
Sentiment Extraction of Earbud Amazon Reviews

1. Introduction

Amazon.com provides an abundant amount of reviews to aid buyers on their decision to purchase the product, but it is very time consuming to read and filter through all the reviews for relevant information. To solve this issue, we propose to automatically extract sentiments from the entire set of reviews and aggregate it in a concise summary.

Our basic approach is to generate a word bank of relevant terms using LDA from the product descriptions and questions asked by purchasers on the product. Next, using SVM and SentiWordNet with Word Sense Disambiguation (WSD), annotated all sentences with a sentiment score and only use the non-neutral sentences in the n-gram generator to produce relevant adj-noun bigrams or trigrams. These terms will be filtered by the word bank produced by the LDA and ranked to produce a set of most relevant terms with their sentiment. The output will be the sentences these terms were produced from.

2. Problem Definition and Methods

2.1. Task Definition

Our problem is to extract sentiment tagged aspects that will give a good representation of the text corpus it was generated from using the method described below. The main questions we want to answer are:

1. What are characteristics of products important to consumers?
2. What are important features reviewers focus on?
3. What makes review content irrelevant or less relevant?

It is an important problem, because if our algorithm accurately extracts the sentiment, we will be saving consumers lots of time on not just Amazon.com but other large marketplaces. Furthermore, this can be applied to other areas involving text summarization, such as automated generation of catchphrases and key points for company ads.

2.2. Algorithm and Methods

To extract relevant terms from the reviews, each sentence of the reviews was made into a POS-tagged parse tree and chunked on various grammars. Through multiple runs of chunking on different grammars, it was determined that chunking into bigrams of $\langle \text{adj} \rangle^* \langle \text{noun} \rangle^*$ and trigrams of $\langle \text{DT} \rangle? \langle \text{adj} \rangle^* \langle \text{noun} \rangle^*$ were most accurate in extracting relevant terms. To offset the thousands of false positives the parse tree generated, the results of the parse tree were filtered by the word bank generated by the LDA. If the word existed in the LDA, it would be given heavier weight. At the end, only words with a weight above a certain threshold were considered relevant terms.

To give each sentence a sentiment score, we first train a bag of words with a combination of the movie review corpus from nltk and a portion of the annotated laptop reviews from Stanford SNAP library with all the unique words stemmed by using Porter Stemmer. With an interface to SentiWordNet using the NLTK WordNet classes, for each sentence, we split it into individual words and for each word in that sentence, we look it up in the SentiWordNet dictionary for its positive and negative sentiment scores. In order to obtain more accurate results, we also used WSD with finding the closest meaning of that particular word by determining the similarity that word has with the other words in the same sentence. The sentiment score of a sentence would be the summation of the sentiment scores of the words inside that sentence. If a word does not exist in the bag of words we obtained from training, it will simply be ignored.

3. Experimental Evaluation

3.1. Methodology

1. Preprocessing stage:
Using the scraper extract all user reviews and split by sentence to generate a list of sentences from all the reviews per product.
2. Processing stage

- 2.1. Determining Important Features of each product
This will be done with a combined static and dynamic technique.
- 2.1.1. Static Technique
Generate important aspects from product descriptions scraped from Amazon.com by using a TFIDF weighting scheme from NLTK. These words, along with the words sound, audio, comfort, bass, highs, mids, tones, durability, quality, blocking, noise, sharp, deep, and tangle will be considered important aspects automatically.
- 2.1.2. Dynamic Technique
Using a parse tree split each sentence into bigrams and trigrams with the grammar format of `<adj><noun>` to represent candidate product aspects of importance to reviewers. We can filter this list using stopword filters and other grammatical based rules to get rid of false positives. We can create an inverted index of these terms and the sentences they occurred in. We can also easily use the list of sentences to find the frequency of aspects in the review corpus.
- 2.2. Using SVM or Naive Bayes annotate each sentence with a sentiment score based on the previously created list of relevant features for the product. If a sentence has a score below a low threshold then it will be discarded.
3. Post-processing stage: Aggregation
- 3.1. We need to summarize the aspect information we have and present sentences from the corpus that best represent the corpus. One way to do this is to rank the aspects based on frequency and then calculate an averaged sentiment score, which can be used to determine which sentiment polarity the majority of the reviews have with respect to that aspect. We can then take a sentence to represent this aspect. A simple scheme could be to find the sentence with the greatest magnitude sentiment score with the same polarity as the aspect average score.

Evaluation: Compare the results with annotated libraries of training and testing data from a project on reviews to find our accuracy on tagging aspect features and sentiment intensity/polarity. After that we will have to do manual surveying from people of the different summarization techniques that we can

use.

To evaluate our results, we ran the different parts of our algorithm on laptop review data that was annotated with relevant terms per sentence and each of those terms was annotated with either positive, negative, or neutral sentiment. Since we wanted to get the overall sentiment of the sentence, each sentence was given a sentiment score based on majority vote of the sentiment of the terms extracted from it. The testing and training data are realistic, because it is also technology review data that focused on extraction of relevant terms and sentiment as a mean of text summarization.

Each part of the algorithm was run separately on the data to get accuracy of the parts. We assume that if the parts are accurate, the overall goal of extracting relevant terms with sentiment will be accurate as well. The preliminary algorithm was (Chetan write stuffs). The algorithm for extracting relevant terms from the parse tree was run multiple times using different grammars to chunk on in order to determine which grammar produced the most accurate results. In general, the tree chunked on `adj,noun` bigrams or `determiner,adj,noun` trigrams. The results are summarized below:

3.2. Results

3.3. Discussion

The high amount of false positives in term selection via parse tree chunking on bigrams and trigrams is due to the oversensitivity of the chunker. The parse tree will chunk on any phrase that fulfills the grammar, including terms such as father or mother, which are nouns but irrelevant to the reviews.

4. Related Work

We will restrict our focus to one type of product (headphones) to make evaluation easier. We see 3 main stages to this project:

5. Future Work

6. Conclusion

7. Resources

Python libraries: Sci-Kit Learn, NLTK, Numpy, Scipy
Data for benchmarking and testing accuracy

- Web scraper written in Node.js to parse relevant

220	data (reviews, product descriptions, etc)	275
221		276
222		277
223	• Annotated data from SemEval2014-Task4 on	278
224	aspect feature based sentiment analysis	279
225	http://alt.qcri.org/semeval2014/task4/	280
226	index.php?id=data-and-tools	281
227		282
228	• Amazon product review data mined by Stanford	283
229	http://snap.stanford.edu/data/	284
230	web-Amazon-links.html	285
231		286
232	Readings on Parsing & Sentiment Analysis:	287
233		288
234	• Stanford paper on parsing & sentiment analysis:	289
235	http://nlp.stanford.edu/~socherr/	290
236	EMNLP2013_RNTN.pdf	291
237		292
238	• Multi-aspect sentiment analysis	293
239	http://www.cs.cornell.edu/home/cardie/	294
240	papers/masa-sentire-2011.pdf	295
241		296
242	• Catching the Drift: Probabilistic Content Models,	297
243	with Applications to Generation and Summariza-	298
244	tion	299
245	http://www.aclweb.org/anthology/N04-1015	300
246		301
247	8. Schedule	302
248	By November 16 th , complete:	303
249		304
250	• the web scraper to scrape product descriptions	305
251		306
252	• the static technique of the processing stage	307
253		308
254	• the annotation of sentences with SVM or Naive	309
255	Bayes	310
256		311
257	By November 22 nd , complete:	312
258		313
259	• the dynamic technique of the processing stage	314
260		315
261	By November 29 th , complete the post-processing	316
262	aggregation stage.	317
263		318
264	By December 4 th , the due date of the poster,	319
265	complete:	320
266		321
267	• ranking of sentiments in terms of relevance to con-	322
268	sumer requirements	323
269		324
270	• charts and tables of data	325
271		326
272	• begin writing the final report	327
273		328
274	Final report due: December 10 th	329