self.tokens[0]
213
214     def setIndices(self, headIndex):
215         self.headIndex = headIndex
216         self.endIndex = headIndex + self.termSize
217
218     def setOffsets(self, offsets):
219         self.offsets = offsets
220
221     def setPolarity(self, polarity):
222         self.polarity = polarity
223
224     class Instance:
225         '''An instance is a sentence, modeled out of XML (pre-specified
           ↳ format, based on the 4th task of SemEval 2014). It contains the
           ↳ text, the aspect terms, and any aspect categories.
           '''
226
227
```

```

228     def __init__(self, element):
229         self.text = element.find('text').text
230         self.id = element.get('id')
231         self.generateTokens()
232         self.aspect_terms = [Aspect('', id=self.id).createEl(e) for es
233             ↪ in
234                 element.findall('aspectTerms') for e in es
235                 ↪ if
236                 es is not None]
237         self.aspect_categories = [Category(term='',
238             ↪ polarity='').createEl(e) for es in
239             ↪ element.findall('aspectCategories')
240                 for e in es if
241                 es is not None]
242         self.updateAspectFields() # Updates Aspect features related to
243             ↪ Tokens, and vice versa
244
245     def generateTokens(self):
246         ''' Generate tokens based on the tokenization of the sentence.
247             '''
248
249         # Tokenize text and create Token object list
250         self.tokenized_text = Tokenizer().tokenize(self.text)
251         self.tokens = [Token(self.tokenized_text[i], i) for i in
252             ↪ range(len(self.tokenized_text))]
253
254         # Update the POS tag for each Token object
255         tagged_text = nltk.pos_tag(self.tokenized_text)
256         for i in range(len(self.tokens)):
257             self.tokens[i].setPOS(tagged_text[i][1])
258
259         # Update the next and previous tokens for each Token object
260         for i in range(len(self.tokens)):
261             token = self.tokens[i]
262             if i == 0 and i == (len(self.tokens) - 1):
263                 token.setPrev(Token())
264                 token.setNext(Token(index=len(self.tokens)))
265             elif i == 0:
266                 token.setPrev(Token())
267                 token.setNext(self.tokens[i+1])
268             elif i == (len(self.tokens) - 1):
269                 token.setPrev(self.tokens[i-1])
270                 token.setNext(Token(index=len(self.tokens)))
271             else:
272                 token.setPrev(self.tokens[i-1])
273                 token.setNext(self.tokens[i+1])
274
275     def updateAspectFields(self):

```



```

270     ''' Update some token-based fields of Aspects, and aspect-based
    ↪ fields of Tokens
271     '''
272     for at in self.aspect_terms:
273         # Find the aspect term within the sentence, then update the
    ↪ indices of the tokens.
274         at.setIndices(self.findHeadIndex(at))
275
276         # Update the tokens' IOB2 fields
277         self.tokens[at.headIndex].setActualIOB2("B")
278         for i in range(at.headIndex+1, at.endIndex):
279             self.tokens[i].setActualIOB2("I")
280
281         # Add a list of the Token objects for the aspect term
282         at.setTokens(self.tokens[at.headIndex:at.endIndex])
283
284     ''' NOTE: No longer needed
285     def predictedFromIOB2(self):
286         Given an instance with predicted IOB2 tags, return a list of
    ↪ predicted Aspects
287
288         term = []
289         termList = []
290         i = 0
291         while i < len(self.tokens):
292             t = self.tokens[i]
293             if t.predictedIOB2 == "B":
294                 term.append(t)
295                 while i+1 < len(self.tokens):
296                     if self.tokens[i+1].predictedIOB2 == "I":
297                         term.append(self.tokens[i+1])
298                         i = i + 1
299                 else:
300                     break
301                 termList.append(term)
302                 term = []
303             i = i + 1
304         return termList
305     '''
306
307     def findHeadIndex(self, at):
308         ''' Two challenges here: we must account for multi-word aspect
    ↪ terms, and we must account for duplicates of the term that
    ↪ may exist in the sentence. '''
309         headToken = at.tokenized_term[0] # The first token of
    ↪ the aspect term(if multiple tokens are in the word/phrase)
310         headCount = self.text.count(headToken) # Count how many times
    ↪ the first word in the aspect term appears in the sentence

```

```

311         index = -1                                # The index we're
           ↪ looking for - will eventually be returned
312
313         # If there is only one occurrence of the aspect term's first
           ↪ word:
314         if headCount == 1:
315             return findNthOccurrence(self.tokenized_text, headToken, 1)
316
317         # If there are multiple occurrences, find the correct occurrence
           ↪ and then find its' index in the token list
318         else:
319             n = 1                                # The nth occurrence of the word is the one
           ↪ we're searching for
320             loc = -1                             # The current location within the sentence
           ↪ string
321             while n <= headCount:
322                 # Find the next occurrence and check if it matches the
           ↪ listed beginning offset.
323                 loc = self.text.find(headToken, loc+1)
324                 if loc == int(at.offsets['from']):
325                     # Find the index in the tokens of the nth occurrence
           ↪ of the term
326                     return findNthOccurrence(self.tokenized_text,
           ↪ headToken, n)
327                 n = n + 1
328             return -1
329
330     def adjustFormat(self):
331         ''' For evaluation purposes. Returns a list of (Term, Indices)
           ↪ tuples, where Indices is a tuple
           '''
332
333         output = []
334         for at in self.aspect_terms:
335             term = " ".join(at.lower_tokenized_term)
336             indices = tuple([token.index for token in at.tokens])
337             output.append((term, indices))
338
339         return output
340
341     def get_aspect_terms(self):
342         return [a.lower_term for a in self.aspect_terms]
343
344     def get_aspect_categories(self):
345         return [c.term.lower() for c in self.aspect_categories]
346
347     def get_predicted_terms(self):
348         return [a.lower_term for a in self.predicted_terms]
349

```

```
350     def get_predicted_categories(self):
351         return [c.term.lower() for c in self.predicted_categories]
352
353     def add_aspect_term(self, term, offsets='', id=''):
354         a = Aspect(term, id)
355         if offsets != '':
356             a.setOffsets(offsets)
357         self.aspect_terms.append(a)
358
359     def add_aspect_category(self, term, polarity=''):
360         c = Category(term, polarity)
361         self.aspect_categories.append(c)
362
363     def add_predicted_term(self, term, id=''):
364         a = Aspect(term, id)
365         self.predicted_terms.append(a)
366
367     def add_predicted_category(self, term, polarity=''):
368         c = Category(term, polarity)
369         self.predicted_categories.append(c)
370
371     class Corpus:
372         '''A corpus contains instances, and is useful for training
373         ↪ algorithms or splitting to train/test files.'''
374
375         def __init__(self, elements):
376             self.corpus = [Instance(e) for e in elements]
377             self.texts = [t.text for t in self.corpus]
378             self.size = len(self.corpus)
379             self.aspect_terms_fd = fd([" ".join(a.tokenized_term) for i in
380             ↪ self.corpus for a in i.aspect_terms])
381             self.top_aspect_terms = freq_rank(self.aspect_terms_fd)
382
383         def __iter__(self):
384             for i in self.corpus:
385                 yield i.tokenized_text
386
387         def top_text_terms(self):
388             ''' Old version of top_aspect_terms
389             '''
390             aspect_terms_fd = fd([a for i in self.corpus for a in
391             ↪ i.get_aspect_terms()])
392             return freq_rank(self.aspect_terms_fd)
393
394         def clean_tags(self):
395             for i in range(len(self.corpus)):
396                 self.corpus[i].aspect_terms = []
```

```
395 def split(self, threshold=0.8, shuffle=False):
396     '''Split to train/test, based on a threshold. Turn on shuffling
    ↪ ↪ for randomizing the elements beforehand.'''
397     clone = copy.deepcopy(self.corpus)
398     if shuffle: random.shuffle(clone)
399     train = clone[:int(threshold * self.size)]
400     test = clone[int(threshold * self.size):]
401     return train, test
402
403 def getPolarityTermDict(self):
404     ''' Returns a dictionary where each aspect term is associated
    ↪ ↪ with a dictionary,
    ↪ '''
405
406     return fd2([(at.term, at.polarity) for at in s.aspect_terms for
    ↪ s in self.corpus])
407
408 def getPolarityCategoryDict(self):
409     ''' Returns a dictionary where each aspect category is
    ↪ ↪ associated with a dictionary
    ↪ '''
410
411     return fd2([(ac.term, ac.polarity) for ac in s.aspect_categories
    ↪ for s in self.corpus])
```

Appendix - Aspect Identification

```
1 import time
2 import math
3 from collections import defaultdict
4 import xml.etree.ElementTree as ET
5 from libraries.structure import Corpus
6 from libraries.structure import fd
7 from libraries.structure import freq_rank
8 from libraries.structure import generate
9
10 from stanford_corenlp_python import jsonrpc
11
12 import nltk
13 import nltk.corpus, nltk.tag
14 from nltk import word_tokenize
15 from nltk.tokenize import WordPunctTokenizer
16 from nltk.tokenize import TreebankWordTokenizer as Tokenizer
17 import nltk.chunk as chunk
18 from nltk.stem.porter import PorterStemmer as Stemmer
19
20 import pycrfsuite
21
22 ##### Association Mining Method #####
23
24 def HuLiu(instances, minSupport = 1.0, minPsupport = 2, maxWordDist =
  ↳ 1.0, maxNonCompact = 1):
25     ''' Hu and Liu's algorithm for aspect term extraction. Returns two
  ↳ arguments: a dictionary containing all predicted terms with
  ↳ their associated p-support, and a list of sentences with the
  ↳ aspect terms in each sentence.
26         instances = a list of Sentence
27         minSupportPercentage = the percentage of sentences the term must
  ↳ appear in to be considered "frequent"
28         minPsupport = the minimum number of sentences in which a
  ↳ candidate term must occur (ignoring any times another candidate term
  ↳ in the sentence subsumes the current candidate term)
29         maxWordDist = the maximum distance allowed between words in a
  ↳ candidate term
```

```

30     maxNonCompact = the maximum number of sentences within the
↳ corpus in which a candidate term can violate the maximum word
↳ distance
31     '''
32
33     # We store the terms by sentence and as a set.
34     terms = dict()                                # Stores terms in a
↳ variety of formats
35     terms['sent'] = dict()                        # Stores the terms of
↳ each sentence separately
36     terms['sent']['(Term,Indices)'] = []         # Stores the terms of
↳ each sentence as a tuple: (FullString, IndicesTuple). FullString
↳ has " " between tokens. IndicesTuple is a tuple containing the
↳ indices of the tokens in String.
37     terms['all'] = dict()                        # Stores terms in one
↳ group, not as a list of sentences
38     terms['all']['set'] = set()                 # The entire set of
↳ distinct terms
39
40     tbc = 0
41     # treebank chunking
42     if tbc:
43         treebank_sents = nltk.corpus.treebank_chunk.chunked_sents()
44         train_chunks = conll_tag_chunks(treebank_sents)
45         u_chunker = nltk.tag.UnigramTagger(train_chunks)
46         ub_chunker = nltk.tag.BigramTagger(train_chunks,
↳ backoff=u_chunker)
47         ubt_chunker = nltk.tag.TrigramTagger(train_chunks,
↳ backoff=ub_chunker)
48         ut_chunker = nltk.tag.TrigramTagger(train_chunks,
↳ backoff=u_chunker)
49         utb_chunker = nltk.tag.BigramTagger(train_chunks,
↳ backoff=ut_chunker)
50         # Find nouns and noun phrases in each sentence - these are
↳ initial candidate terms. Nouns are lists of (String, Index)
↳ tuples
51         nounsBySentence = [nounsAndPhrasesInSentence(s,
↳ chunker=ub_chunker, tbc=True) for s in instances]
52     else:
53         # Find nouns and noun phrases in each sentence - these are
↳ initial candidate terms. Nouns are lists of (String, Index)
↳ tuples
54         nounsBySentence = [nounsAndPhrasesInSentence(s, ne=True) for s
↳ in instances]
55
56     # Include combined pairs and triples of nouns / phrases within
↳ sentences as candidate terms, then get a dictionary of their
↳ frequencies (support)

```

```

57     temp = [[t2list(term) for term in (s + getPairs(s) + getTriples(s))]
58             ↪ for s in nounsBySentence]
59
60     support = fd([term[0] for s in temp for term in s])
61
62     # Get all frequent candidate terms - those that meet the minimum
63     ↪ support.
64     terms['sent']['(Term,Indices)'] = [[term for term in s if
65             ↪ support[term[0]] >= minSupport] for s in temp]
66
67     # Update support
68     support = fd([term[0] for s in terms['sent']['(Term,Indices)'] for
69             ↪ term in s])
70
71     # Store the set of current candidate terms
72     terms['all']['set'] = set(support.keys())
73
74     nonCompact = dict.fromkeys(terms['all']['set'], 0)
75     ↪ # Stores occurrences of non-compact form for each term
76     pSupport = fd([term[0] for s in terms['sent']['(Term,Indices)'] for
77             ↪ term in removeSubsets(s)]) # Stores p-support of each term
78     isSubset = dict.fromkeys(terms['all']['set'])
79
80     for term in terms['all']['set']:
81         isSubset[term] = False
82         for term2 in terms['all']['set']:
83             if term in term2:
84                 if term != term2:
85                     isSubset[term] = True
86                     continue
87
88     # Check to see if the distance between words exceeds maxDist
89     for sentence in terms['sent']['(Term,Indices)']:
90         for term in sentence:
91             indices = term[1]
92             if len(indices) <= 1:
93                 # Term has only one word - skip to next term
94                 continue
95             max = maxDist(indices)
96             if max > maxWordDist:
97                 nonCompact[term[0]] = nonCompact[term[0]] + 1
98
99     # Remove terms that appear in non-compact form more than
100    ↪ "maxNonCompact" times. Also, remove terms below the minimum
101    ↪ p-support threshold that are a subset of some other term.
102    newTerms = set()
103
104    sub = 0
105
106    nc = 0

```

```

97     for term in terms['all']['set']:
98         if nonCompact[term] > maxNonCompact:
99             # Condition violated, term is removed
100            nc += 1
101            continue
102        if pSupport[term] >= minPsupport:
103            # Term meets minimum p-support; term is kept
104            newTerms.add(term)
105        else:
106            if isSubset[term]:
107                # Term is a part of another aspect term; term is removed
108                sub += 1
109                continue
110            else:
111                # Term is not part of another term; term is kept
112                newTerms.add(term)
113
114    terms['all']['set'] = newTerms
115
116    # Update terms['sent']['(Term,Indices)']
117    newTermSents = []
118    for sentence in terms['sent']['(Term,Indices)']:
119        newSent = []
120        for term in sentence:
121            if term[0] in terms['all']['set']:
122                newSent.append(term)
123        newTermSents.append(newSent)
124    terms['sent']['(Term,Indices)'] = newTermSents
125
126    # Update p-support values
127    pSupport = fd([term[0] for s in terms['sent']['(Term,Indices)'] for
128    ↪ term in removeSubsets(s)]) # Stores p-support of each term
129
130    return support, terms['sent']['(Term,Indices)']
131
132 def nearestNoun(sentence, adjIndex):
133     ''' Returns the nearest noun to a given term in a sentence. Sentence
134     ↪ is an Instance, adjIndex is the index of the adjective in the
135     ↪ sentence.
136     Returns None if there are no nouns in the sentence'''
137     nouns = [(token.text, token.index) for token in sentence if
138     ↪ token.POS2 == "NN"]
139     if len(nouns) == 0:
140         return None
141     adjIndex = adj[1]
142     nearest = None
143     minDist = float("inf")
144     for noun in nouns:

```



```

141         dist = abs(adjIndex - noun[1])
142         if dist < minDist:
143             minDist = dist
144             nearest = noun
145     return nearest
146
147 def maxDist(indices):
148     ''' Takes a list of indices (ex: [1, 3, 4, 5]
149         Returns the max distance between any 2 adjacent indices '''
150     maxDist = 0
151     for i in range(len(indices)-1):
152         dist = indices[i+1] - indices[i]
153         if dist > maxDist:
154             maxDist = dist
155     return maxDist
156
157 def neChunker(instance):
158     tagged = [(token.text, token.POS) for token in instance.tokens]
159     rawChunks = nltk.chunk.ne_chunk(tagged)
160     (tags, chunks) = zip(*(conll_tag_chunks([rawChunks])[0]))
161     return chunks
162
163 def nounsAndPhrasesInSentence(instance, chunker='', reg=False,
164     ↪ tbc=False, ne=False):
165     ''' Input: A sentence instance
166         Returns a list of lists, each containing the (String, Index)
167     ↪ tuples corresponding to the tokens of a noun or noun phrase
168         '''
169     # Get tagged sentence in the form of a list of (token, POS) tuples
170     tagged = [(token.text, token.POS) for token in instance.tokens]
171     # Create a list of (String, Index) tuples for each noun / noun
172     ↪ phrase
173     nouns = []
174     if reg==True:
175         # Find noun phrases using regex
176         pattern = r"""
177             NBAR:
178             {<NN.*|JJ>*<NN.*>} # Nouns and Adjectives, terminated with
179     ↪ Nouns
180             NP:
181             {<NBAR>}
182             {<NBAR><IN><NBAR>} # Above, connected with in/of/etc...
183             """
184         NPChunker = nltk.RegexpParser(pattern)
185         tagged = NPChunker.parse(tagged)
186         # cTagged = chunk.ne_chunk(tagged)
187         nounIndices = chunkParse(tagged)

```

```

185     for n in nounIndices:
186         nList = []
187         for i in n:
188             nList.append((instance.tokens[i].lower_text,
189                 ↪ instance.tokens[i].index))
189             nouns.append(nList)
190 elif tbc==True and chunker == '':
191     print "ERROR: Chunker must be provided."
192     exit()
193 else:
194     if ne == True:
195         chunks = neChunker(instance)
196     elif tbc == True:
197         (words, tags) = zip(*tagged)
198         (tags2, chunks) = zip(*chunker.tag(tags))
199     else:
200         print "Error"
201         exit()
202     n = []
203
204     # Iterate over tokens
205     for i in range(len(instance.tokens)):
206         # Add nouns outside of noun phrases
207         if chunks[i] == '0':
208             if tagged[i][1].startswith('NN'):
209                 nouns.append([(instance.tokens[i].lower_text, i)])
210             # Start or continue building noun phrase
211             else:
212                 n.append((instance.tokens[i].lower_text, i))
213
214             # Check if current token is the last token
215             if i+1 >= len(instance.tokens):
216                 # If we were building a noun, add it
217                 if len(n) > 0:
218                     nouns.append(n)
219                     n = []
220             #
221             elif chunks[i+1].startswith('I') == False:
222                 nouns.append(n)
223                 n = []
224
225     return nouns
226
227 def getPairs(terms):
228     ''' Given a list of terms stored as lists of (String, Index) tuples,
229     ↪ return all pairs (e.g. [1,2], [2,3], [3,4], etc.)
230     '''
231     pairs = []

```

```

231     if len(terms) >= 2:
232         for i in range(len(terms)-1):
233             pairs.append(terms[i] + terms[i+1])
234     return pairs
235
236 def getTriples(terms):
237     ''' Given a list of terms stored as lists of (String, Index) tuples,
238     ↪ return all triples (e.g. [1,2,3], [2,3,4], etc.)
239     Input example: [ [('Microsoft',2), ('Office',3)],
240     ↪ [ ('Word',5)], [ ('Key',8), ('Board',9)] ]
241     Output example: [ [('Microsoft',2), ('Office',3), ('Word',5)],
242     ↪ [ ('Word',5), ('Key',8), ('Board',9)] ]
243     '''
244     triples = []
245     if len(terms) >= 3:
246         for i in range(len(terms)-2):
247             triples.append(terms[i] + terms[i+1] + terms[i+2])
248     return triples
249
250 def removeSubsets(terms):
251     ''' Given a list of (Term, Indices) tuples, return a list of (Term,
252     ↪ Indices) tuples without any subsets (strings that are substrings
253     ↪ of another string in the list)
254     '''
255     newTerms = []
256     for i in range(len(terms)):
257         subset = False
258         for j in range(len(terms)):
259             if i != j and set(terms[i][1]) < set(terms[j][1]):
260                 subset = True
261         if subset == False:
262             newTerms.append(terms[i])
263     return newTerms
264
265 def t2list(term):
266     ''' Given a term stored as a list of (String, Index) tuples, return
267     ↪ a tuple: (Full_String, Index_Tuple)
268     '''
269     if len(term) == 0:
270         return term
271     t = zip(*term)
272     return (" ".join(t[0]), t[1])
273
274 def chunkParse(cTagged):

```

```

269     ''' Given a sentence tagged with chunks (using nltk.pos_tag and
    ↪ RegexpParser), return a list. Each element is a list itself,
    ↪ containing the indices of each noun / noun phrase;
    ↪ single-element lists are a single index and correspond to
    ↪ single-word nouns. Multi-element lists store a list of indices
    ↪ in sequence, and are noun phrases.
    '''
270
271     index = 0
272     parsed = []
273     for i in range(len(cTagged)):
274         if type(cTagged[i]) is nltk.tree.Tree:
275             parsed.append(range(index, index + len(cTagged[i])))
276             index = index + len(cTagged[i])
277         else:
278             if cTagged[i][0].startswith('NN'):
279                 parsed.append([index])
280                 index = index + 1
281     return parsed
282
283 def conll_tag_chunks(chunk_sents):
284     tag_sents = [nltk.chunk.tree2conlltags(tree) for tree in
    ↪ chunk_sents]
285     return [[(t, c) for (w, t, c) in chunk_tags] for chunk_tags in
    ↪ tag_sents]
286
287 ##### Conditional Random Fields #####
288
289 def crf(train, test, alg = ""):
290     # Convert sentences to appropriate feature/label format
291     train_features, train_labels = generate(train)
292     test_features, test_labels = generate(test)
293
294     # Train CRF
295     trainer = pycrfsuite.Trainer()
296     for x,y in zip(train_features, train_labels):
297         trainer.append(x,y)
298
299     if alg != "":
300         trainer.select(alg)
301     trainer.set_params({
302         # 'c1': 1.0, # coefficient for L1 penalty
303         # 'c2': 1e-3, # coefficient for L2 penalty
304         # 'max_iterations': 50, # stop earlier
305
306         # include transitions that are possible, but not
    ↪ observed
307         'feature.possible_transitions': False
308     })

```

```

309     trainer.train('conll2002-esp.crfsuite')
310     trainer.logparser.last_iteration
311
312     # Test CRF
313     tagger = pycrfsuite.Tagger()
314     tagger.open('conll2002-esp.crfsuite')
315
316     predicted_labels = [tagger.tag(s) for s in test_features]
317     print "# Sentences: " + str(len(test_labels))
318     confusion = dict()
319     confusion['B'] = {'actual':0, 'predicted':0, 'truePos':0,
320 ↪ 'falseNeg':0, 'falsePos':0}
321     confusion['I'] = {'actual':0, 'predicted':0, 'truePos':0,
322 ↪ 'falseNeg':0, 'falsePos':0}
323     confusion['O'] = {'actual':0, 'predicted':0, 'truePos':0,
324 ↪ 'falseNeg':0, 'falsePos':0}
325     for s in range(len(test_labels)):
326         for t in range(len(test_labels[s])):
327             actual = test_labels[s][t]
328             pred = predicted_labels[s][t]
329             confusion[actual]['actual'] = confusion[actual]['actual'] +
330 ↪ 1
331             confusion[pred]['predicted'] = confusion[pred]['predicted']
332 ↪ + 1
333             if actual == pred:
334                 confusion[actual]['truePos'] =
335 ↪ confusion[actual]['truePos'] + 1
336             else:
337                 confusion[actual]['falseNeg'] =
338 ↪ confusion[actual]['falseNeg'] + 1
339                 confusion[pred]['falsePos'] =
340 ↪ confusion[pred]['falsePos'] + 1
341     return predicted_labels
342
343     ##### ATE Evaluation Methods #####
344
345     # Evaluate ATE methods
346     def evaluateAspectTerms(trueFD, trueSent, predictedFDs, predictedSents,
347 ↪ methodNames = [], toScreen = False):
348         ''' Evaluate ATE methods
349             '''
350
351         distinct = {"TP":{}, "FN":{}, "FP":{}, "P":{}, "R":{}, "F":{}}
352         instances = {"TP":{}, "FN":{}, "FP":{}, "P":{}, "R":{}, "F":{}}
353         weighted = {}
354         if len(methodNames) == 0:
355             methodNames = [str(i) for i in range(1,len(predictedFDs)+1)]

```

```

348     for i in range(len(predictedFDs)):
349         predictedFD = predictedFDs[i]
350         predictedSent = predictedSents[i]
351
352         # Distinct
353         TP, FN, FP, P, R, F =
354             ↪ evaluateAspectTermsDistinct(set(trueFD.keys()),
355             ↪ set(predictedFD.keys()))
356         distinct["TP"][methodNames[i]] = TP
357         distinct["FN"][methodNames[i]] = FN
358         distinct["FP"][methodNames[i]] = FP
359         distinct["P"][methodNames[i]] = P
360         distinct["R"][methodNames[i]] = R
361         distinct["F"][methodNames[i]] = F
362
363         # Instances
364         TP, FN, FP, P, R, F = evaluateAspectTermsInstances(trueSent,
365             ↪ predictedSent, trueFD, predictedFD)
366         instances["TP"][methodNames[i]] = TP
367         instances["FN"][methodNames[i]] = FN
368         instances["FP"][methodNames[i]] = FP
369         instances["P"][methodNames[i]] = P
370         instances["R"][methodNames[i]] = R
371         instances["F"][methodNames[i]] = F
372
373         # Average weighted precision
374         weighted[methodNames[i]] = evaluateAspectTermsWeighted(trueFD,
375             ↪ predictedFD)
376
377     if toScreen:
378         print "----- Distinct: -----"
379         for m in methodNames:
380             print m
381             print "TP: " + str(distinct["TP"][m])
382             print "FN: " + str(distinct["FN"][m])
383             print "FP: " + str(distinct["FP"][m])
384             print "P: " + str(distinct["P"][m])
385             print "R: " + str(distinct["R"][m])
386             print "F: " + str(distinct["F"][m])
387             print ""
388
389         print "----- Instances: -----"
390         for m in methodNames:
391             print m
392             print "TP: " + str(instances["TP"][m])
393             print "FN: " + str(instances["FN"][m])
394             print "FP: " + str(instances["FP"][m])
395             print "P: " + str(instances["P"][m])

```

```

392         print "R: " + str(instances["R"][m])
393         print "F: " + str(instances["F"][m])
394         print ""
395
396         print "----- Average Weighted Precision:
397         ↪ -----"
398         for m in methodNames:
399             print m
400             print str(weighted[m])
401             print ""
402
403     return distinct, instances, weighted
404
405 # Evaluate by distinct aspect terms
406 def evaluateAspectTermsDistinct(trueSet, predictedSet, toScreen =
407     ↪ False):
408     ''' Input: two sets of distinct aspect terms (the predicted set and
409     ↪ the true set) for a corpus
410     Output: Evaluation metrics for distinct aspect terms'''
411
412     truePos = predictedSet.intersection(trueSet)
413     falseNeg = trueSet.difference(predictedSet)
414     falsePos = predictedSet.difference(trueSet)
415     TP = len(truePos)
416     FN = len(falseNeg)
417     FP = len(falsePos)
418     if TP == 0:
419         P = 0.0
420         R = 0.0
421         F = 0.0
422     else:
423         P = float(TP)/(TP + FP)
424         R = float(TP)/(TP + FN)
425         F = 2.0*P*R/(P+R)
426
427     '''
428     print "\nEvaluate by Distinct Term: "
429     print "Predicted: " + str(len(ptSet))
430     print "Actual: " + str(len(atSet))
431     '''
432
433     if toScreen:
434         print "True Positive: %f -- False Negative: %f -- False
435         ↪ Positive: %f (P = %d, R = %d, F = %d)" % TP, FN, FP, P, R, F
436
437     return TP, FN, FP, P, R, F
438
439

```

```

435 def evaluateAspectTermsInstances(trueTermsBySent, predictedTermsBySent,
↳ trueFD="", predictedFD="", toScreen = False):
436     ''' Input: predictedTerms (a list of sentences, where each sentence
↳ is expressed as a list of its predicted aspect terms in (Term,
↳ Indices) format) and actualTerms (same as predictedTerms, but
↳ with the human-annotated terms). If true and predicted frequency
↳ dictionaries are not specified, they are computed.
437     Output: Evaluation metrics for instances of aspect terms'''
438
439     if trueFD == "":
440         trueFD = fd([term[0] for s in trueTermsBySent for term in s])
441     if predictedFD == "":
442         predictedFD = fd([term[0] for s in predictedTermsBySent for term
↳ in s])
443
444     # Below is a version of the code that gives frequencies for true
↳ positive, false negative, and false positive for each word:
445
446     truePosList = [] # List for terms in both actual and predicted
447     falseNegList = [] # List for terms in actual but not predicted
448     falsePosList = [] # List for terms in predicted but not actual
449
450     # Check whether terms are in actual, predicted, or both
451     for i in range(len(trueTermsBySent)):
452         # Get true and predicted sets of term indices tuples from
↳ sentence
453         trueIndices = set([term[1] for term in trueTermsBySent[i]])
454         predictedIndices = set([term[1] for term in
↳ predictedTermsBySent[i]])
455
456         # Update TP, FN, and FP using sets
457         truePosList.extend([term[0] for term in trueTermsBySent[i] if
↳ term[1] in trueIndices.intersection(predictedIndices)])
458         falseNegList.extend([term[0] for term in trueTermsBySent[i] if
↳ term[1] in trueIndices.difference(predictedIndices)])
459         falsePosList.extend([term[0] for term in predictedTermsBySent[i]
↳ if term[1] in predictedIndices.difference(trueIndices)])
460
461     truePos = fd(truePosList)
462     falseNeg = fd(falseNegList)
463     falsePos = fd(falsePosList)
464     '''
465     print "-----TruePos-----"
466     print list(truePos.keys())[:100]
467     print "-----FalseNeg-----"
468     print list(falseNeg.keys())[:100]
469     print "-----FalsePos-----"
470     print list(falsePos.keys())[:100]

```



```

471     '''
472
473     TP = sum(truePos.values())
474     FN = sum(falseNeg.values())
475     FP = sum(falsePos.values())
476
477     if TP == 0:
478         P = 0.0
479         R = 0.0
480         F = 0.0
481     else:
482         P = float(TP)/(TP + FP)
483         R = float(TP)/(TP + FN)
484         F = 2.0*P*R/(P+R)
485
486     '''
487     print "\nEvaluate by Instance: "
488     print "Predicted: " +
↪ str(sum(train_data.predicted_terms_fd.values()))
489     print "Actual: " + str(sum(train_data.aspect_terms_fd.values()))
490     '''
491     if toScreen:
492         print "True Positive: %f -- False Negative: %f -- False
↪ Positive: %f (P = %d, R = %d, F = %d)" % TP, FN, FP, P, R, F
493
494     return TP, FN, FP, P, R, F
495
496 def evaluateAspectTermsWeighted(trueFD, predictedFD, toScreen = False):
497     ''' Inputs: dictionaries of term frequencies, both true and
↪ predicted. toScreen specifies whether to print output or not
498     Outputs: the average weighted precision
499     '''
500     trueFD_sorted = freq_rank(trueFD)
501     trueRanked = {trueFD_sorted[i]:(i+1) for i in
↪ range(len(trueFD_sorted))}
502     awp = avgWeightedPrecision(trueRanked, freq_rank(predictedFD))
503     if toScreen:
504         print "Average weighted precision: " + str(awp)
505     return awp
506
507 def weightedPrecision(trueSet, predictedFreqRank, m):
508     ''' Inputs: a set of true aspect terms, a list of predicted terms
↪ (in order of decreasing frequency), and a parameter m.
509     Output: the weighted precision of the first m predicted terms.
510     '''
511     predicted = predictedFreqRank[0:m]
512     wp = 0.0
513     denom = 0.0

```

```

514     for i in range(m):
515         denom = denom + 1.0/(i+1.0)
516         if predicted[i] in trueSet:
517             wp = wp + 1.0/(i+1.0)
518     return wp/denom
519
520 def weightedRecall(trueRanked, predictedFreqRank, m):
521     ''' Inputs: a dictionary of true aspect terms, where values are
522     ↪ their frequency rank (ex: the kth most frequent term has value
523     ↪ k), a list of predicted terms, and a parameter m.
524     Output: the weighted recall of the first m predicted terms.
525     '''
526     predicted = predictedFreqRank[0:m]
527     wr = 0.0
528     denom = 0.0
529     # Compute the numerator
530     for i in range(m):
531         if predicted[i] in trueRanked:
532             # Sum the reciprocal of the
533             wr = wr + 1.0/trueRanked[predicted[i]]
534     # Compute the denominator
535     for i in range(len(trueRanked)):
536         denom = denom + 1.0/(i+1.0)
537     return wr/denom
538
539 def avgWeightedPrecision(trueRanked, predictedFreqRank):
540     ''' Inputs: a dictionary of true aspect terms, where values are
541     ↪ their frequency rank (ex: the kth most frequent term has value
542     ↪ k), a list of predicted terms, and a parameter m.
543     Output: the weighted recall of the first m predicted terms.
544     '''
545     awp = 0.0
546     wr = [weightedRecall(trueRanked, predictedFreqRank, m) for m in
547     ↪ range(1, len(predictedFreqRank) + 1)]
548     wp = [weightedPrecision(trueRanked, predictedFreqRank, m) for m in
549     ↪ range(1, len(predictedFreqRank) + 1)]
550     for i in range(11):
551         r = i/10.0
552         max = 0.0
553         for m in range(0, len(predictedFreqRank)):
554             if wr[m] >= r:
555                 if wp[m] > max:
556                     max = wp[m]
557         awp = awp + max
558     return awp
559
560 ##### Data Processing #####

```

```

556 def IOB2toAspectTerms(IOB2labels, sentences):
557     ''' Input: A list of IOB2 labels corresponding to sentences, and a
    ↪ list of sentences (Instance objects)
558         Output: support (a dictionary of aspect terms, each token
    ↪ separated with " ", and frequencies) and predictedSentences (a list
    ↪ of sentences, each stored as a list of (Term, Indices) tuples.
559         '''
560     predictedTermsBySentence = []
561     for i in range(len(sentences)):
562         predictedTermsBySentence.append([])
563         labels = IOB2labels[i]
564         sentence = sentences[i]
565         term = ""
566         indices = []
567         for j in range(len(sentence.tokens)):
568             token = sentence.tokens[j]
569             if labels[j] == "B":
570                 term = token.lower_text
571                 indices.append(token.index)
572             if labels[j] == "I":
573                 term = term + " " + token.lower_text
574                 indices.append(token.index)
575             if ((j+1) == len(sentence.tokens)) and (len(term) > 0):
576                 predictedTermsBySentence[i].append((term,
    ↪ tuple(indices)))
577                 term = ""
578                 indices = []
579                 continue
580             if (labels[j] != "O") and (labels[j+1] != "I"):
581                 predictedTermsBySentence[i].append((term,
    ↪ tuple(indices)))
582                 term = ""
583                 indices = []
584     support = fd([term[0] for s in predictedTermsBySentence for term in
    ↪ s])
585     return support, predictedTermsBySentence

```

Appendix - Sentiment Analysis

```
1 import time
2 from collections import defaultdict
3 import xml.etree.ElementTree as ET
4 from libraries.structure import Corpus
5 from libraries.structure import fd
6 from libraries.structure import freq_rank
7 from libraries.structure import generate
8
9 from stanford_corenlp_python import jsonrpc
10
11 import nltk
12 from nltk import word_tokenize
13 from nltk.tokenize import WordPunctTokenizer
14 from nltk.tokenize import TreebankWordTokenizer as Tokenizer
15 from nltk.corpus import sentiwordnet as swn
16 import nltk.chunk as chunk
17 from nltk.stem.porter import PorterStemmer as Stemmer
18 from nltk.corpus import wordnet as wn
19 from nltk.sentiment.vader import SentimentIntensityAnalyzer
20
21 import pycrfsuite
22
23 ##### Extracting adjectives from text
24 → #####
25
26 def getOpinionAdjs(instances, terms):
27     ''' Input: a list of Instances (corpus) and terms (organized as a
28     → list of sentences, where each sentence is a list of terms with
29     → the format (Term, Indices)
30     Output: the nearest opinion adjective to each term, organized as
31     → a dictionary where frequency is the value
32     '''
33     # Create a dictionary of opinion adjectives
34     opinionAdjs = defaultdict(lambda:0)
35     for i in range(len(terms)):
36         sentence = terms[i]
37         for term in sentence:
```

```

34         adj = nearestAdj(instances[i], term)
35         if adj != None:
36             opinionAdjs[adj[0]] = opinionAdjs[adj[0]] + 1
37     return opinionAdjs
38
39 def getTermsFromAdjs(sentences, terms, opinionAdjs):
40     ''' INCOMPLETE
41         '''
42     # Update candidate terms based on set of adjectives
43     for i in range(len(terms['sent']['(Term,Indices)'))):
44         # Check if there are any candidate terms - if so, move on
45         if len(terms['sent']['(Term,Indices)'][i]) > 0:
46             continue
47         sentence = sentences['ins'][i]
48         adjs = [(token.text, token.index) for token in sentence.tokens
49                 ↪ if (token.POS2 == "JJ" and token.text in opinionAdjs)]
49         for a in adjs:
50             noun = nearestNoun(sentence, a)
51             # If there are multiple opinion adjectives in a sentence,
52             ↪ they may return the same nearest noun
52             if noun != None and noun not in
53             ↪ terms['sent']['(Term,Indices)'][i]:
53                 terms['all']['set'].add(noun)
54                 terms['sent']['(Term,Indices)'][i].append(noun)
55     return terms
56
57 def nearestAdj(sentence, term):
58     ''' Returns the nearest adjective as a tuple (Adj, Index) to a given
59     ↪ term in a sentence. Sentence is a list of Token objects, term is
60     ↪ a (Term, Indices) tuple.
61     Returns "" if there are no viable adjectives in the sentence'''
62     # Find all adjectives in sentence
61     adjs = [(token.text, token.index) for token in sentence if
62             ↪ token.POS2 == "JJ"]
62     if len(adjs) == 0:
63         return None
64     # Find the term's "average" index value (ex: a term with indices [1,
65     ↪ 2, 3, 5] would have a center of 3.75)
65     termIndices = map(list, term[1])
66     termAvgIndex = sum(termIndices) / float(len(termIndices))
67     # Find the adjective closest to the term's "average" index
68     nearest = ""
69     minDist = float("inf")
70     for adj in adjs:
71         # Don't count an adjective if it's already within the candidate
72         ↪ term
72         if adj in term:
73             continue

```

```

74         dist = abs(termAvgIndex - adj[1])
75         if dist < minDist:
76             minDist = dist
77             nearest = adj
78     return nearest
79
80     ##### Extracting polarity scores #####
81
82 def polaritiesByCluster(polarityDict, clusters):
83     clusterPolarityDict = {}
84     for c in clusters:
85         clusterPolarityDict[c] = defaultdict(lambda: 0)
86         for term in clusters[c]:
87             for polType in polarityDict[term]:
88                 clusterPolarityDict[c][polType] +=
89                 ↪ polarityDict[term][polType]
90     return clusterPolarityDict
91
92 def getTermPolarities(instances):
93     ''' Input: A list of instances
94         Output: A dictionary of terms, where each term contains a
95     ↪ dictionary with counts for each polarity category (positive,
96     ↪ negative, neutral, conflict)
97     '''
98     return [(aspect.term, aspect.polarity) for aspect in
99     ↪ instance.aspect_terms] for instance in instances]
100
101 def getCategoryPolarities(instances):
102     ''' Input: A list of instances
103         Output: A list of sentences, where each sentence is a list of
104     ↪ (category, polarity) tuples, where polarity is one of (positive,
105     ↪ negative, neutral, conflict)
106     '''
107     return [(category.term, category.polarity) for category in
108     ↪ instance.aspect_categories] for instance in instances]
109
110 ##### VADER #####
111
112 def vader(sia, sentence):
113     polarity = sia.polarity_scores(sentence)
114     return polarity
115
116 def vaderAdjusted(sia, sentence):
117     polarity = vader(sia, sentence)
118     if polarity['compound'] < 0:
119         return 'negative'
120     elif polarity['neu'] == 1.0:
121         return 'neutral'

```

```

115         #elif abs(polarity['pos'] - polarity['neg']) < 0.05:
116         #return 'conflict'
117     elif polarity['pos'] > polarity['neg']:
118         return 'positive'
119     else:
120         return 'negative'
121
122 def vaderAdjusted2(sia, instance, aspect):
123     polarity = vader(sia, context(instance, aspect))
124     if polarity['compound'] < 0:
125         return 'negative'
126     elif polarity['neu'] == 1.0:
127         return 'neutral'
128     #elif abs(polarity['pos'] - polarity['neg']) < 0.05:
129     #return 'conflict'
130     elif polarity['pos'] > polarity['neg']:
131         return 'positive'
132     else:
133         return 'negative'
134
135 def context(instance, aspect, r=12):
136     avgIndex = aspect.headIndex + (aspect.termSize - 1)/2
137     beg = max(avgIndex - r, 0)
138     end = min(avgIndex + r, len(instance.tokens))
139     return " ".join([token.text for token in instance.tokens[beg:end]])
140
141 def vaderTermPolarities(instances, adjusted = True):
142     ''' Input: A list of instances
143         Output: A list of sentences, where each sentence contains a list
144     ↪ of (term, polarity) tuples. These polarities are estimated from the
145     ↪ VADER sentiment analyzer.
146     '''
147     polarities = []
148     sia = SentimentIntensityAnalyzer()
149     for instance in instances:
150         '''
151         if adjusted:
152             p = vaderAdjusted(sia, instance.text)
153         else:
154             p = vader(sia, instance.text)
155         polarities.append([(aspect.term, p) for aspect in
156     ↪ instance.aspect_terms])
157         '''
158     polarities.append([(aspect.term, vaderAdjusted2(sia, instance,
159     ↪ aspect)) for aspect in instance.aspect_terms])
160     return polarities
161
162 def vaderCategoryPolarities(instances, adjusted = True):

```

```

159     ''' Input: A list of instances
160         Output: A list of sentences, where each sentence contains a list
↳ of (category, polarity) tuples. These polarities are estimated from
↳ the VADER sentiment analyzer.
161         '''
162     polarities = []
163     sia = SentimentIntensityAnalyzer()
164     for instance in instances:
165         if adjusted:
166             p = vaderAdjusted(sia, instance.text)
167         else:
168             p = vader(sia, instance.text)
169         polarities.append([(category.term, p) for category in
↳ instance.aspect_categories])
170     return polarities
171
172     ##### Ratings #####
173
174     def computeRatingsVader(polarityDict):
175         ''' Inputs: a dictionary of aspect terms/categories or clusters,
↳ where each value is a dictionary describing aggregate polarity
↳ scores (ex: {"keyboard":{"positive":4.534, "negative":2.386, "
176             '''
177         ratings = {}
178         for term in polarityDict:
179             p = polarityDict[term]
180             ratings[term] = 4.0*p["positive"] / (p["positive"] +
↳ p["negative"]) + 1
181         return ratings
182
183     def computeRatings(polarityDict):
184         ''' Inputs: a dictionary of aspect terms/categories or clusters,
↳ where each value is a dictionary describing polarity counts (ex:
↳ {"keyboard":{"positive":5, "negative":7, "neutral":2, "conflict":1}})
185
186         Outputs: Ratings are scored as follows:  $4 * ( (P + 0.5*C) / (P + N$ 
↳  $+ 0.5*C) ) + 1$ 
187         '''
188         ratings = {}
189         for t in polarityDict:
190             p = polarityDict[t]
191             ratings[t] = 4.0*(float(p["positive"] + 0.5*p["conflict"]) /
↳ (p["positive"] + p["negative"] + p["conflict"])) + 1
192         return ratings
193
194     ##### ABSA Evaluation Methods #####
195
196     def evaluatePolarities(trueBySent, predictedBySent):

```



```

197     # Create dictionary where confusion[i][j] is the count where a
    ↪ term/category with true polarity i is predicted to have polarity
    ↪ j.
198     confusion = defaultdict(lambda:defaultdict(lambda:0))
199     tot = 0
200     for i in range(len(trueBySent)):
201         for j in range(len(trueBySent[i])):
202             confusion[trueBySent[i][j][1]][predictedBySent[i][j][1]] +=
    ↪ 1
203             tot += 1
204     polTypes = ['positive', 'negative', 'neutral'] #, 'conflict']
205     tot -= sum([confusion['conflict'][j] for j in polTypes])
206     print confusion
207     accuracy = sum(confusion[i][i] for i in polTypes) / float(tot)
208     print accuracy
209     precision = {i:float(confusion[i][i])/sum([confusion[i][j] for j in
    ↪ polTypes]) for i in polTypes}
210     recall = {i:float(confusion[i][i])/sum([confusion[j][i] for j in
    ↪ polTypes]) for i in polTypes}
211     f = {i:(2.0*precision[i]*recall[i]/(precision[i]+recall[i])) for i
    ↪ in polTypes}
212     print "Precision:"
213     print precision
214     print "Recall:"
215     print recall
216     print "F-measure:"
217     print f
218
219     def evaluateRatings(trueRatings, predictedRatings):
220         ''' Input: true and predicted ratings in a dictionary (keys are
    ↪ terms or cluster labels, values are ratings)
221             Output: Evaluation metrics
222             '''
223         if trueRatings.keys() != predictedRatings.keys():
224             print "Error: keys don't match"
225
226         diffs = [abs(trueRatings[t] - predictedRatings[t]) for t in
    ↪ trueRatings]
227         MSE = sum([d^2 for d in diffs])/float(len(trueRatings))
228         print "Number of terms/clusters: %d", len(trueRatings)
229         print "MSE: %f", MSE

```

Biography

Sean Byrne was born in 1994 in the state of Pennsylvania. He attended Lehigh University for his undergraduate education, and was highly involved in the Industrial & Systems Engineering department through various projects with professors and as a member of the ISE Student Council. He graduated from Lehigh University with a Bachelor of Science in Industrial & Systems Engineering and a Bachelor of Science in Mathematics in May 2016. He is now completing a Master of Science degree in Industrial & Systems Engineering through the President's Scholars program, and will graduate in May 2017.