

# Automatic Labelling of Customer Reviews

By

Salomey Afua Addo (salomey@aims.edu.gh)  
African Institute for Mathematical Sciences (AIMS) - Ghana

Supervised by: Professor Axel Ngonga  
Paderborn University, Germany

June 2018

*AN ESSAY PRESENTED TO AIMS-GHANA IN PARTIAL FULFILMENT OF THE REQUIREMENTS FOR THE AWARD OF  
MASTER OF SCIENCE IN MATHEMATICAL SCIENCES*



# DECLARATION

This work was carried out at AIMS-Ghana in partial fulfilment of the requirements for a Master of Science Degree.

I hereby declare that except where due acknowledgement is made, this work has never been presented wholly or in part for the award of a degree at AIMS-Ghana or any other University.

Scan your signature

Student: Salomey Afua Addo

Scan your signature

Supervisor: Prof. Axel Ngonga

# 14 ACKNOWLEDGEMENTS

15 I thank God almighty for giving me life and wisdom to pursue mathematics. Much appreciation  
16 goes to Prof. Axel Ngonga, my supervisor, Dr. Prince Koree Osei, my mentor and AIMS-Ghana  
17 president and Ms. Abla Azalekor, my tutor.

# 18 DEDICATION

19 I dedicate this work to all machine learning lovers. We hope to see Artificial Intelligence gain  
20 root in Africa.

21

# Abstract

22

23

24

25

26

27

28

29

30

31

32

33

The growth of Information, Communication and Technology (ICT) coupled with the emergence of Internet as a medium of communication has led to the proliferation of Electronic Commerce. Customers are given the opportunity to express their shopping experiences in the form of reviews or star ratings after shopping online. These feedback from customers might help improve products to suit market demands. However, mining customer feedback is an arduous task because there exists a large amount of them and they are unstructured in nature (i.e., natural language). Therefore, in this work we employed Natural Language Processing (NLP) tools such as word embeddings and unsupervised learning techniques such as clustering algorithms to build a pipeline that helps obtain summary words (labels) for customer reviews. Labels obtained for customer reviews would help online market owners like Amazon to uncover the products that customers like or dislike, producers would improve their product quality and customers would make informed purchasing decisions.

34

**keywords:** ICT, Electronic Commerce, NLP, word embeddings

35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58

# Contents

<b>Declaration</b>	<b>i</b>
<b>Acknowledgements</b>	<b>ii</b>
<b>Dedication</b>	<b>iii</b>
<b>Abstract</b>	<b>iv</b>
<b>1 Introduction</b>	<b>1</b>
1.1 Motivation . . . . .	1
1.2 Approach . . . . .	2
<b>2 Preliminaries</b>	<b>4</b>
2.1 Word Representations . . . . .	4
2.2 State-of-the-art Word Representations . . . . .	5
<b>3 Approach</b>	<b>13</b>
3.1 Clustering . . . . .	13
3.2 Pointwise Mutual Information (PMI) . . . . .	18
3.3 Word-Cluster Matrix . . . . .	19
<b>4 Implementation</b>	<b>20</b>
4.1 Read Amazon data . . . . .	20
4.2 Extracting positive and negative reviews based on star ratings . . . . .	20
4.3 Pre-processing of text . . . . .	21
4.4 Vectorising the reviews using Embeddings from Language Modelling (ELMo) . . . . .	22
4.5 Clustering review vectors using DBSCAN . . . . .	22
4.6 Extracting the words in the reviews in the clusters . . . . .	24
4.7 PPMI Matrix . . . . .	26
4.8 Sort words in the labels . . . . .	26

59	<b>5 Results and Discussion</b>	<b>29</b>
60	5.1 First setting . . . . .	29
61	<b>References</b>	<b>31</b>
62	<b>APPENDICES</b>	<b>32</b>
63	5.2 Python codes for the entire pipeline for our approach . . . . .	32
64	5.3 Cluster images for both positive and negative reviews . . . . .	38
65	5.4 Labels for the first six selected clusters . . . . .	39

66 

# List of Figures

67	1.1 A customer review on a dress on Amazon . . . . .	1
68	1.2 A customer review on a phone on Amazon . . . . .	1
69	1.3 The pipeline for our algorithm . . . . .	3
70	4.1 Clustering results for 20 positive reviews . . . . .	23
71	4.2 Clustering results for 20 negative reviews . . . . .	24
72	4.3 word-cluster count matrix for 20 negative reviews . . . . .	27
73	4.4 PPMI matrix for 20 negative reviews . . . . .	27
74	5.1 Clustering results for 70 positive reviews . . . . .	38
75	5.2 Clustering results for 70 negative reviews . . . . .	39



# 76 List of Tables

77	2.1 A co-occurrence matrix for the vocabulary $V$ . . . . .	9
----	---	---

# 1. Introduction

## 1.1 Motivation

Electronic Commerce, commonly known as e-commerce, involves buying and selling of products or rendering services over computer networks, such as the Internet [1]. Online customer reviews may be defined as product evaluations posted by customers on a company's or a retailer's website [2]. For instance, on Amazon<sup>1</sup>, the online retailer and cloud computing corporation's website, customers are offered the opportunity to express their shopping experiences in the form of numerical star ratings ranging from 1 to 5 stars and open-ended customer authored comments about the products. Consequently, there exists a wealth of customer reviews online. The Amazon corporation has a collection of more than 130 million customer reviews [3]. These reviews contain important information that could be useful to stakeholders such as producers and customers. For instance, a producer could improve upon his product design based on the comments customers make about the product. Also customers make informed purchasing decisions based on the reviews made by other customers about the product they intend buying. Below are examples of customer reviews on Amazon Marketplace.

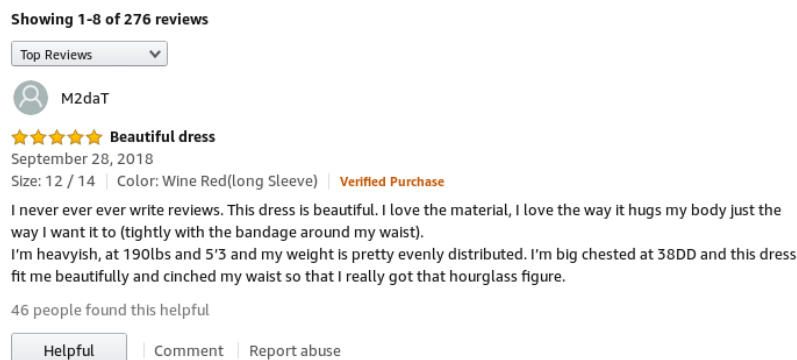


Figure 1.1: A customer review on a dress on Amazon

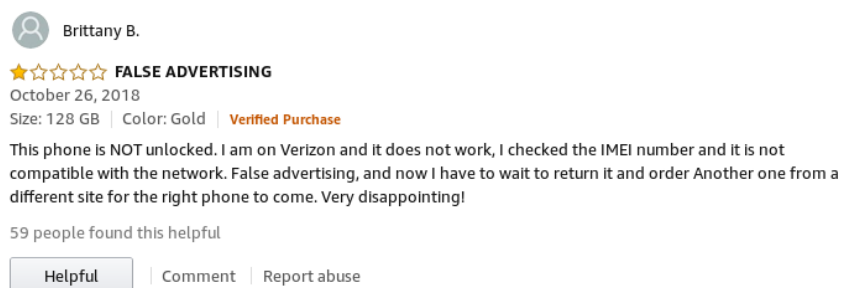


Figure 1.2: A customer review on a phone on Amazon

Given the enormous amount of customer reviews available, it is tedious and time consuming to

---

<sup>1</sup><https://www.amazon.com/>

manually inspect and analyse these reviews. Due to these constraints, producers and customers are unable to accrue the benefits from the enormous body of customer reviews available online. To address this issue, in this thesis we present an approach that can automatically assign labels to clusters of customer reviews. This approach helps the producers and consumers to have a summary or overview of the customers feedback without manually analysing the vast pool of reviews online.

## 1.2 Approach

Our goal is to build a pipeline that automatically labels customer reviews. Firstly, we partition the customer reviews into two groups based on the star rating given for each review. Here star rating less than or equal to 2 is considered as negative, while any review with star rating greater than 2 is considered as a positive review. We assume that the first group contains all the negative reviews and the second group contains all the positive reviews. We present an example of negative and positive reviews below.

	star_rating	review_body
2	Was soon looking forward to it...	
2	Honestly I didn't care for it...	
1	Movie is great - problem was...	
2	Poor quality video...	
1	A financial loss for me ...	
1	I hate powerpuff girls because ...	
1	Perhaps the saddest thing is...	

	star_rating	review_body
5	Fast shipping. Pleasure to deal...	
5	Kevin Kline is the most versatile,	
4	great movie	
3	Y'know what this reminded me of?	
5	Wonderfully accurate portrayal...	
4	In great shape for how old it is!	

Secondly, we pre-process the words in each review by tokenizing the words in the corpus, removing non-alphabetic characters and punctuation marks such as "?", ",", and "!". Also, we remove frequently occurring words which carry low information (stop words) such as "and", "or" and "the".

Thirdly, we use Embeddings from language model (ELMo) to convert reviews into vectors. We will provide further details and a definition for Embeddings from Language Model in the next chapter.

Fourthly, we perform thematic clustering on each set of review vectors (vectors for positive and negative reviews) using Density Based Spatial Clustering of Applications with Noise (DBSCAN).

Lastly, we find the Pointwise Mutual Information (PMI) scores between words and the clusters obtained. This approach helps us to determine the most informative words in a cluster. The most informative words are used as labels for each cluster.

We present the algorithm for our pipeline below:

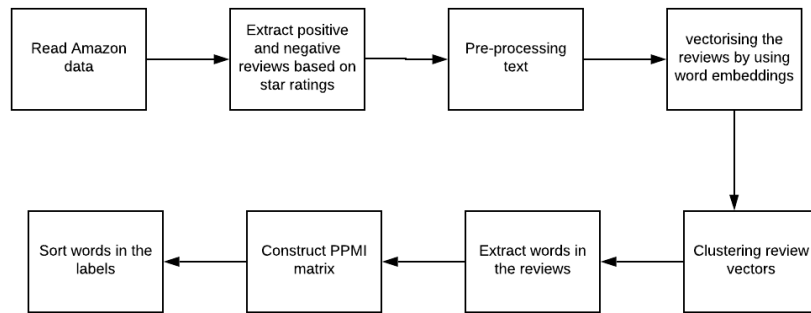


Figure 1.3: The pipeline for our algorithm

1. read the tsv file containing the Amazon reviews,
2. extract both positive and negative reviews based on the star ratings given by the customers. Put each review category into a list,
3. write a function that pre-processes each review by removing non-alphabetic characters and stop words, and returns a list containing lists of tokenised reviews,
4. create a for loop that picks each list of tokenised reviews from the bigger list, and converts each review into a vector using embeddings. Particularly, we used ELMo in this thesis,
5. run a clustering algorithm on the vectors. We used DBSCAN for our clustering,
6. extract the words in each cluster,
7. compute Positive Pointwise Mutual Information (PPMI) between words and clusters to determine the most informative words in a cluster.
8. Sort the labels based on their PPMI scores.

This thesis is organized as follows: Section 2 provides the details of the state-of-art approaches for converting natural language in the form of words into vectors. Section 3 presents clustering analysis and PPMI matrix. Section 4 presents a detailed implementation of our pipeline. Section 5 presents the discussion of the results obtained for the work.

## 2. Preliminaries

Language modelling is a task pivotal to the study of Natural Language Processing (NLP). It helps carry out tasks such as question answering [4], textual entailment [5] and semantic role labelling [6] among others. The objective of Language modelling is to train language models to learn the likelihood of occurrence of words based on previous sequences of words used in text. Language models can be trained on character level, word level, sentence level or paragraph level [7]. Computers do not have the ability to understand human natural language, hence to be able to create a platform for computers to explore the knowledge in human natural language, we encode words in a numerical or vectorial form. This form of representing words is called word representations or word embeddings.

Before we proceed to how word representations have evolved over time, let us clarify two ways of referring to words in a text. The first way of referring to words in a text, is by the term "word token" (see e.g. 2.0.0.1) and the second way is by the term "word type" (see e.g. 2.0.0.1). A word token refers to the total number of word units in a text or corpus irrespective of how often they may be repeated. While word type refers to the number of distinct words in a text or corpus. For instance, consider the following example of a text:

**Example 2.0.0.1.** As shades of colors tell the stories of life; we tell the african stories from an african perspective.

In Example 2.0.0.1, the word "tell" appeared twice, and the same thing happens for the words "african" and "stories". Hence, there are 18 word tokens and 13 word types in this particular example of text.

Now let us proceed to how word representation has evolved over time.

### 2.1 Word Representations

Hitherto, word types in a text were represented by integers, thus this approach of representing word types is referred to as discrete representations [7]. In the discrete representation approach, non-negative integers are used to represent word types in a text. One advantage of discrete representations is that the set of non-negative integers has an infinite number of elements and hence the body of words used in a particular language could be continuously expanded as we encounter new word types [7]. The assignment of integers to word types could be in an alphabetical order or based on the order of occurrence of the word type [7].

However, one limitation of discrete representation is that for instance, two word types with similar meanings might be assigned integers whose difference is huge. When these two integers are plotted on a number line, they may appear to be far apart although the words may be semantically related.

To perform complex NLP tasks, the computer should be able to determine the sense of a word and the similarity among words in a text corpus [7]. Therefore, words are transformed into vectors in

order to capture word sense in each of the dimensions of the vector [7]. For example, a collection of word types can be referred to as a class, (e.g., colours) and we can assign 1 to a specific dimension of the vector if the word type belongs to the word class, or 0 otherwise.

Moreover, we can use dimensions to capture the sense of magnitude. For example, we may assign 1,000,000 to the word "sun" and 1 to the word "atom". Each of the dimensions of a vector associated to a word type is called a feature. Features can be derived from a word's vector by use of automated algorithms or they could be designed by experts [7].

One way of representing word types as vectors is by the one-hot encoding vectors. This method involves assigning a binary value (0 or 1) to each of the dimensions of word vectors. In the one-hot encoding approach each word type has its own dimension. We assign 1 to the dimension of the word type and assign 0 to the other remaining dimensions of the one-hot encoding vector [7]. For the sake of illustration let us consider the following text comprised of two sentences:

**Example 2.1.0.1.** Ama is a beautiful girl. Ama is a pretty girl.

The vocabulary set denoted as  $V$  for this text is  $V = \{\text{Ama, is, a, beautiful, pretty, girl.}\}$  Now let us create a one-hot encoded vector for each word type in the vocabulary,  $V$ .

**Example 2.1.0.2.**  $\text{Ama}=[1,0,0,0,0,0]$ ,  $\text{is}=[0,1,0,0,0,0]$ ,  $\text{a}=[0,0,1,0,0,0]$ ,  $\text{beautiful}=[0,0,0,1,0,0]$ ,  $\text{pretty}=[0,0,0,0,1,0]$ ,  $\text{girl}=[0,0,0,0,0,1]$ .

The length of the one-hot encoded vector in Example 2.1.0.2 is equal to the size of the vocabulary  $V$ . Hence the vectors we obtained in Example 2.1.0.2 lie in a 6 dimensional vector space where each word type lies in each dimension. When we take the dot product of any two distinct vectors of the basis, we obtain zero, e.g.

$$\text{beautiful} \cdot \text{pretty} = 0 \quad (2.1.1)$$

and

$$\text{Ama} \cdot \text{is} = 0. \quad (2.1.2)$$

According to one-hot encoding technique, the Equation 2.1.1 implies that the word types "beautiful" and "pretty" have nothing in common just as in Equation 2.1.2 the word types "Ama" and "is" have nothing in common. In reality, this is false.

Another limitation is that, one-hot encoded vectors are highly dimensional, sparse and they could be computationally expensive especially when the vocabulary size increases [7]. Therefore, one-hot encoded vectors are analogous to integer-based representation—they don't convey semantic relationship among words [7]. Hence, they are not useful in performing NLP tasks. In view of this we consider some state-of-the-art approach of producing word embeddings that carry semantic meanings and are greatly useful in various NLP tasks.

## 2.2 State-of-the-art Word Representations

Currently, the state-of-the-art approach to producing word embeddings is based on neural networks. According to Patterson and Gibson [8], "neural network is a computational model that

shares some properties with the animal brain in which many simple units are working in parallel with no centralized control unit". Mathematically, neural networks are mathematical models defining a function  $f : X \mapsto Y$ . The state-of-the-art approach to producing word embeddings such as Word2Vec [9], GloVe [10] and ELMo [11] use neural networks to find the vector mappings for each word in a given corpus.

**2.2.1 Word2Vec.** Word2Vec introduced by Mikolov et al. [9] in 2013, is a model used to produce word embeddings. The main idea behind it is that words with similar context (neighbouring words) have similar meaning. This confirms John Firth, an English linguist's quote, "You shall know a word by the company it keeps" [12]. The Word2Vec model is a two layer neural network which takes a corpus of text as its input and produces a high dimensional vector as output, with each word type in the corpus being represented by a vector. Word2Vec produces vectors such that words with similar context in the corpus are close to each other in the vector space [9]. Let us illustrate a bit more on Word2Vec with the next example.

**Example 2.2.1.1.** Consider the following text: the capital city of Ghana is Accra and the capital city of Togo is Lome.

In the vector space produced by Word2Vec, the difference between "Ghana" and "Accra" is the same as the difference between "Togo" and "Lome". Also, the vectors for "Ghana" and "Togo" would have close proximity whereas the vectors for "Ghana" and "onions" would be far apart in the vector space. Amazingly, mathematical operations could be applied to the word vectors to yield reasonable results [9]. For instance,

$$Ghana - Accra + Togo = Lome. \quad (2.2.1)$$

There are two approaches or algorithms to convert words to vectors using Word2Vec. They are the skip-gram based model approach and the continuous bag of words approach.

To start with, the skip-gram model approach seeks to predict the surrounding words of a target word in a corpus, given a fixed window size (the radius of words surrounding the target word to the right and left). By illustration let us consider the sentence

**Example 2.2.1.2.** Ghana's chocolate is the best in the world.

Let the word "best" be the target word. Suppose we fix a window size of 3, the context words to the left is "the", "is" and "Ghana's" while the context words to the right is "in", "the" and "world". The continuous bag of words approach is the reverse of the skip-gram based model approach. It predicts the target word given a context.

We would like to proceed by examining the mathematics behind Word2Vec. The skip-gram based model approach predicts the words of each context independently given the target word. The objective of Word2Vec is to maximise the probability of the context given the target word across the corpus [9]. The probability of predicting a context given its target word across the corpus can be mathematically expressed as follows [9]:

$$\prod_{t=1}^T \prod_{-m \leq j \leq m, j \neq 0} p(w_{t+j} \mid w_t; \theta). \quad (2.2.2)$$

Where

- $\theta$  denotes the hyper-parameters of the model (word representations of the target word),
- $m$  stands for the window size (radius around the target word),
- $T$  is the size of the corpus,
- $j$  denotes the index of the words that constitute the context of the target word  $w_t$  for a given  $t$ ,
- $t$  denotes the index of a target words,
- $w_t$  denotes the target word given the index  $t$ , where  $t = 1, \dots, T$ .

After applying the logarithm function to the expression 2.2.2, we obtain

$$\sum_{t=1}^T \sum_{-m \leq j \leq m, j \neq 0} \log p(w_{t+j} | w_t; \theta). \quad (2.2.3)$$

Let us take the average of the terms in the mathematical Expression 2.2.3 so that instead of taking the probability of the whole corpus we could take the average probability, that is probability of each word. Let us minimise the objective function by introducing a negative sign

$$J(\theta) = -\frac{1}{T} \sum_{t=1}^T \sum_{-m \leq j \leq m, j \neq 0} \log p(w_{t+j} | w_t; \theta). \quad (2.2.4)$$

Equation 2.2.4 now becomes the objective function we seek to minimise.

Similarly, for the continuous bag of words approach, the context words are summed in order to predict the target word. In this case, the objective function of the corresponding minimisation problem that underpins Word2Vec is given by:

$$Z(\theta) = -\frac{1}{T} \sum_{t=1}^T \log p(w_t | \sum_{-m \leq j \leq m, j \neq 0} w_{t+j}; \theta). \quad (2.2.5)$$

In the Word2Vec model, each word type has two vector representations: the context embedding vector and the target embedding vector [9]. The probability distribution is defined by the softmax function [8]. The probability distribution for skip-gram based model is

$$p(w_c | w_t) = \frac{\exp(v_{w_c}^T u_{w_t})}{\sum_{w=1}^W \exp(v_w^T u_{w_t})}. \quad (2.2.6)$$

Also, the probability distribution for continuous bag of words is

$$p(w_t | w_c) = \frac{\exp(u_{w_t}^T v_{w_c})}{\sum_{w=1}^W \exp(u_{w_t}^T v_w)}. \quad (2.2.7)$$



Where  $u_{w_t}$  represents the word vectors for the target words and  $v_{w_c}$  represents the word vectors for the context words.

The softmax function would always yield a positive real number since the numerator is made up of an exponential function and the denominator is made up of the sum of exponential functions. Moreover, the softmax function would yield a real number between 0 and 1 due to the normalization factor in the denominator.

In a nutshell, Word2Vec is the state-of-the-art approach for building word vectors. It scans each word in the corpus, predicts the context of each word and captures co-occurrence of words one at a time. This approach is quite inefficient supposing our corpus is really large. It is suggested that an optimal way is to scan through the entire corpus once and count the frequency of co-occurring words [10]. Also, instead of making updates on samples one at a time, it is more efficient to make one update step that captures the entire word count. To tackle these limitations, we introduce the next state-of-the-art approach for producing word vectors called Global Vectors for Word Representation (GloVe) [10].

**2.2.2 GloVe.** There are two main approaches to produce word representations. These are:

1. matrix factorization method,
2. local window context method.

The matrix factorization method seeks to decompose large matrices that contain statistical information about the corpus, (e.g., window based co-occurrence matrix), by utilizing low rank approximations such as the Singular Value Decomposition (SVD) [10].

The rows and the columns of the co-occurrence matrix corresponds to word types in the corpus and the entries corresponds to the frequency of a given word as it occurs in the context of another word.

Let us take an example of a co-occurrence matrix. Given the three sentences below:

1. I like potatoes.
2. I enjoy pizza.
3. I like reggae.

The vocabulary set  $V$  becomes

$$V = \{I, \text{like}, \text{enjoy}, \text{potatoes}, \text{pizza}, \text{reggae}, .\} \quad (2.2.8)$$

From Table 2.1 we could see that "like" and "enjoy" co-occur with the word "I". This gives a clue about their similarity. Increasing the size of the corpus would enable us appreciate these similarity relationships much better.

The rows of the co-occurrence matrix could be used to form word vectors. However, there are limitations to the use of a co-occurrence matrix. Firstly, the matrix increases in size as we

count	I	like	enjoy	potatoes	pizza	reggae	.
I	0	2	0	0	0	0	0
like	2	0	0	1	0	1	0
enjoy	1	0	0	0	1	0	0
potatoes	0	1	0	0	0	0	1
pizza	0	0	1	0	0	0	1
reggae	0	1	0	0	0	0	1
.	0	0	0	1	1	1	0

Table 2.1: A co-occurrence matrix for the vocabulary V

increase the vocabulary. Hence, the word vectors generated would have many dimensions, thus requiring lots of storage space and computational cost may increase. Also, one may encounter sparsity issues when using these vectors to train machine learning models, such as classification models. To solve this problem we may employ a dimensionality reduction method such as Principal Component Analysis (PCA) [13].

One main limitation of the matrix factorisation approach is that it gives disproportionate importance to frequently occurring words in the corpus, for example words like "the", "is" and "and" among others. This gives the wrong impression that these words (they are also called stop words) share common semantic relationships with the words they co-occur with.

On the other hand, window-based methods seek to learn word representations that enables us to predict words within local context windows [10]. An example of this approach is the skip-gram based model, which we have explained earlier under Word2Vec approach.

Unlike the matrix factorisation approach for learning word vectors, window-based methods are deficient in the information carried by the statistics of co-occurring words of the corpus, therefore it fails to leverage the rich amount of repetitions in the data.

In summary, the matrix factorisation approach complements the window-based approach. The former is fast to train while the latter is slow to train. Also, matrix factorisation makes good use of the statistics of co-occurrence of words in the corpus thereby capturing word similarity well, although it gives undue importance to some frequently occurring words and it performs poorly on specific tasks such as question answering [4]. On the other hand, the window-based approach makes inefficient use of co-occurrence statistics but performs well in specific tasks such as named entity recognition and word analogies. Besides similarity relationships, the window based approach also captures syntactic relationships among words [10].

GloVe, introduced in 2014, was built to take advantage of the efficiencies of the two main approaches for learning word vectors.

The objective function of GloVe is of the form,

$$J(\theta) = \frac{1}{2} \sum_{i,j=1}^W f(P_{ij}) (u_i^T v_j - \log P_{ij})^2. \quad (2.2.9)$$

Where

- $\theta$  denotes all the parameters of the model
- $P_{ij}$  denotes the co-occurrence matrix of the words in the corpus
- $W$  denotes the size of the vocabulary.

GloVe seeks to minimize the distance between the inner product ( $u_i^T v_j$ ) and the log count of these words,  $P_{ij}$ . Also,  $f$  allows us to weight the frequent counts of co-occurrences, that is we try to suppress the count of frequent words (stop words) in the objective function.

**2.2.3 Embeddings from Language Models (ELMo).** So far we have examined two major state-of-the-art approaches for generating word vectors- Word2Vec and GloVe. These word embeddings represent word types using a vector. However, they do not capture the polysemous (words with more than one meaning) nature of some words in natural language.

Let us consider the following sentences:

1. He fans the fire using a book cover.
2. His fans love him although he is polygamous.

In the first sentence "fans" is a verb which describes the act of rekindling a fire by blowing air while in the second sentence "fans" refers to the followers of a celebrity. Although the word "fans" is spelt the same in both sentences, they have different meanings in different contexts. By using Word2Vec or GloVe it is highly likely that the word "fans" would be represented by the same vector in both sentences. Thus context-independent vectors do not sufficiently model natural language.

In view of this, we introduce a novel contextual word representation-Embeddings from Language Models (ELMo), which was introduced by Peters et al. [14] in 2018. ELMo is based on neural networks, specifically it is a function of all the internal layers of a bidirectional language model (biLM) [14], hence it is often called deep contextualised word representation. One difference between ELMo and the other context-independent word embeddings is that it produces word vectors for word tokens instead of word types.

Let us go deeper into language models. Language models seek to find the probability distribution of the next word in a sequence given the sequence of the preceding words. A language model is trained by reading sentences in a corpus in the forward and backward direction. Thus we have forward language model and backward language model.

Given a sequence of  $N$  tokens  $(t_1, t_2, \dots, t_N)$  in a corpus, a forward language model seeks to predict the next word token given the sequence of previous words. While the backward language model is the reverse, given the future context it seeks to find the probability of the previous word token. Mathematically, we express the forward language model as

$$p(t_1, t_2, \dots, t_N) = \prod_{k=1}^N p(t_k \mid t_1, t_2, \dots, t_{k-1}) \quad (2.2.10)$$

while the backward language model is expressed as

$$p(t_1, t_2, \dots, t_N) = \prod_{k=1}^N p(t_k \mid t_{k+1}, t_{k+2}, \dots, t_N). \quad (2.2.11)$$

337 A bidirectional language model combines both forward and backward language models in predict-  
338 ing word tokens in a corpus.

Our objective in this model is to jointly maximize the log likelihood of both forward and backward language models. Hence we obtain

$$\sum_{k=1}^N \left( \log p(t_k \mid t_1, \dots, t_{k-1}, \theta_x, \vec{\theta}_{\text{LSTM}}, \theta_s) + \log p(t_k \mid t_{k+1}, \dots, t_N, \theta_x, \overleftarrow{\theta}_{\text{LSTM}}, \theta_s) \right). \quad (2.2.12)$$

339 Where

- 340 •  $\theta_x$  denotes the parameter for token representation
- 341 •  $\theta_s$  denotes the parameter for the softmax layer
- 342 •  $\vec{\theta}_{\text{LSTM}}$  denotes the parameter for forward language model
- 343 •  $\overleftarrow{\theta}_{\text{LSTM}}$  denotes the parameter for backward language model.

344 Both the forward and backward directions share the same representation parameter and softmax  
345 parameter while we keep separate parameters for each of the Long Short Term Memory (LSTM)  
346 for each direction.[14]

Given  $L$  number of stacked LSTMs, for each token we train  $L$ -layer LSTM forward and backward language model, thus we generate  $2 \times L + 1$  different vector representations. Therefore, the word embedding for input token  $k$  is

$$R_k = \{x_k^{LM}, \vec{h}_{k,j}^{LM}, \overleftarrow{h}_{k,j}^{LM} \mid j = 1, \dots, L\} \quad (2.2.13)$$

$$= \{h_{k,j}^{LM} \mid j = 0, \dots, L\}. \quad (2.2.14)$$

347 Here  $x_k^{LM} = h_{k,0}^{LM}$  is a context-independent token representation. We pass the context-independent  
348 vector through the  $L$ - layers of forward and backward LSTMs. Each LSTM layer returns a  
349 context-dependent representation  $\vec{h}_{k,j}^{LM}, \overleftarrow{h}_{k,j}^{LM}$  at each word token position  $k$  [14].

350 The combination of the intermediate layer representations in the biLM is used to perform specific  
351 tasks such as sentiment analysis [4], named entity recognition [11], co-reference resolution [4],  
352 among others [14].

To obtain an ELMo word representations customised or tailored for a specific task we find the linear combination of the  $h_{k,j}^{LM}$  and some parameters. Hence we have

$$ELMo_k^{\text{task}} = \gamma^{\text{task}} \sum_{j=0}^L s_j^{\text{task}} h_{k,j}. \quad (2.2.15)$$

353 Where  $\gamma^{\text{task}}$  allows the task model to scale the entire ELMo vector and  $s_j^{\text{task}}$  is the softmax-  
354 normalized weight for each layer [14].

355 Not only is ELMo a contextual word vector, it models complex characteristics such as syntactic  
356 and semantic word use. It is shown that the higher-level LSTM states models the semantic  
357 aspects of words whereas the lower level LSTM captures the context-dependent aspect of word  
358 structure such as word syntax [14]. Also, ELMo succeeds in modelling polysemous language due  
359 to the fact that it returns unique vectors for each word token.

360 Unlike Word2Vec and GloVe which are more or less like lookup tables, ELMo word vectors include  
361 both word type vectors and neural network parameters [7]. This implies that word type vectors  
362 are not redundant however they could be concatenated with context-dependent embeddings to  
363 perform well on downstream tasks [14].

## 3. Approach

### 3.1 Clustering

Clustering involves grouping a set of objects in such a way that objects that belong to the same group tend to be more similar to each other compared to objects from other groups. Unlike classification where objects are grouped based on some pre-existing labels, clustering does not require any class labelling or pre-existing categorization. Clustering partitions data sets according to their intrinsic categorization. Hence it is referred to as an unsupervised technique for grouping data. A group of similar data points is called a cluster. Clusters are formed based on the similarity among data points. A similarity is an estimate of the strength of the relationship among data points [15]. In most clustering algorithms distance functions such as Euclidean distance, Cosine similarity, Minkowski distance among others are often employed as a measure of similarity [15]. Hence, data points that belong to the same cluster has a minimum distance value while data points across different clusters are a certain maximum distance apart.

When our dataset becomes large and the dimension of the data set increases, it gives rise to the following requirements [16] on cluster algorithms:

1. determination of the input parameter,
2. discovery of clusters with arbitrary shapes,
3. efficiency on large data base.

There are two traditional approaches to clustering algorithms [17]. These are,

1. hierarchical clustering algorithm,
2. partitioning clustering algorithm.

We have two types of hierarchical clustering approach, viz., agglomerative hierarchical clustering and divisive hierarchical clustering [17]. The agglomerative type of hierarchical clustering involves considering each data point as an individual cluster. We merge similar clusters at each iteration until we form  $k$  number of clusters,  $k \geq 1, k \in \mathbb{N}$ . However, the divisive hierarchical clustering algorithm involves considering the entire dataset as one mass cluster. At each iteration, non-similar data points are separated from the cluster. Here each separated data point is considered as a cluster on its own. One setback of the hierarchical clustering algorithm is the need to define the termination condition that indicates when to merge or separate clusters [16]. Examples of hierarchical clustering approach are, Balance Iterative Reducing and Clustering (BIRCH) and Cluster Using REpresentatives (CURE) [18].

Unlike hierarchical clustering, partitioning clustering algorithm involves grouping the data points into sets, each of which is a partition of the data and is comprised of  $k$  number of clusters. Each

partition or group on its own represents a cluster and must contain at least one object which belongs to exactly one partition [19]. This type of algorithm has no notion of outliers. There are four main types of partitioning clustering methods, viz. K-means Algorithm [19], K-medoids Algorithm [20] and Clustering for Large Applications based on RANdomized Search (CLARANS) [21]. Let us examine the method by K-means.

The K-means clustering algorithm was introduced by James Macqueen in 1967 [19]. The objective of the K-means is to obtain a minimum distance between the centroid or mean of the cluster and its neighbouring data points.. Thus the objective function used for the K-Means is the sum of deviations between a data point and its centroid using an appropriate distance function. Below is the mathematical expression of the objective function.

$$J = \sum_{j=1}^k \sum_{i=1}^n \|x_{ij} - c_j\|^2 \quad (3.1.1)$$

Where

- J is the objective function,
- k is the number of clusters,
- n is the number of data points,
- $x_{ij}$  is the  $i$ th data point in the  $j$ th cluster,
- $c_j$  is the mean for data points in cluster  $j$ ,
- $x \mapsto \|x\|$  is the Euclidean distance function.

For a given set of data points  $x_1, x_2, \dots, x_n$ , the centroid, which we denote by  $C$ , is defined as follows: The centroid  $C$ , is computed as

$$C = \frac{1}{n} \sum_{i=1}^n x_i. \quad (3.1.2)$$

Where  $C$  is the average of all the points in the cluster. The coordinates of the centroid of each cluster is the arithmetic mean for each dimension over all the points in the cluster.

Let us consider the algorithmic steps for K-means clustering. Let  $x = \{x_1, x_2, \dots, x_N\}$  be the set of data points and  $C = \{c_1, c_2, \dots, c_k\}$  be the set of centroids. We present the procedure for finding the optimal clusters using K-Means approach.

1. Randomly select " $k$ " data points as cluster centroids,
2. compute the distance between each data point and the clusters' centroids,
3. assign each data point to the nearest centroid, i.e., data point whose distance from the cluster center is minimum as compared to all other centroids,

4. recompute the new cluster centroid using the formula

$$c_j = \frac{1}{n_j} \sum_{i=1}^{n_j} x_i \quad (3.1.3)$$

where  $n_j$  is the number of data points in the  $j$ th cluster and  $x_1, \dots, x_{n_j}$  are the data points in this cluster.

5. recompute the distance between each data point and the newly obtained cluster centroids.

6. stop if no data point is reassigned to a new cluster. Otherwise, repeat 3, 4, and 5.

K-means is easy to understand and implement, computationally faster than hierarchical clustering methods and yields the best results when the data set is well separated from each other. However, its limitations are that they require a priori specification of the  $k$  number of centroids and it performs poorly when the data set is noisy. Besides, the algorithm fails when the data set is non-linear.

We introduce K-medoids which is another approach for partitioning clustering approach. We have already stated that K-means is sensitive towards outliers hence we introduce another partitioning clustering approach which is less sensitive to outliers. The differences between K-medoids and K-means is that while K-means takes the means of the data points as the representative of each cluster, K-medoids takes the medoid-which is the most centrally located object in the cluster, as the representative of the cluster. Also, in K-means the representative of a cluster may not necessarily belong to the data set, however each medoid used in k-medoids algorithm belongs to the data set. The objective for k-medoids is similar to that of K-means, in that we seek to minimise the distance between the points and medoids in each cluster. The strategy is to find  $k$  clusters given  $n$  data points by first randomly selecting  $k$  number of medoids for each cluster. The remaining data points are assigned to each cluster based on their closeness to the medoid of each cluster.

The objective function is expressed mathematically as

$$E = \sum_{j=1}^k \sum_{i=1}^n \|x_{ij} - m_j\|^2. \quad (3.1.4)$$

Where

- $E$  is the objective function,
- $k$  is the number of clusters,
- $n$  is the number of data points,
- $x_{ij}$  is the  $i$ th data point in the  $j$ th cluster,
- $m_j$  is the medoid for cluster  $j$ ,



- $x \mapsto \|x\|$  is the Euclidean distance function.

Let us consider the algorithmic steps for k-Medoids clustering. Let  $x = \{x_1, x_2, \dots, x_N\}$  be the set of data points and  $M = \{m_1, m_2, \dots, m_k\}$  be the set of medoids. We present the procedure for finding the optimal clusters using k-Medoids approach.

1. Randomly select " $k$ " cluster medoids,
2. compute the distance between each data point and cluster centroids,
3. assign each data point to the nearest medoid, i.e., data point whose distance from the cluster medoid is minimum as compared to all other medoids,
4. compute the cost, that is the sum of distances of points to their medoids,
5. for each medoid  $m$ , for each non-medoid data point  $o$ , swap  $m$  and  $o$ , re-assign each data point to the closest medoid,
6. recompute the sum of the distances of points to their medoids,
7. if the total cost increases in the previous step undo the swap,
8. if the total cost decreases, swap the medoid  $m$  with  $o$  to form the new set of  $k$ -Mediod,
9. repeat 5,6,7 and 8 until there is no change in the sum of the distances of points to their medoids.

The two main approaches to clustering algorithm-hierarchical and partitioning clustering methods fail to meet all the requirements when the database is large. Hence, we consider the density based approach which was explored by Jain in 1988 [22]. The density based notion of constructing clusters involves partitioning a dataset into non-overlapping cells and constructing a histogram based on the cells created. Here cluster centers corresponds to the cells with high frequency counts (high density) whereas areas of low density may be classified as "noise". This is a density based approach to distinguishing between clusters and noise.

**3.1.1 Density Based Spatial Clustering of Applications with Noise (DBSCAN).** This algorithm was introduced by Ester et al. [16]. It relies on the density based notion of clustering to discover the clusters and noise in a dataset. The main idea behind DBSCAN is that every point in the cluster must contain a minimum number of points in its neighbourhood determined by a specified radius. DBSCAN has two input parameters, performs better than other algorithms on a large dataset and discovers arbitrary cluster shapes in any dimensional space [16]. Before we explore how DBSCAN is implemented let us consider some definitions and lemmata pertaining to density based notion of clustering.

**Definition 3.1.1.1** (*Eps*-neighbourhood of a point). The *Eps*-neighbourhood of a point  $p$  denoted by  $N_{Eps(p)}$ , is defined by  $N_{Eps(p)} = \{q \in D : \text{dist}(p, q) \leq Eps\}$ [16].

Where  $p$  and  $q$  are points in the dataset  $D$ , and  $Eps$  is the radius of the neighbourhood with respect to the point  $p$  and  $dist()$  is an arbitrary distance measure.

There are two types of points in a cluster. These are;

- core points,
- border points.

A core point in a cluster has at least a minimum number of points denoted as  $MinPts$  in its  $Eps$ -neighbourhood while the border point has less than the minimum number of points in its  $Eps$ -neighbourhood. All neighbours in the  $Eps$ -strip of a core point is considered to belong to the core point's cluster. This notion is termed as direct density reachability. Let us take a formal definition.

**Definition 3.1.1.2** (directly density-reachable). A point  $p$  is directly density-reachable from a point  $q$  with respect to  $Eps$ ,  $MinPts$  if

- $p \in N_{Eps}(q)$  and
- $|N_{Eps}(q)| \geq MinPts$  (core point condition) [16].

If any of the points, say  $q$  in the neighbourhood of a core point, say  $p$ , is a core point, then the neighbourhood of  $q$  is included in the cluster formed by  $p$ . This notion is termed as density-reachability.

**Definition 3.1.1.3** (density-reachable). A point  $p$  is density reachable from a point  $q$  with respect to  $Eps$  and  $MinPts$  if there is a chain of points  $p_1, \dots, p_n$ ,  $p_1 = q$ ,  $p_n = p$  such that  $p_{i+1}$  is directly density-reachable from  $p_i$  [16].

All points within a cluster is said to be density-connected. Let us take a formal definition of density-connectedness.

**Definition 3.1.1.4** (density-connected). A point  $p$  is density connected to a point  $q$  with respect to  $Eps$  and  $MinPts$  if there is a point  $o$  such that both,  $p$  and  $q$  are density-reachable from  $o$  [16].

Let us consider a formal definition of a cluster.

**Definition 3.1.1.5** (cluster). Let  $D$  be a dataset of points. A cluster  $C$  with respect to  $Eps$  and  $Minpts$  is a non-empty subset of  $D$  satisfying the following conditions:

- $\forall p, q : \text{if } p \in C \text{ and } q \text{ is density reachable from } p \text{ with respect to } Eps \text{ and } MinPts, \text{ then } q \in C.$
- $\forall p, q : p \text{ is density-connected to } q$  [16].

Any point that does not belong to a cluster is considered as noise. Let us take a formal definition of a noise

**Definition 3.1.1.6** (noise). Let  $C_1, \dots, C_k$  be the clusters of the dataset  $D$  with respect to parameters  $Eps_i$  and  $Minpts_i$ ,  $i = 1, \dots, k$ . Then we define the noise as the set of points in the dataset  $D$  not belonging to any cluster  $C_i$ , i.e.  $\text{noise} = \{p \in D : \forall i : p \notin C_i\}$  [16].

## 3.2 Pointwise Mutual Information (PMI)

According to Bouma [23], "Pointwise Mutual Information is a measure of how much the actual probability of a particular co-occurrence of events  $p(x, y)$  differs from what we would expect it to be on the basis of the probabilities of the individual events and the assumption of independence  $p(x)p(y)$ ". The formula for computing PMI for the random variables  $X, Y$  given a pair of outcomes  $x$  and  $y$  is defined as:

$$PMI(x, y) = \log_2 \frac{p(x, y)}{p(x)p(y)} \quad (3.2.1)$$

where  $p(x)$  and  $p(y)$  are marginal probabilities and  $p(x, y)$  is the joint probability. The numerator of Equation (3.2.1) expresses how often  $x$  co-occurs with  $y$  and the denominator expresses the probability of  $x$  and  $y$  occurring together, given they occurred independently [24]. Pointwise mutual information is one of the efficient ways to estimate the association between words [25]. It estimates how much more two words co-occur in our word corpus than we would have expected them to occur a priori by chance. [24].

The formula for PMI between a target word  $w$  and a context word  $c$  is given as [25]

$$PMI(w, c) = \log_2 \frac{p(w, c)}{p(w)p(c)}.$$

The values of PMI are elements of  $\mathbb{R}$ , that is it assumes both positive and negative values. If PMI is negative it implies that there is a weak association between the target and context words. Also, if PMI value is positive it implies that there is a strong association between words. PMI value being 0 implies that the target word and context word are independent, that is,  $p(w, c) = p(w)p(c)$ ,  $\log_2(1) = 0$ .

Negative PMI values could be misleading if the size of our corpus is not enormous [24]. In view of this, let us replace all negative PMI values with 0. This approach is called Positive Pointwise Mutual Information (PPMI) [25]. The formula for PPMI is given as

$$PPMI(w, c) = \max \left( \log_2 \frac{p(w, c)}{p(w)p(c)}, 0 \right). \quad (3.2.2)$$

In our approach, to find the PPMI of each word in a cluster we need to first compute a word-cluster matrix.

### 3.3 Word-Cluster Matrix

The frequency of the occurrence of  $t$  words in the entire document is represented by a word-cluster matrix. The word-cluster matrix constitutes  $n$  clusters and  $t$  rows. Each column represents a cluster of reviews obtained from DBSCAN and each row represents the unique words in the entire document. Each entry  $a_{ij}$  of the matrix,  $A$ , represents the frequency count of the  $i$ th word in the  $j$ th cluster. The PPMI of each entry in the matrix can be computed as:

$$ppmi_{ij} = \begin{cases} P_{ij} & \text{if } P_{ij} > 0 \\ 0 & \text{otherwise} \end{cases} \quad (3.3.1)$$

where

$$P_{ij} = \log \left( \frac{p_w}{p_u p_s} \right), \quad (3.3.2)$$

$$p_w = \frac{a_{ij}}{\sum_{i=1}^t \sum_{j=1}^n a_{ij}}, \quad (3.3.3)$$

$$p_u = \frac{\sum_{j=1}^n a_{ij}}{\sum_{i=1}^t \sum_{j=1}^n a_{ij}}, \quad (3.3.4)$$

$$p_s = \frac{\sum_{i=1}^t a_{ij}}{\sum_{i=1}^t \sum_{j=1}^n a_{ij}}. \quad (3.3.5)$$

Here,  $p_w$  is the probability that the word  $w_i$  occurs in the cluster  $j$  with respect to the entire document,  $p_u$  is the probability of the word  $w_i$  occurring in the entire document and  $p_s$  is the probability of a cluster  $s_j$  in the entire document.

## 4. Implementation

We present the implementation of our approach for obtaining labels automatically given any customer reviews.

### 4.1 Read Amazon data

We first read the tab separated text file, without quote or escape characters containing customers' details on Amazon. This data set contains columns for marketplace, customer\_id, review\_id, product\_id, product\_parent, product\_title, product\_category, star\_rating, helpful\_votes, total\_votes, vine, verified\_purchase, review\_headline, review\_body and review\_date<sup>1</sup>. Find below an example of how to read a tsv file in Python.

```
obs=[]
filename = "amazon_reviews_us_Video_v1_00.tsv"
file = open(filename, mode='r')
for line in file.readlines():
    obs.append(line.split('\t'))
```

### 4.2 Extracting positive and negative reviews based on star ratings

Next, we created a data frame for the dataset using Pandas<sup>2</sup>. We extracted the negative review\_body (the review text) and positive review\_body based on the star ratings given by the customers for each product purchased. Customer reviews with star ratings less than or equal to two are grouped as negative reviews while reviews with star ratings greater than two are grouped as positive reviews. The reason why we first grouped the positive and negative reviews is to help stakeholders such as producers to easily identify products that are successful on Amazon or products that do not meet the customers' satisfaction. This approach helps the producers to improve on their products in order to meet market demands. We present a sample of negative reviews extracted from Amazon customer reviews on musical instrument.

star_rating	review_body
2	Note— Does not Fit Epiphone SG Special! Mine i...
2	Bridge pickup was broken. I replace d the pick...
1	I was hoping it would work well, but tried a s...
1	It works well while it can. Mine failed when l...

<sup>1</sup><https://s3.amazonaws.com/amazon-reviews-pds/tsv/index.txt>

<sup>2</sup><https://pandas.pydata.org/pandas-docs/stable/>

```

563 1      Really bad. Bought as a midi trigger kit but t...
564 1      sorry but wasnt made right
565 2      Novelty item. Fun. The delay on it is bad enou...
566 2      did not fit the Asus
567 1      The company accidentally sent me an alto inste...
568 2      Makes lots of Noise while unit has nice feature...

```

569 We also present 10 samples of positive reviews extracted from the Amazon reviews on musical  
 570 instrument.

```

571 star_rating      review_body
572 3      Works very good, but induces ALOT of noise.
573 5      Nice headphones at a reasonable price.
574 3      removes dust. does not clean
575 5      I purchase these for a friend in return for pl...
576 5      This is an awesome mic!
577 5      Used to cool equipment inside credenzas. Work...
578 5      Well built Ukulele! My daughter loves it!
579 5      Had to replace a new light after a lightning s...
580 5      I've owned multiple fixed boards over the year...
581 5      Consistent quality, no strong odours, works we...

```

## 582 4.3 Pre-processing of text

583 To be able to perform any NLP task on the negative and positive reviews we first had to pre-  
 584 process the reviews. We tokenised sentences in the documents into samples  $S_i, i = 1, \dots, n$ ,  
 585 where  $n$  is the total number of reviews in either the positive or negative reviews. Next, we tokenise  
 586 each word in  $S_i$  and converted all words in  $S_i$  into lower case. We also removed all non-alphabetic  
 587 characters such as "#\*&" and punctuation marks such as "?!,." from our document by using  
 588 python module regex. Also, we removed words which are commonly used such as "and", "of",  
 589 "in", "the", "a" and "an". These words are known as stop words. Our goal is to use statistical  
 590 analysis to extract words that best describes customers' reviews, for this reason we remove words  
 591 and characters that are statistically insignificant for our work. We present below the reviews  
 592 before text pre-processing.

```

593 ["Sounds okay, but I had numerous problems with this microphone.
594 <br /><br />First, sections were just loose to begin with. Fortunately,
595 I found I could tighten the bottom section and that improved.
596 <br /><br />Next, both controls (the On/Off switch and the Volume control)
597 fell off within in the first few minutes of working with them.
598 To clarify, the caps or knobs fell off. The switch itself was still intact.
599 It felt very cheap.
600 <br /><br />Next, while it sounded good at first, I suddenly started
601 getting some loud 60Hz hum. This was isolated to where the USB cord

```

```

602 connected to the bottom of the mic.
603 Reconnecting everything firmly
604 reduced the hum but I don't know for how long.]

```

605 Let us consider the text after pre-processing.

```

606 [['sounds', 'okay', 'numerous', 'problems',
607   'microphone', 'br', 'br', 'first',
608   'sections', 'loose', 'begin', 'fortunately',
609   'found', 'could', 'tighten', 'bottom',
610   'section', 'improved', 'br', 'br',
611   'next', 'controls', 'switch', 'volume',
612   'control', 'fell', 'within', 'first',
613   'minutes', 'working', 'clarify', 'caps',
614   'knobs', 'fell', 'switch', 'still',
615   'intact', 'caps', 'knobs', 'tiny',
616   'rubber', 'like', 'pieces', 'could',
617   'easily', 'lost', 'felt', 'cheap',
618   'br', 'br', 'next', 'sounded',
619   'good', 'first', 'suddenly', 'started',
620   'getting', 'loud', 'hz', 'hum',
621   'isolated', 'usb', 'cord', 'connected',
622   'bottom', 'mic', 'reconnecting', 'everything',
623   'firmly', 'reduced', 'hum', 'know', 'long']]

```

## 624 4.4 Vectorising the reviews using Embeddings from Lan- 625 guage Modelling (ELMo)

626 After pre-processing our text we vectorised each review in the document by using ELMo. ELMo  
627 generates a three layered vector representation of each review with a fixed number of features,  
628 that is 1024 features. We extracted the top layer of the ELMo vector for our work because  
629 according to Peters et al. [14] semantic information is captured in the top layer of ELMo vectors.  
630 Below is an example of a review vector for the first negative review.

## 631 4.5 Clustering review vectors using DBSCAN

632 We perform DBSCAN clustering from the vector array produced by ELMo. The DBSCAN algo-  
633 rithm groups together review vectors that are close to each other based on a Euclidean distance  
634 and a minimum number of points. Review vectors in regions with points less than the minimum  
635 number of points are marked as outliers. The DBSCAN algorithm has two main parameters,  
636 these are *eps* and *minpts*. The *eps* specifies how close points should be in a cluster. A small *eps*

value would render most points in the data set to be considered as noise, however a large *eps* will cause clusters to merge, hence unsimilar points would be in the same cluster. Ester et al. [16] suggested we choose an arbitrary point  $p$  and find the distance to the fourth nearest neighbour in the data set. This distance will be our *eps* value. Also, we may fix *eps* based on domain knowledge [26], that is the task or application of clustering can inform how one could select the *eps* value. In this case it is advisable to vary the *minpts* to improve clustering results.

The second parameter *minpts* can be derived from the dimension of the data set. Sander et al. [27] proposed that we set *minpts* to twice the dimension of the data set. However, it is also advisable to increase *minpts* for improved clustering results when the data set is noisy or large [26]. Figure 5.2 is an image of our clustering result for 20 positive reviews. We employed a dimensionality reduction tool PCA to present a two dimensional distribution of the review vectors.

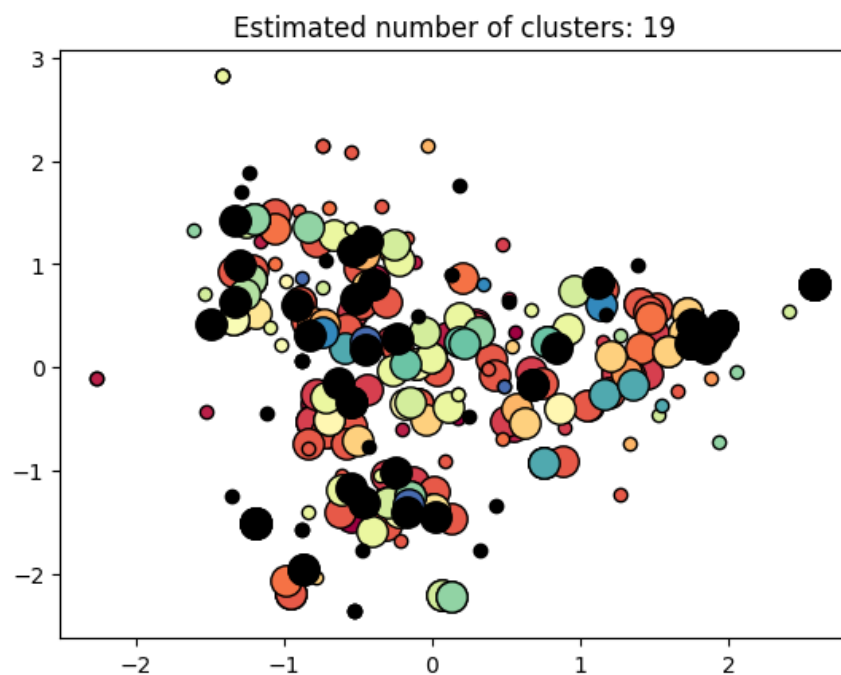


Figure 4.1: Clustering results for 20 positive reviews

Figure 4.2 is an image of our clustering result for 20 negative reviews.



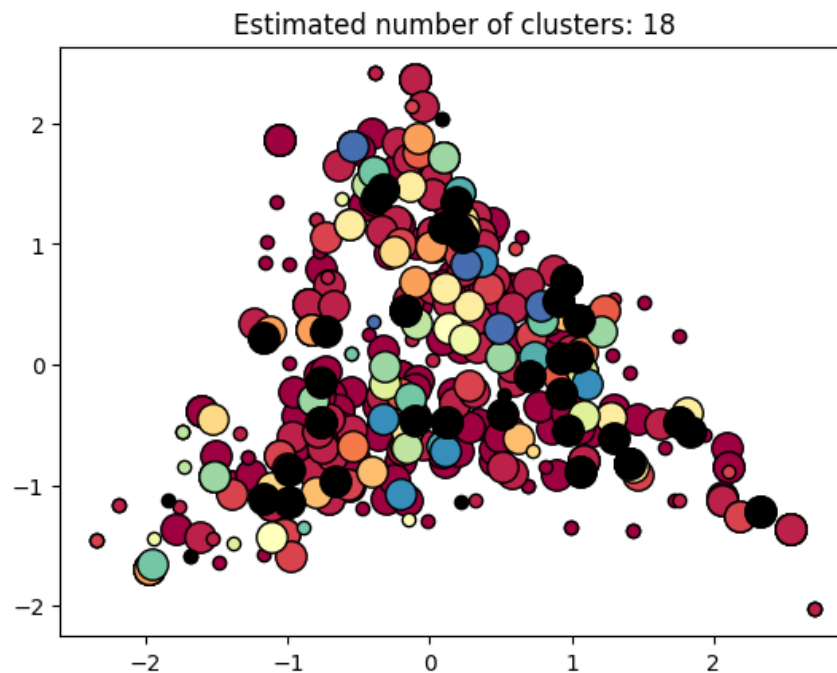


Figure 4.2: Clustering results for 20 negative reviews

## 4.6 Extracting the words in the reviews in the clusters

We extract the reviews in each cluster by using the indices of the cluster labels. DBSCAN clustering algorithm on scikit-learn [28] has attribute called labels. Cluster label of value  $-1$  represents a noise point and all other non-negative cluster label value represents a point (either core point or border point) in the cluster. We extracted the indices of the core points and border points from the list of label values, thereafter we extracted the reviews that correspond with the indices of the core points and border points in the cluster. We present the two clusters for our positive reviews.

Cluster 0:

```
[ 'love lights bright vibrant colors remote worked change colors flash
come v adapter either ',
'son drummer band bought surprise loves says fits drum gear room spare
thanks time purchase chalk one mom',
'works great happy product thanks',
'fantastic piano keyboard fair price got incredible sound quality
effects high quality really fun mess around keys feel
realistic well weighted touch screen icing cake' ]
```

Cluster 5:

```
[ 'leans little stand straight ',
'excellent product ',
'nice quality product ',
```

672 'tried loads tuners however reviewers said best bother anything else get',  
673 'guitar excellent finish nice finish frets ruff edges nice would gave  
674 stars action strings way high tried lower hit pickup could lowed enough  
675 would come undone pickup still high still hit strings use pickup mounting  
676 plate one vintage witch matched binding great much lower height black  
677 one came bridge slammed bottom still strings high looking lower  
678 bridge playable great looking guitar much better looking fender  
679 squire vintage standard wise  
680 cheep bridge hardware bridge design bad hardware used pickup  
681 great sound cheep changing well  
682 tele dreams']

683 We present the two clusters for our negative reviews.

684 Cluster 3:

685 ['item terrible expecting great quality begin bridge neck warped  
686 tuning pegs  
687 right cost repairs violin would proper luthier get damn thing  
688 playable better ordering guitar center', 'records warped took package  
689 pretty disappointed',  
690 'made seatbelt material almost friction unless hold tightly hands  
691 times strap causes guitar slide place apparently people like like  
692 strap hold guitar nicely balanced waste energy trying keep guitar line',  
693 'bought music studio within first week one broke thought  
694 anomaly broken br br problem pretty simple made cheap plastic turn locking  
695 points necessary mic stand plastic inserts simply crumble unless  
696 purchasing whip small children la catholic nuns waste money inexpensive  
697 pack reason',  
698 'looked great put came weeks temporary',  
699 'waste money piece junk',  
700 'good hope another materials']

701

702 Cluster 10:

703 ['firs snark purchased works great  
704 particular one turned  
705 consistently drained battery  
706 time go turn dead disappointed',  
707 'got even stay upright heavy  
708 support even making proper  
709 adjustments  
710 cheap worth excited price  
711 nope shoulda saved money  
712 put nicer one',  
713 'part tang guess called  
714 snapped playing well really

```

715     hard guess broke
716     playing snap went looks nice
717     everything great price guess
718     would almost
719     recommend anyone experience
720     one buy beginner though good
721     learn learned next one get
722     higher quality ',
723     'difficult use uncomfortable
724     sound quality good returning
725     update
726     returned headset grand electronics
727     july request request result
728     review contacted said would
729     replace could update review
730     response
731     new headset nothing update ']
```

## 732 4.7 PPMI Matrix

733 To compute the PPMI matrix, we first constructed a word-cluster count matrix by using CountVec-  
 734 torizer, an object of scikit-learn <sup>3</sup> which converts a collection of text documents to a matrix of  
 735 token counts. Each entry of the word-cluster count matrix specifies the number of times a word  
 736 appears in a given cluster. We proceed by computing the PPMI matrix by using Equation 3.2.2.  
 737 Figure 4.3 is an example of word-cluster count matrix for negative reviews.

738 Figure 4.4 is an example of a PPMI matrix for negative reviews.

## 739 4.8 Sort words in the labels

740 Lastly, we collected all words with non-negative PPMI scores in each cluster. These words are  
 741 now our label vectors for each cluster. We sort the words in the labels we obtained in descending  
 742 order based on their PPMI scores. The most informative words that could be used to summarize  
 743 the clusters have high PPMI scores as compared to the other words in the label vectors. Below  
 744 is an example of label vectors and their respective PPMI scores.

```

745 Cluster_labels 11:
746     {'german ': 5.101337170814387,
747     'student ': 5.101337170814387,
748     'reed ': 5.101337170814387,
749     'reeds ': 5.101337170814387,
```

---

<sup>3</sup>[https://scikit-learn.org/stable/modules/generated/sklearn.feature\\_extraction.text.CountVectorizer.html](https://scikit-learn.org/stable/modules/generated/sklearn.feature_extraction.text.CountVectorizer.html)

	0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17
<b>aa</b>	0	3	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
<b>aaa</b>	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
<b>abandon</b>	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
<b>able</b>	4	2	0	0	0	0	1	0	1	0	0	1	0	0	0	0	0	0
<b>absolute</b>	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
<b>absolutely</b>	2	2	0	0	0	0	0	0	0	0	0	0	0	0	0	1	0	0
<b>absurd</b>	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	0	0
<b>abuse</b>	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
<b>abysmal</b>	0	0	0	0	0	0	0	1	0	0	0	0	0	0	0	0	0	0

Figure 4.3: word-cluster count matrix for 20 negative reviews

	0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17
<b>aa</b>	0.000000	1.246067	0.0	0.0	0.0	0.0	0.000000	0.000000	0.000000	0.0	0.0	0.000000	0.0	0.0	0.0	0.000000	0.0	0.0
<b>aaa</b>	0.000000	1.246067	0.0	0.0	0.0	0.0	0.000000	0.000000	0.000000	0.0	0.0	0.000000	0.0	0.0	0.0	0.000000	0.0	0.0
<b>abandon</b>	0.000000	1.246067	0.0	0.0	0.0	0.0	0.000000	0.000000	0.000000	0.0	0.0	0.000000	0.0	0.0	0.0	0.000000	0.0	0.0
<b>able</b>	0.058802	0.000000	0.0	0.0	0.0	0.0	2.487219	0.000000	1.41548	0.0	0.0	2.904113	0.0	0.0	0.0	0.000000	0.0	0.0
<b>absolute</b>	0.000000	1.246067	0.0	0.0	0.0	0.0	0.000000	0.000000	0.000000	0.0	0.0	0.000000	0.0	0.0	0.0	0.000000	0.0	0.0
<b>absolutely</b>	0.000000	0.329776	0.0	0.0	0.0	0.0	0.000000	0.000000	0.000000	0.0	0.0	0.000000	0.0	0.0	0.0	1.479298	0.0	0.0
<b>absurd</b>	0.000000	0.000000	0.0	0.0	0.0	0.0	0.000000	0.000000	0.000000	0.0	0.0	0.000000	0.0	0.0	0.0	3.088736	0.0	0.0
<b>abuse</b>	0.869733	0.000000	0.0	0.0	0.0	0.0	0.000000	0.000000	0.000000	0.0	0.0	0.000000	0.0	0.0	0.0	0.000000	0.0	0.0
<b>abysmal</b>	0.000000	0.000000	0.0	0.0	0.0	0.0	0.000000	3.471421	0.000000	0.0	0.0	0.000000	0.0	0.0	0.0	0.000000	0.0	0.0
<b>ac</b>	0.869733	0.000000	0.0	0.0	0.0	0.0	0.000000	0.000000	0.000000	0.0	0.0	0.000000	0.0	0.0	0.0	0.000000	0.0	0.0

Figure 4.4: PPMI matrix for 20 negative reviews

```

750 'clarinets ': 5.101337170814387,
751 'distance ': 5.101337170814387,
752 'st ': 5.101337170814387,
753 'somewhere ': 5.101337170814387,
754 'unscrew ': 5.101337170814387,
755 'intermittently ': 5.101337170814387,
756
757 'american ': 5.101337170814387,
758 'counterpart ': 5.101337170814387,
759 'clarinet ': 4.695872062706222,
760 'perhaps ': 4.408189990254441,
761 'excellent ': 4.408189990254441,
762 'acceptable ': 4.408189990254441,
763 'clicking ': 4.408189990254441,
764 'best ': 4.002724882146277,

```

```
765 'disappointment ': 3.715042809694496,  
766 'driver ': 3.491899258380286,  
767 'screw ': 3.155427021759073,  
768 'say ': 3.0218956291345505,  
769 'side ': 3.0218956291345505,  
770 'size ': 3.0218956291345505,  
771 'wrong ': 2.904112593478167,  
772 'brand ': 2.904112593478167,  
773 'give ': 2.904112593478167,  
774 'able ': 2.904112593478167,  
775 'purchase ': 2.850045372207891,  
776 'though ': 2.7987520778203407,  
777 'ever ': 2.616430521026386,  
778 'year ': 2.616430521026386,  
779 'recommend ': 2.616430521026386,  
780 'ordered ': 2.5363878133528495,  
781 'even ': 2.462279841199128,  
782 'know ': 2.393286969712176,  
783 'using ': 2.393286969712176,  
784 'makes ': 2.26812382675817,  
785 'enough ': 2.2109654129182217,  
786 'small ': 2.2109654129182217,  
787 'noise ': 2.0102947174560706,  
788 'much ': 1.8824613459461856,  
789 'would ': 1.856144037628812,  
790 'use ': 1.8432406327929043,  
791 'fit ': 1.8432406327929043,  
792 'return ': 1.6356012680146599,  
793 'used ': 1.4904192581701619,  
794 'could ': 1.4904192581701619,  
795 'one ': 0.4473768206568631}
```

## 5. Results and Discussion

### 5.1 First setting

We implement the algorithm of our approach with Amazon reviews on musical instruments. This data is saved in a tsv file of size 475.2MB and was downloaded from Amazon.com<sup>1</sup>. There are 904,765 reviews in this dataset, after extracting the positive and negative reviews based on star ratings we obtained 106271 negative reviews and 798494 positive reviews. We used an Inspiron 5558 laptop with memory size 8Gib with CPU capacity of 2700 MHz to run our Python codes. Due to the inefficiency of our laptop to run huge dataset, we implemented our algorithm on 70 customer reviews on musical instruments. It took 920.4337 seconds to run the algorithm on 70 positive reviews while it took 164.5895 seconds for 70 negative reviews. We obtained 73 clusters for the negative reviews and 49 clusters for the positive reviews. Each cluster contains an average of 8 reviews. The unique words extracted from the clusters obtained was used in constructing a word-cluster matrix. Each entry of the word-cluster matrix represents the number of occurrence of the  $i$ th word in the  $j$ th cluster. After computing the PPMI matrix by using Equation 3.2.2, we eliminated words with 0 PPMI scores. The remaining words in the word-cluster matrix were ranked according to their PPMI scores. The PPMI score of a word is a measure of the significance of the word in a cluster. Therefore, by our approach we were able to extract key words from customer reviews on Amazon. See the Appendix for some results.

---

<sup>1</sup>[https://s3.amazonaws.com/amazon-reviews-pds/tsv/amazon\\_reviews\\_us\\_Musical\\_Instruments\\_v1\\_0.tsv.gz](https://s3.amazonaws.com/amazon-reviews-pds/tsv/amazon_reviews_us_Musical_Instruments_v1_0.tsv.gz)

# References

- [1] Martin Kütz. *ntroduction to E-Commerce: Combining Business and Information Technology*. 2016.
- [2] M MUDAMBI Susan and S David. What makes a helpful online review? a study of customer reviews on amazon. com. *MIS Quarterly*, 34(1):185–200, 2010.
- [3] inc. Amazon Web Services. Amazon customer reviews from us customers for the product "watches". Retrieved from [https://s3.amazonaws.com/amazon-reviews-pds/tsv/amazon\\_reviews\\_us\\_Watches\\_v1\\_00.tsv.gz](https://s3.amazonaws.com/amazon-reviews-pds/tsv/amazon_reviews_us_Watches_v1_00.tsv.gz), Accessed May, 2019.
- [4] Kenton Lee, Luheng He, Mike Lewis, and Luke Zettlemoyer. End-to-end neural coreference resolution. *arXiv preprint arXiv:1707.07045*, 2017.
- [5] Qian Chen, Xiaodan Zhu, Zhenhua Ling, Si Wei, Hui Jiang, and Diana Inkpen. Enhanced lstm for natural language inference. *arXiv preprint arXiv:1609.06038*, 2016.
- [6] Luheng He, Kenton Lee, Mike Lewis, and Luke Zettlemoyer. Deep semantic role labeling: What works and what’s next. In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 473–483, 2017.
- [7] Noah A Smith. Contextual word representations: A contextual introduction. *arXiv preprint arXiv:1902.06006*, 2019.
- [8] Josh Patterson and Adam Gibson. *Deep learning: A practitioner’s approach*. " O’Reilly Media, Inc.", 2017.
- [9] Tomas Mikolov, Ilya Sutskever, Kai Chen, Greg S Corrado, and Jeff Dean. Distributed representations of words and phrases and their compositionality. In *Advances in neural information processing systems*, pages 3111–3119, 2013.
- [10] Jeffrey Pennington, Richard Socher, and Christopher Manning. Glove: Global vectors for word representation. In *Proceedings of the 2014 conference on empirical methods in natural language processing (EMNLP)*, pages 1532–1543, 2014.
- [11] Matthew E Peters, Waleed Ammar, Chandra Bhagavatula, and Russell Power. Semi-supervised sequence tagging with bidirectional language models. *arXiv preprint arXiv:1705.00108*, 2017.
- [12] W Haas. J.r. firth: Papers in linguistics, 1934–1951. xii, 233 pp., 11 plates. london, etc. *Bulletin of The School of Oriental and African Studies-university of London*, 21, 10 1958. doi: 10.1017/S0041977X00060559.
- [13] Ian Jolliffe. *Principal component analysis*. Springer, 2011.
- [14] Matthew E Peters, Mark Neumann, Mohit Iyyer, Matt Gardner, Christopher Clark, Kenton Lee, and Luke Zettlemoyer. Deep contextualized word representations. *arXiv preprint arXiv:1802.05365*, 2018.

- [15] Christopher D Manning, Christopher D Manning, and Hinrich Schütze. *Foundations of statistical natural language processing*. MIT press, 1999.
- [16] Martin Ester, Hans-Peter Kriegel, Jörg Sander, Xiaowei Xu, et al. A density-based algorithm for discovering clusters in large spatial databases with noise. In *Kdd*, volume 96, pages 226–231, 1996.
- [17] Leonard Kaufman and Peter J Rousseeuw. Partitioning around medoids (program pam). *Finding groups in data: an introduction to cluster analysis*, pages 68–125, 1990.
- [18] Yogita Rani<sup>1</sup> and Harish Rohil. A study of hierarchical clustering algorithm. *ter S & on Te SIT-2*, page 113, 2013.
- [19] James MacQueen et al. Some methods for classification and analysis of multivariate observations. In *Proceedings of the fifth Berkeley symposium on mathematical statistics and probability*, volume 1, pages 281–297. Oakland, CA, USA, 1967.
- [20] Yuzhen Zhao, Xiyu Liu, and Jianhua Qu. The k-medoids clustering algorithm by a class of p system. *JOURNAL OF INFORMATION & COMPUTATIONAL SCIENCE*, 9(18):5777–5790, 2012.
- [21] Raymond T Ng and Jiawei Han. Clarans: A method for clustering objects for spatial data mining. *IEEE Transactions on Knowledge & Data Engineering*, (5):1003–1016, 2002.
- [22] Anil K Jain and Richard C Dubes. *Algorithms for clustering data*. 1988.
- [23] Gerlof Bouma. Normalized (pointwise) mutual information in collocation extraction. *Proceedings of GSCL*, pages 31–40, 2009.
- [24] James H Martin and Daniel Jurafsky. *Speech and language processing: An introduction to natural language processing, computational linguistics, and speech recognition*. Pearson-/Prentice Hall Upper Saddle River, 2009.
- [25] Kenneth Ward Church and Patrick Hanks. Word association norms, mutual information, and lexicography. *Computational linguistics*, 16(1):22–29, 1990.
- [26] Erich Schubert, Jörg Sander, Martin Ester, Hans Peter Kriegel, and Xiaowei Xu. Dbscan revisited, revisited: why and how you should (still) use dbscan. *ACM Transactions on Database Systems (TODS)*, 42(3):19, 2017.
- [27] Jörg Sander, Martin Ester, Hans-Peter Kriegel, and Xiaowei Xu. Density-based clustering in spatial databases: The algorithm gdbscan and its applications. *Data mining and knowledge discovery*, 2(2):169–194, 1998.
- [28] Fabian Pedregosa, Gaël Varoquaux, Alexandre Gramfort, Vincent Michel, Bertrand Thirion, Olivier Grisel, Mathieu Blondel, Peter Prettenhofer, Ron Weiss, Vincent Dubourg, et al. Scikit-learn: Machine learning in python. *Journal of machine learning research*, 12(Oct):2825–2830, 2011.
- [29] Robert M Fano and David Hawkins. Transmission of information: A statistical theory of communications. *American Journal of Physics*, 29:793–794, 1961.



# APPENDICES

## 5.2 Python codes for the entire pipeline for our approach

```
# import the needed package for the tasks
import numpy as np
import pandas as pd
import re
from nltk.corpus import stopwords
import matplotlib.pyplot as plt
from sklearn.feature_extraction.text import CountVectorizer
from allennlp.commands.elmo import ElmoEmbedder
from sklearn.cluster import DBSCAN
from sklearn import metrics
from sklearn.preprocessing import StandardScaler

import time
import random

then = time.time() #Time before the operations start

#DO YOUR OPERATIONS HERE

# parameters for the pipeline
num_reviews = 5
eps = 0.3
min_samples = 4

# read the tsv file
obs=[]
filename = "/home/salomey/Salomay/ESSAY_PHASE_A/Data_for_project/amazon_rev
file = open(filename , mode='r')
for line in file.readlines():
    obs.append(line.split('\t'))
```

926  
927  
928  
929  
930  
931  
932  
933  
934  
935  
936  
937  
938  
939  
940  
941  
942  
943  
944  
945  
946  
947  
948  
949  
950  
951  
952  
953  
954  
955  
956  
957  
958  
959  
960  
961  
962  
963  
964  
965  
966  
967  
968  
969  
970

```
# visualizing the data in a dataframe using pandas
data=pd.DataFrame(obs,columns=["marketplace ","customer_id" ,"review_id ","p

# drop the first row
data = data.drop([0], axis=0)
#print("These are the amazon review data for video",data)

# extracting positive reviews based on star ratings>2
positive_review= data[data['star_rating']>'2']

# extracting the positive review and converting into a list
positives = positive_review.review_body.tolist()
#print("These are the positive reviews:",positives)

# extracting negative reviews based on star ratings<=2
negative_df=data[data['star_rating']<='2']

# converting negative reviews into a list
negatives = negative_df.review_body.tolist()
#print("These are the negative reviews:",negatives)
# pre-processing reviews(remove stop words and non-alphabetic characters)

def process(review , remove_stopwords=True):

    review_text = re.sub("[^a-zA-Z]"," ", review)

    words = review_text.lower().split()

    if remove_stopwords:
        stops = set(stopwords.words("english"))
        words = [w for w in words if not w in stops]

    cleaned_reviews=[]

    for word in words:
        cleaned_reviews.append(word)
```

```

971
972     return(cleaned_reviews)
973
974 cleaned_review_body = []
975 for review in negatives:# change "negatives" to "positives" if you are inter
976     cleaned_review_body.append( " ".join(process(review)))
977
978 tokkenn=[]
979 for r in cleaned_review_body:
980     tokkenn.append(process(r))
981 #print("these are the pre-processed words", tokkenn)
982 # converting the words into vectors using ELMo
983
984 elmo = ElmoEmbedder()
985 word_embeddings = []
986 for reviews in tokkenn[:num_reviews]:
987     word_embeddings.append(elmo.embed_sentence(reviews)[0]) # select the top
988
989 #print("These are the word embeddings for the reviews",word_embeddings)
990 #print(len(word_embeddings))
991 X=(np.concatenate(word_embeddings))
992 #print(type(X))
993 #print("This is the shape of our vector X",np.shape(X))
994
995 #####
996
997 # Now we cluster our vector X using sklearn DBSCAN
998
999 X = StandardScaler().fit_transform(np.asarray(X))
1000 plt.rcParams.update({'figure.max_open_warning': 0})
1001
1002
1003 db = DBSCAN(eps, min_samples, algorithm="kd_tree").fit(X)
1004
1005 core_samples_mask = np.zeros_like(db.labels_, dtype=bool)
1006 core_samples_mask[db.core_sample_indices_] = True
1007 labels = db.labels_
1008 #print("these are labels",labels)
1009
1010 #print("These are the core sample indices",db.core_sample_indices_)
1011
1012 n_clusters_ = len(set(labels)) - (1 if -1 in labels else 0)
1013 n_noise_ = list(labels).count(-1)
1014 #print('Estimated number of clusters: %d' % n_clusters_)
1015 #print('Estimated number of noise points: %d' % n_noise_)

```

```

1016
1017 # Black removed and is used for noise instead.
1018 unique_labels = set(labels)
1019 fig = plt.figure()
1020 colors = [plt.cm.Spectral(each)
1021            for each in np.linspace(0, 1, len(unique_labels))]
1022 for k, col in zip(unique_labels, colors):
1023     if k == -1:
1024         # Black used for noise.
1025         col = [0, 0, 0, 1]
1026
1027     class_member_mask = (labels == k)
1028
1029     xy = X[class_member_mask & core_samples_mask]
1030     plt.plot(xy[:, 0], xy[:, 1], 'o', markerfacecolor=tuple(col),
1031              markeredgecolor='k', markersize=14)
1032
1033     xy = X[class_member_mask & ~core_samples_mask]
1034     plt.plot(xy[:, 0], xy[:, 1], 'o', markerfacecolor=tuple(col),
1035              markeredgecolor='k', markersize=6)
1036
1037
1038 plt.title('Estimated number of clusters: %d' % n_clusters_)
1039
1040 fig.savefig('cluster_image.png', bbox_inches='tight')
1041
1042 # extracting the reviews from the clusters
1043
1044 index2 = []
1045 counter = 0
1046 for i in labels:
1047     if i != -1: # remove the noise labels
1048         index2.append(counter)
1049     counter+=1
1050
1051 label2 = []
1052 for j in index2:
1053     k = labels[j]
1054     label2.append(k) # these are the labels for points in each cluster
1055
1056
1057 d = {"Index":index2, "Labels":label2}
1058
1059 df = pd.DataFrame(d)
1060

```

```

1061 s = pd.Series(["Labels"])
1062
1063
1064 # extracting the indices of the labels in the cluster
1065 index_kk = []
1066
1067 for i in range(len(label2)):
1068     p= df[df['Labels']==i]
1069
1070     index_kk.append(p.Index.tolist())
1071 #print("index for each label:",index_kk )
1072 index_labels = [x for x in index_kk if x != []]
1073 #print ("These are the indices for each label",index_labels)
1074
1075 # extracting the reviews in each cluster
1076 empty_lst = []
1077 for i in index_labels:
1078     empty_lst.append([cleaned_review_body[j] for j in i])
1079
1080 #print ("these are the reviews found in each cluster",empty_lst)
1081
1082 for i in range(n_clusters_):
1083     #print("Cluster %d:" % i, ' %s' % empty_lst[i])
1084     print("Cluster %d:" % i, ' %s' % len(empty_lst[i]))
1085
1086 # all reviews in a cluster is now a single token
1087 clust = []
1088 for i in range(len(empty_lst)):
1089     t= ' '.join(empty_lst[i])
1090     clust.append(t)
1091 #print("These are the reviews in a cluster",clust)
1092
1093 #####
1094
1095
1096 # building word-cluster count matrix
1097
1098 def fn_tdm_df(docs, xColNames = None):
1099     ''' create a term document matrix as pandas DataFrame
1100     with **kwargs you can pass arguments of CountVectorizer
1101     if xColNames is given the dataframe gets columns Names'''
1102
1103     #initialize the vectorizer
1104     vectorizer = CountVectorizer()
1105     x1 = vectorizer.fit_transform(clust)

```

```

1106     #create dataframe
1107     df = pd.DataFrame(x1.toarray().transpose(), index = vectorizer.get_featu
1108
1109     if xColNames is not None:
1110         df.columns = xColNames
1111
1112     return df
1113
1114 for i in range(len(clust[:5])):
1115
1116     count_matrix=fn_tdm_df(docs=clust, xColNames=None)
1117     print("This is the count matrix",count_matrix.head())
1118
1119 # building word-cluster PMI matrix
1120
1121 def pmi(df):
1122     '''
1123     Calculate the positive pointwise mutual information score for each entry
1124     https://en.wikipedia.org/wiki/Pointwise_mutual_information
1125     We use the  $\log( p(y|x)/p(y) )$ , y being the column, x being the row
1126     '''
1127     # Get numpy array from pandas df
1128     arr = df.as_matrix()
1129
1130     #  $p(y|x)$  probability of each t1 overlap within the row
1131     row_totals = arr.sum(axis=1).astype(float)
1132     prob_cols_given_row = (arr.T / row_totals).T
1133
1134     #  $p(y)$  probability of each t1 in the total set
1135     col_totals = arr.sum(axis=0).astype(float)
1136     prob_of_cols = col_totals / sum(col_totals)
1137
1138     # PMI:  $\log( p(y|x) / p(y) )$ 
1139     # This is the same data, normalized
1140     ratio = prob_cols_given_row / prob_of_cols
1141     ratio[ratio==0] = 0.00001
1142     _pmi = np.log(ratio)
1143     _pmi[_pmi < 0] = 0
1144
1145     return _pmi
1146
1147
1148
1149
1150 # PMI matrix into a data frame

```

```

1151 c=pmi(count_matrix)
1152 vectorizer = CountVectorizer()
1153 x1 = vectorizer.fit_transform(clust)
1154 PMI_df = pd.DataFrame(c, index = vectorizer.get_feature_names())
1155
1156 # extracting the labels of each cluster with a positive pmi score into a list
1157 List = []
1158 for c in range(len(empty_lst)):
1159     col=PMI_df[c]
1160     List.append(col[col>0].sort_values(ascending=False).to_dict())
1161 #print('These are the label vectors for each cluster')
1162
1163 for i in range(n_clusters_):
1164     print("Cluster_labels %d:\n" % i, ' %s' % List[i].keys())
1165
1166
1167
1168 now = time.time() #Time after it finished
1169
1170 print("It took: ", now-then, " seconds")

```

## 1171 5.3 Cluster images for both positive and negative re- 1172 views

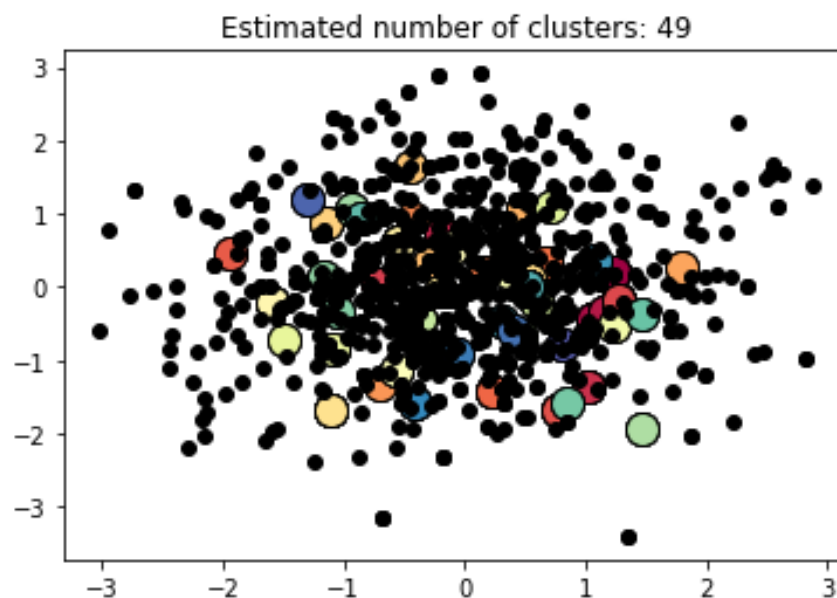


Figure 5.1: Clustering results for 70 positive reviews

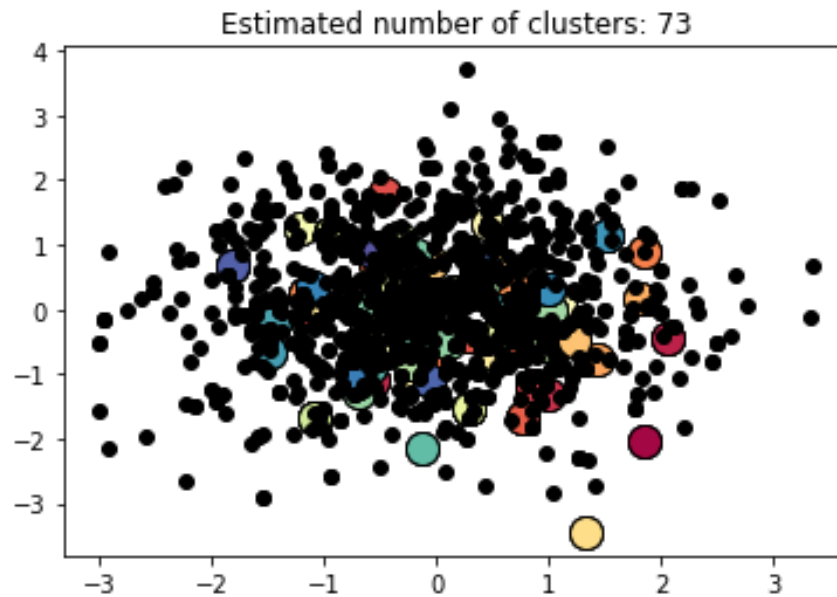


Figure 5.2: Clustering results for 70 negative reviews

## 5.4 Labels for the first six selected clusters

Cluster\_labels 0:

```
dict_keys(['cast', 'factory', 'phillips', 'originals', 'randomly',
'rd', 'recorders', 'word', 'speaking', 'substitute', 'ping',
'grainy', 'backlash', 'possess', 'events', 'welll', 'dean',
'washed', 'die', 'strips', 'operate', 'turns', 'smoothly', 'channels',
'pickup', 'mics', 'didnt', 'chinese', 'pain', 'wasting', 'front',
'dime', 'bother', 'dont', 'mediocre', 'results', 'fell', 'crap',
'tuners', 'manufacturer', 'hoping', 'built', 'directly', 'almost',
'returning', 'barely', 'far', 'plastic', 'useless', 'iphone',
'cheaply', 'cut', 'thin', 'still', 'away', 'stay', 'pick', 'getting',
'system', 'bridge', 'better', 'violin', 'apart', 'guitars',
'junk', 'line', 'replace', 'unless', 'properly', 'comes', 'screws',
'pretty', 'way', 'low', 'think', 'different', 'send', 'cheap',
'much', 'well', 'buying', 'tune', 'broken', 'hold', 'ok', 'disappointed',
'made', 'good', 'unit', 'also', 'put', 'work', 'small', 'bad',
'using', 'mic', 'use', 'used', 'money', 'broke', 'first', 'new',
'time', 'even', 'two', 'back', 'one', 'guitar', 'got', 'sound', 'quality'])
```

Cluster\_labels 1:

```
dict_keys(['weights', 'reducer', 'germany', 'steel', 'unplug',
'designed', 'behringer', 'opened', 'removed', 'plate', 'tuner',
'produce', 'noise', 'pedal', 'china', 'turned', 'ends', 'inside', 'xlr', ...])
```



```

1197 'connection', 'install', 'cheaply', 'plug', 'connector', 'last', 'amp', 'v
1198 'bottom', 'junk', 'happy', 'nothing', 'apart', 'looked', 'line', 'ever', '
1199 'plastic', 'makes', 'touch', 'less', 'day', 'buying', 'thing', 'recommend'
1200 'guitar', 'fit', 'bad', 'much', 'mic', 'even', 'made', 'great', 'back', 'c
1201 Cluster_labels 2:
1202 dict_keys(['advise', 'descriptions', 'starter', 'souvenir',
1203 'someday', 'solution', 'shoddy', 'seal', 'changed', 'harmonica',
1204 'chassis', 'promised', 'coolest', 'department', 'distances',
1205 'astounding', 'nilly', 'drilled', 'msd', 'motu', 'easier', 'mkii',
1206 'exist', 'frustrate', 'internet', 'industrial', 'graders', 'hindsight',
1207 'teach', 'pwma', 'harmonicas', 'ticking', 'alvarez', 'arts',
1208 'willy', 'wireless', 'asking', 'upper', 'air', 'tombo', 'tremolo',
1209 'thus', 'wt', 'octave', 'tech', 'pyle', 'seriously', 'paper',
1210 'amplifier', 'ukulele', 'respond', 'neither', 'heats', 'brought',
1211 'workmanship', 'learning', 'electronic', 'initial', 'sustain',
1212 'assistance', 'everywhere', 'worst', 'materials', 'student',
1213 'playable', 'weird', 'night', 'none', 'making', 'external',
1214 'owned', 'flute', 'three', 'probably', 'unplayable', 'takes',
1215 'sounding', 'school', 'serious', 'wasted', 'experience', 'difference',
1216 'products', 'middle', 'beginner', 'drivers', 'musical', 'awful',
1217 'basically', 'project', 'worthless', 'mm', 'twice', 'hot', 'tell',
1218 'heard', 'units', 'output', 'amps', 'th', 'manufacturer', 'instrument',
1219 'band', 'find', 'kit', 'pass', 'pads', 'save', 'returned', 'dollars',
1220 'ever', 'hard', 'audio', 'days', 'whole', 'note', 'static',
1221 'sounded', 'instead', 'purchase', 'notes', 'warranty', 'take',
1222 'full', 'support', 'see', 'metal', 'play', 'away', 'called',
1223 'ago', 'decent', 'feels', 'years', 'playing', 'maybe', 'sounds',
1224 'much', 'worked', 'toy', 'speakers', 'brand', 'second', 'played',
1225 'help', 'jack', 'gave', 'ordered', 'couple', 'junk', 'matter',
1226 'apart', 'minutes', 'color', 'trying', 'started', 'defective',
1227 'might', 'pay', 'least', 'like', 'received', 'properly', 'expensive',
1228 'also', 'sent', 'even', 'lot', 'holes', 'far', 'price', 'stay',
1229 'year', 'bad', 'looking', 'give', 'go', 'tune', 'money', 'looks',
1230 'day', 'noise', 'buying', 'thing', 'end', 'never', 'ok', 'disappointed',
1231 'get', 'bought',
1232 'something', 'would', 'another', 'really', 'cheap', 'well', 'use',
1233 'used', 'new', 'good'])
1234 'waste', 'keep', 'least', 'unit', 'probably', 'expensive', 'properly',
1235 'however', 'sent', 'problem', 'cable', 'many', 'less', 'high',
1236 'away', 'different', 'strings', 'amp', 'well', 'working', 'looks',
1237 'without', 'since', 'guitar', 'first', 'need', 'poor', 'make',
1238 'recommend', 'two', 'work', 'quality', 'long', 'fine',
1239 'something',
1240 'works', 'product', 'also', 'worked', 'amazon', 'sound', 'really',
1241 'much', 'money', 'would', 'even', 'time', 'great'])

```

Cluster\_labels 4:

```
dict_keys(['ac', 'posting', 'curling', 'smoking', 'glanced',
'dollar', 'elected', 'prepped', 'needle', 'suppor', 'exit',
'machines', 'lemon', 'flawlessly', 'knobs', 'fluke', 'standard',
'samson', 'apparently', 'sink', 'appears', 'arm', 'seconds',
'guide', 'shipped', 'unplugged', 'saw', 'smoke', 'particular',
'singing', 'shot', 'shock', 'owned', 'needless', 'nd', 'sad',
'point', 'start', 'unusable', 'shipping', 'fact', 'items', 'electronics',
'mine', 'mount', 'turning', 'supposed', 'star', 'machine', 'adapter',
'tight', 'replacement', 'build', 'device', 'weight', 'update',
'felt', 'read', 'take', 'things', 'obviously', 'quickly', 'either',
'help', 'stand', 'seem', 'review', 'anything', 'keep', 'say',
'turn', 'lot', 'plugged', 'every', 'wrong', 'touch', 'bit',
'looking', 'sure', 'enough', 'power', 'last', 'give', 'thought',
'disappointed', 'item', 'something', 'would', 'months', 'problem',
'worked', 'another', 'fit', 'small', 'bad', 'tried', 'well',
'one', 'could', 'buy', 'back', 'get', 'quality', 'like', 'good',
'product']])
```

Cluster\_labels 5:

```
dict_keys(['flying', 'picking', 'rattling', 'icicle', 'fit',
'someting', 'prescreen', 'embarrassed', 'soundtrack', 'special',
'plig', 'tambourines', 'taping', 'recordings', 'lo', 'pitched',
'converter', 'warm', 'marimba', 'analog', 'man', 'audition',
'timpani', 'watching', 'painted', 'buzzy', 'percussionist', 'classical',
'local', 'instruments', 'recording',
'lp', 'laptop', 'jingles', 'quite', 'gig', 'meinl', 'everywhere',
'mid', 'fixed', 'albeit', 'duct', 'foot', 'epiphone', 'chrome',
'wasting', 'digital', 'tambourine', 'stock', 'regular', 'noisy',
'provided', 'shop', 'quick', 'sg', 'settings', 'particular',
'apart', 'lower', 'noticed', 'okay', 'leather', 'awful', 'past',
'fell', 'project', 'middle', 'device', 'toy', 'higher', 'zero',
'version', 'song', 'simple', 'large', 'fall', 'show', 'build',
'guess', 'inside', 'try', 'wish', 'advertised', 'metal', 'needed',
'things', 'body', 'pieces', 'feel', 'kept', 'buying', 'seems',
'garbage', 'feedback', 'computer', 'matter', 'couple', 'find',
'nothing', 'terrible', 'least', 'phone', 'music', 'however',
'pretty', 'makes', 'plastic', 'like', 'come', 'wrong', 'purchase',
'old', 'many', 'think', 'bit', 'bad', 'cheap', 'years', 'mic',
'broken', 'go', 'noise', 'thought', 'first', 'end', 'need',
'set', 'quality', 'received', 'right', 'something', 'way', 'play',
'sound', 'fit', 'tried', 'much', 'money', 'product', 'would',
'could', 'made', 'two', 'bought', 'back', 'guitar', 'got']])
```