2020-03-17

# DICTIONARY AND RULE BASED (HYBRID) APPROACH OF AMHARIC SENTIMENT ANALYSIS FOR KANA TV

## MEKONNEN, YOHANNES

Wisdom at the source of the Blue Nile

**BAHIR DAR UNIVERSITY**

**BAHIR DAR INSTITUTE OF TECHNOLOGY**

**SCHOOL OF RESEARCH AND POSTGRADUATE STUDIES**

**FACULTY OF COMPUTING.**

Thesis Title: **DICTIONARY AND RULE BASED (HYBRID) APPROACH OF AMHARIC SENTIMENT ANALYSIS FOR KANA TV**

**BY: YOHANNES MEKONNEN**

JULY, 2018

BAHIRDAR, ETHIOPIA

# DICTIONARY AND RULE BASED (HYBRID) APPROACH OF AMHARIC SENTIMENT ANALYSIS FOR KANA TV

## BY: YOHANNES MEKONNEN

## A THESIS SUBMITTED TO THE SCHOOL OF RESEARCH AND GRADUATE STUDIES OF BAHIR DAR INSTITUTE OF TECHNOLOGY, BDU IN PARTIAL FULFILLMENT FOR THE DEGREE OF MASTERS IN THE COMPUTER SCIENCE IN THE FACULTY OF COMPUTING.

BahirDar, Ethiopia

July, 2018

## DECLARATION

I, undersigned, declare the thesis comprises my own study. In compliance with globally accepted practices, I have acknowledged and refereed all materials used in this work. I understand that non-faithfulness to the principles of academic honesty and integrity, misrepresentation/ fabrication of any idea/data/fact/source will constitute sufficient ground for disciplinary action by the University and can also suggest penal action from the sources which have not been properly cited or acknowledged.

Name of the student_____Yohannes Mekonnen_____ Signature _____

Date of submission: _11/07/2018_

Place:   Bahir Dar

This thesis has been submitted for examination with my approval as a university advisor.

Advisor Name: _Dr. Tesfa Tegegne_____

Advisor's Signature: _____

# Bahir Dar University

# Bahir Dar Institute of Technology-

# School of Research and Graduate Studies

# Faculty of Computing

## THESIS APPROVAL SHEET

| Yohannes Mekonnen | | 24/06/2018 |
|---|---|---|
| Name of student | Signature | Date |

The following graduate faculty members certify that this student has successfully presented the necessary written final thesis and oral presentation for partial fulfillment of the thesis requirements for the Degree of Master of Science in Computer Science.

**Approved By:**

Advisor:
Dr. Tesfa Tegegne                                                    24/06/2018

Name                              Signature                          Date

External Examiner:
Mesfin Kifle (PhD)                                                   26/06/2018

Name                              Signature                          Date

Internal Examiner:
Dr. Gebeyehu Belay                                                   03/07/2018

Name                              Signature                          Date

Chair Holder:

Bhabani Shankar D. M.             BSDasmRapti       July 3, 2018

Name                              Signature                          Date

Faculty Dean:

Dawed Nesru (Msc.)                                                   July 3, 2018

Name                              Signature                          Date

iv

Dedicated to my Gift, to my Love, Mar

በሆነው ሁሉ የሥጦታየ የማር እጅ አለበት !!!

# ACKNOWLEDGEMENTS

# ABSTRACT

Sentiment Analysis (SA) is an ongoing field of research in text mining field. SA is the computational treatment of opinions, sentiments and subjectivity of text. [1]

The Comments which is given by the viewers/non-viewers of the program that reflect whether the program is positive (positive incremental) or negative (negative decrement) or neutral. SA can be analyzing the given text into predefined categories as positive, incremental positive, negative, decrement negative or neutral based on the sentiment terms that appear with in the opinionated documents. These comments need to be explored, analyses and organized for better decision making.

The earlier related researcher not enough consider the POS tagging which is very important to identify the polarity of the sentiment. And did not consider the irony and stairs of expressions. And also they only considered positive and negative polarities but important to consider the inverter words that will change the polarity. In our paper the gaps basically held by using NLP Techniques.

Sentiment analysis system applied to solve the polarity by using Rule based and Dictionary approach. The comments collected from the viewers/non-viewers from website/Facebook page, focusing group discussion, and by distributing an open ended questioner. The experiments are conducted using 1022 (one thousands and twenty two) sentiment comments with four target research area. The average Accuracy, Precession, Recall and f-score respectably are 0.85, 0.95, 0.89 and 0.91. Experimental results using viewers of comments shows that the effectiveness of the system.

.

**Keywords:**   hybrid approach, polarity, optioned Documents, Sentiment Analysis, Focusing group discussion, NLP Technique.

# TABLE OF CONTENTS

# LIST OF ACRONYMS

| LR | Literature review |
|----|----|
| SA | Sentiment Analysis |
| NLP | Natural Language Processing |
| GI | General Inquirer |
| POS | Part Of Speech |
| IMDB | Internet Movie Data Base |
| SVM | Support Vector Machine |
| ML | Machine Learning |
| UN | Unclassified |
| ASCII | American Standard Code for Information Interchange |
| NLTK | Natural Language Toolkit |
| Dec | Decrement |
| Inc | Incremental |

# LIST OF FIGURES

# LIST OF TABLES

# LIST OF APPENDIXES

# LIST OF ALGORITHMS

# CHAPTER ONE

# 1. INTRODUCTION

## 1.1. Background

Natural language is a language that has evolved naturally as a means of communication among people and one of the behavior of human being. NLP is an area of research and application that explores how computers can be used to understand and manipulate natural language text or speech to do useful things. [2]

**Sentiment analysis** (sometimes known as **opinion mining** or emotion AI) refers to the use of natural language processing, text **analysis**, computational linguistics, and biometrics to systematically identify, extract, quantify, and study affective states and subjective information. Sentiment analysis which is also known as emotion mining, attitude mining or subjectivity mining, is a hot research discipline which is concerned with the computational study of opinions, sentiments and emotions expressed in an opinionated text. A basic task in sentiment analysis is classifying the polarity of a given text at the document, sentence, or feature/aspect level whether the expressed opinion in a document, a sentence or an entity feature/aspect is positive, negative, or neutral. There are three levels on which sentiment analysis can be performed: [3]

**Document level**

The task at this level is to classify whether a whole opinion document expresses a positive or negative sentiment. For example, given a product review, the system determines whether the review expresses an overall positive or negative opinion about the product. This task is commonly known as document-level sentiment classification. This level of analysis assumes that each document expresses opinions on a single entity (e.g., a single product). Thus, it is not applicable to documents which evaluate or compare multiple entities.

**Sentence level**

The task at this level goes to the sentences and determines whether each sentence expressed a positive, negative, or neutral opinion. Neutral usually means no opinion.

This level of analysis is closely related to subjectivity classification which distinguishes sentences (called objective sentences) that express factual information from sentences (called subjective sentences) that express subjective views and opinions. However, we should note that subjectivity is not equivalent to sentiment as many objective sentences can imply opinions, e.g., "We bought the car last month and the windshield wiper has fallen off."

**Entity and Aspect level**

Aspect level performs finer-grained analysis. Aspect level was earlier called feature level (feature-based opinion mining and summarization) Instead of looking at language constructs (documents, paragraphs, sentences, clauses or phrases), aspect level directly looks at the opinion itself. It is based on the idea that an opinion consists of a sentiment (positive or negative) and a target (of opinion).

**<u>About KANA broadcasting/TV</u>**

**Kana TV** ("ቃና") is a private commercial television station operating in Ethiopia. Kana TV is part of *Moby* Media Group. The channel operates solely in the Amharic language with half of the content being local and half being foreign content that is dubbed. It has quickly earned the title of being the most watched television channel in Ethiopia. This is mostly due to the popularity of its dubbed foreign TV dramas. Among the most famous are Turkish dramas such as 'Kuzey Güney 'ኩዚ ጉኒ (North South), 'Yetekema Hiwot' የተቀማ ሕይወት (Lively Life) and 'Tikur Fiker' ጥቁር ፍቅር 'Kara Para Aşk' (Black Love). [4]

The channel has drawn criticism from conservative commentators who argue that over-consumption of foreign soap operas dubbed in Amharic will corrupt Ethiopian culture. Kana Television started its broadcast in March 2016 and has offices located in Addis Ababa, Ethiopia. The channel broadcasts on Nile Sat. [4]

Figure 1: Sample Satellite dishes for viewing different program [4]

A significant part of our comment collecting has been to find what the people think and feel about different subjects. To perform comments from the kana TV program, sentiment analysis will be a solution to categorize as positive, negative or neutral comments.

SA is widely applied to voice of the customer materials such as reviews and survey responses, online and social media, and healthcare materials for applications that range from marketing to customer service to clinical medicine.

Generally speaking, sentiment analysis aims to determine the attitude of a speaker, writer, or other subject with respect to some topic or the overall contextual polarity or emotional reaction to a document, interaction, or event. The attitude may be a judgment or evaluation (see appraisal theory) [5], affective state (that is to say, the emotional state of the author or speaker), or the intended emotional communication (that is to say, the emotional effect intended by the author or interlocutor).

## 1.2. Motivation

Currently Kana TV viewers are many in Ethiopia, which give positive and negative comments by using social media websites about the program. The comments are very important for the kana TVs production Managers or owners to approve the program.

The motivation behind this study is the hot issues - good and bad comments which are given by the viewers/non-viewers of kana TV program and analyze opinions for the purpose of decision making. Therefore, it is necessary to conduct a research on sentiment analysis for kana TV that will help to be good program.

## 1.3. Statement of the problem

At the present time kana TV program viewers/non-viewers respond through websites/Facebook page, good and bad comments. Most of the time, KANA TV program is prepared movies and different programs into Amharic languages for Ethiopians. During viewing this program, comments will be taken from different aspects. Most people said that the kana TVs program is very interesting than other Ethiopian TV programs; especially programs broadcasted by Ethiopian broad Casting Corporation are not attractive. On the contrary, others argue that Kana TV program influences and spoils the norm and the culture of the society.

However, there is a claim that teenagers are more attracted and become addicted the program. Some even claim that some of the programs broadcasted in KANA become a discussion point at work places, schools etc.

Nowadays Companies and other profit and non-profit organizations have accumulated a vast amount of data on how their employees or customers feel about the products and services they receive from the aforementioned organizations.

In Ethiopia organizational situation different organization uses Twitter, Facebook, LinkedIn, etc. to gather information about how prefer about their services in Amharic language but in the meantime the information is accumulated through this tools are vast. Therefore we need a sentiment analysis system to extract important information from information gathered. That is why we are conducted this study.

Hence, sentiment analysis for kana TV has been proposed using NLP Techniques. This study answers the following questions:

1. How computationally identify and categorize sentiments in a peace of Amharic texts?
2. How to classify the polarity of a given text at the word, sentence or an aspect level?
3. What types of comments are suggested by the viewers/non-viewers of the program?
4. How to expressed opinion in a word, sentence or an aspect level is positive, negative or neutral?
5. What is the performance of sentiment analysis of the system compared to manually identify?

## 1.4. Objectives

### 1.4.1. General objective

✓ The General objective is to build Amharic sentiment analysis system for the program that assigns positive, incremental positive, negative, decrement negative or neutral polarity using dictionary and rule based approach.

### 1.4.2. Specific objectives

**To achieve the general objective, the following specific objectives are addressed.**

✓ Develop Amharic lexicon dictionaries
✓ To develop Amharic sentiment analysis system prototype.
✓ Transliterating the Amharic text into an equivalent English text.
✓ Develop a dictionary for the negative, Dec. negative, positive and Inc. Positive sentiments.
✓ To Compare each if there is a positive, negative, incremental or decrement in the dictionary
✓ Analysis of the general structure of Amharic statements related to sentiments such as identifying negative, positive and neutral statements.

## 1.5. Significance of the study

"What other people think and feel" has always been an important piece of information for most of us during the decision-making process. [6] A basic task in sentiment analysis is classifying the *polarity* of a given text at the document, sentence, or feature/aspect level—whether the expressed opinion in a document, a sentence or an entity feature/aspect is positive or negative. The significance behind this paper is to identify the polarity of a given text either positive or negative by collecting data from kana TV websites about the program by focusing group discussion and an open ended questioner. And will be able to automatically analyze the sentiment of huge amount of collected reviews prior to making decisions. Therefore, this research can help the Kana TV broadcast corporation to improve their services in the future. The results of the research can be used as an input to the development of full-fledged opinion mining system for Amharic language. Another important significance of the study is that, the output can be used as an input data for recommender and opinion retrieval/search systems. And the system can be used to answer people's opinion or feeling questions. Generally the significance of the study for the society, the researcher and the owners of KANA TV program Managers.

## 1.6. Scope and limitations

The scope of this study is to analysis the sentiment of Kana TV viewers/non-viewers. All sentiments will be incorporated in this study. The paper doesn't consider the total concept of the sentence and document. Comments in English will be converted into Amharic equivalent meaning, but comments in other languages are out of the scope of this study. One of the main limitation is the data set is small.

## 1.7. Methodology

### 1.7.1. Literature Review

Reviewing of similar works (published papers, journal articles and other materials) to develop a deeper emphatic about sentiment analysis using different approaches.

And also literature reviews on different approaches in doing sentiment analysis of different languages, such as rule based, dictionary based, and statistical based approaches etc.

### 1.7.2. Data Collection

Comments from Kana TV website will be collected to analyze the statement of the viewers. In addition, an open ended questionnaire will be conducted to collect the views of Kana TV viewers and non-viewers.

### 1.7.3. Tools

We use python 3.4.3 to analyze the sentiment of Kana TV viewers.

### 1.7.4. Prototyping

To test the proposed system, I have developed a prototype. I construct a lexicon of comments according to the procedures and guidelines for executing the prototype.

## 1.8. Thesis organization

This thesis report is consisting of Introduction, Literature review, under LR related works will be presented, and design and implementation, Experimentation and results and Conclusion and Recommendations. The first chapter gives the general ideas and details about sentiment analysis, levels of sentiment analysis and kana television with sentiment analysis. The second chapter presents literature reviews on different literatures regarding sentiment analysis with its approaches and different natural language techniques as well as previous related works for both English and Non-English documents in varies domains. Chapter three focus on design and implementation the describes analyzing of Amharic language structure such as Amharic writing system, Amharic character characteristics and achieve the goal of the research and also the General system architecture presented in this chapter Defining a structure for the text, pre-processing, normalization, steaming clearly listed. The four chapters discuss the experimentation and results of the findings of how these experiments and methodologies were implemented.

Finally, chapter five deals with the conclusion and the recommendation haggard from the findings of the research.

# CHAPTER TWO

# 2. LITERATURE REVIEW

## *2.1.* Introduction

Sentiment analysis (SA), also known as opinion mining, has attracted an increasing interest. It is a hard challenge for language technologies, and achieving good results is much more difficult than some people think. The task of automatically classifying a text written in a natural language into a positive or negative feeling, opinion or subjectivity [7] is sometimes so complicated that even different human annotators disagree on the classification to be assigned to a given text. Personal interpretation by an individual is different from others, and this is also affected by cultural factors and each person's experience. And the shorter the text, and the worse written, the more difficult the task becomes, as in the case of messages on social networks like Twitter or Face book. [8] Sentiment analysis is becoming a popular area of research and social media analysis, especially around user reviews and tweets. It is a special case of text mining generally focused on identifying opinion polarity, and while it's often not very accurate, it can still be useful.

Sentiment Analysis is one of the interesting applications of text analytics. Although the term is often associated with sentiment classification of documents, generally speaking it refers to the use of text analytics approaches applied to the set of problems related to identifying and extracting subjective material in text sources. [9] It is concerned with the identification of opinions in a text and their classification as positive, negative and neutral. Sentiment analysis refers to a broad area of natural language processing, computational linguistics and text mining that aims to determine the attitude of a speaker or writer with respect to some topic [10]. Defined sentiment mining as a recent disciple at the cross roads of information retrieval, text mining and computational linguistics which tries to detect the opinions expressed in the natural language texts. Sentiment mining is a complex field as it involves the processing and interpretation of natural language. Hence it must deal with natural languages' inherently ambiguous natures, the importance of context, and other complications that do not lend themselves to automation. [11]

The main activities needed for building a sentiment mining system are: development of linguistic resources (e.g. build a lexicon of subjective terms), classification of text (entire documents, sentences) based on their content (e.g. classifying a news article either as positive or negative in relation to the subject), extraction of opinion expression from text, including relations with the rest of content (e.g. recognizing an opinion, who is expressing it, who/what is the target of the opinion), mining tools and visualization tools to extract meaningful information from the mined articles based on the sentiment tags [12].

## 2.2. Types of opinion words

Let us discuss about some types of opinion words. There are two main types of opinion [13]: regular opinions and comparative opinions. Regular opinions are often referred to simply as opinions in the research literature. A comparative opinion expresses a relation of similarities or differences between two or more entities, and/or some of the shared aspects of the entities. A comparative opinion is usually expressed using the superlative or comparative form of an adjective or adverb, although not always. On other hand, regular opinion is a positive or negative sentiment, attitude, emotion or appraisal about an entity or an aspect of the entity from an opinion holder. Positive, negative and neutral are called opinion orientations (also called sentiment orientations, semantic orientations, or polarities). We can also classify opinions based on how they are expressed in text, explicit opinion and implicit (or implied) opinion. An explicit opinion is a subjective statement that gives a regular or comparative opinion. An implicit opinion is an objective statement that implies a regular or comparative opinion. Such an objective statement usually expresses a desirable or undesirable fact. "Explicit opinions are easier to detect and to classify than implicit opinions. Much of the current research has focused on explicit opinions. Relatively less work has been done on implicit opinions [14].

According to [15], there are three types of opinion words. These are personal emotion (e.g. happy, delighted, proud, sad, angry, horrified, etc.), appreciation (flexible, stable, efficient, reduced, ideal, backward, poor, highest etc.) and judgment (e.g. active, decisive, caring, dedicated, intelligent, negligent, evil, etc.).

## 2.3. Sentiment classification or Levels of Analysis

1. document level
2. sentence level
3. Aspect/feature level.

**Document level:** In this level, sentiment is extracted from the entire review, and whole opinion is classified based on the overall sentiment of the opinion holder (or analyzing the overall sentiment expressed in the text). The goal is to classify a comment as positive, negative, or neutral.

**Example:**

"I bought a new model lap top a few days ago. It is such a good lap top, although a little large. The touch screen is adorable. The quality is clear too. I simply like it.

Is the review classification positive or negative? Document level classification works best when the document is written by a single person and expresses an opinion/sentiment on a single entity. [16]

Selama [17] has studied document level sentiment mining for opinionated Amharic text using general and domain specific opinion terms. His model has components such as: pre-processing, sentiment word detection, weight manipulation, polarity classification and polarity strength. According to the researcher, all positive sentiment terms are tagged in the lexica by „+" and given a default value of +2 at run time while all the negative sentiment terms are tagged by „-" and given a default value of -2. Before the final average polarity weight is calculated, the polarity propagation is done which is used to modify the initial value of the sentiment terms. This modification of the initial value or weight is done only if the sentiment word is linked to a modifier term (negations or intensifiers). Accordingly the polarity of the review is determined by the resultant sum of polarity value of opinion terms in the reviews. If the resultant sum is greater than zero then the review is categorized as positive. Similarly, if the resultant sum is less than zero then the review is categorized as negative. Also if the resultant sum is equal to zero then the review is taken as neutral.

According to the study [17], tests on the prototype were done using movie and newspaper reviews where the obtained result with these test data is encouraging.

Another work on document level sentiment analysis is the work in [18] which considered overall sentiment for the classification of a document. The aim of their work was to examine whether it suffices to treat sentiment classification simply as a special case of topic-based categorization with two topics, positive and negative, or whether special sentiment categorization methods need to be developed. They experimented with three standard algorithms: Naïve Bayes, Maximum Entropy, and Support Vector Machines classifications. Having movie reviews as data, they found that standard machine learning techniques definitely outperform human-produced baselines. However, from their experimental result they realized that sentiment categorization using machine learning techniques is more difficult than topic based classification. Work in [19] on the classification of reviews was the closest to their work.

Similarly a document level sentiment analysis which presented a simple unsupervised learning algorithm for classifying reviews as recommended (thumbs up) or not recommended   (thumbs down). The researcher considered an average semantic orientation of the phrases in the review that contain adjectives or adverbs in order to classify a review. [20] According to the researcher, the semantic orientation of a phrase is calculated as the mutual information between phrase and the word excellent minus the mutual information between phrase and the word poor. Accordingly, semantic orientation of the review is recommended if the average semantic orientation of its phrases is positive. The author collected around 410 reviews in four different domains from opinions for the evaluation and found the proposed algorithm achieved an average accuracy of 74%. [21] According to the researcher, two main limitations were observed: required longer time for query processing which can be improved by computer with higher processor, and level of accuracy for movie domain becomes down which can be improved by the combination of proposed algorithm with a supervised classification algorithm.

Work in [22] also studied document level sentiment analysis. They have developed a methodology for extracting small investor sentiment from stock message boards.

They have applied statistical and natural language processing techniques to extract opinion from the boards. According to them, messages are classified into one of the three types: bullish, bearish and neutral by using five classifiers: Naïve classifier, Vector distance classifier, Discriminant-based classifier, Adjective-Adverb phrase classifier, and Bayesian classifier.

The researchers mentioned that some of these classifiers are language dependent while others are independent. From their performance evaluation, their idea was worth pursuing.

**Sentence level:**

This process usually involves two steps:

• Subjectivity classification of a sentence into one of two classes: objective and subjective.
• Sentiment classification of subjective sentences into two classes: positive and negative.

An objective sentence presents some factual information, while a subjective sentence expresses personal feelings, views, emotions, or beliefs. Subjective sentence identification can be achieved through different methods such as Naïve Bayesian classification. However, just knowing that sentences have a positive or negative opinion is not sufficient. This is an intermediate step that helps filter out sentences with no opinions and helps determine to an extent if sentiments about entities and their aspects are positive or negative. A subjective sentence may contain multiple opinions and subjective and factual clauses.

**Example**
"iPhone sales are doing well in this bad economy."

Sentiment classifications at both the document and sentence levels are useful, but they do not find what people like or dislike, nor do they identify opinion targets.

Subjectivity identification is the first task in sentence level opinion mining [23].

They have investigated the idea of creating a subjectivity classifier that uses lists of subjective words learned by two bootstrapping algorithms: Meta-Bootstrapping and Basilik algorithms that were designed to learn words that belong to a semantic category. According to the authors both algorithms need seed words and UN annotated text corpus as input. They have mentioned a relevant annotation scheme that was developed for a U.S government sponsored project with a team of 10 researchers in 2002 for the identification and characterization of subjective words in a sentence. According to them, sentences are labeled by two classifiers: first for subjective sentences and second for objective sentences. The sentences that are not clearly classified into any category are left unlabeled and omitted at this stage. Both of the classifiers are based on preset list of words that indicate sentence subjectivity.

The subjective classifier looks for the presence of words from the list, while the objective classifier tries to locate sentences without those words. According to the results presented by the researchers, their classifiers achieved around 77% recall with 81% precision during the tests. Hatzivassiloglou and Wiebe [24] also studied subjectivity identification at sentence level by considering the impact of adjectives. Semantically oriented adjectives and gradable adjectives were used for subjectivity classification and prediction. Their work was largely exploratory regarding interaction of different characteristics of adjectives in the prediction of subjective sentences. Incorporating their investigation to machine learning models for the prediction of subjectivity was part of their plan.

Another work on sentence level is work of Kim and Hovy [25]. They have presented a system that automatically determines sentiment words within a sentence and combining them by having a pre-defined topic. Their system works in four steps: First it selects sentences that contain both the topic phrase and holder candidates. Next, the holder-based regions of opinion are delimited. Then the sentence opinion classifier calculates the polarity of all opinion bearing words individually. Finally, the system combines them to produce the holder's sentiment for the whole sentence. They experimented with various models of classifying and combining sentiment at word and sentence levels with promising results.

Additional work on sentence level sentiment analysis was the work of Michael Gammon et al [26]. They have developed a prototype system called Pulse which has been used to extract taxonomy of major categories (makes) and minor categories (models) of cars in order to process according to two dimensions of information: sentiment and topic detection from large quantities of car reviews text database. According to them, first the sentences for a make and model of car have been assigned to clusters and have received a sentiment score from the sentiment classifier then visualization component displays the clusters and the keyword labels that were produced for the sentences associated with that car. Their classifier considered three classes: positive, negative and "others". The researchers stated that the technique they have described was simple but effective for clustering sentences and a bootstrapping approach have been applied for sentiment classification. Their experiment indicates that they achieved some of the best results on the negative and "other" classes.

Automatic spelling correction is the future work in addition to identification of sentiment vocabulary and sentiment orientation with minimal customization cost for a new domain.

The current work on sentence level opinion mining is the work of Barbosa and Fang [27].

The researchers have proposed the use of meta-information about the words on tweets and characteristics of how tweets are written in order to classify sentence of the tweets as subjective or objective, and in order to determine the subjective tweets as positive or negative. Their approach uses biased and noisy labels as input to build the models. According to them their solution is more effective and robust while more detailed analysis of sentences is a future work.

**Aspect/Feature level:** in this level, commented features are identified and extracted and the sentiment towards these features is determined from the source data. Sentiment classification at both document level and sentence level are not enough to tell what people like and/or dislike, because a positive opinion on an object does not mean that the opinion holder likes everything similarly a negative opinion on an object does not mean that the opinion holder dislikes everything [28] .

In general the tasks of sentiment mining are: determining document subjectivity, determining document polarity and determining strength of document orientation [29]. Determining document subjectivity: deciding whether a given text has a factual nature or expresses an opinion on its subject matter. This amounts to performing binary text categorization under categories of objective and subjective. Determining document polarity: decides if a given subjective text expresses positive, negative or neutral opinion on its subject matter. Determining strength of document orientation: decides whether the positive opinion expressed by a text on its subject matter is weakly positive, mildly positive or strongly positive.

Similarly decides whether the negative opinion expressed by a text on its subject matter is weakly negative, mildly negative or strongly negative.



Figure 2: Sentiment analysis granularity levels

## 2.4. Steps to Sentiment analysis

The sentiment analysis is a complex process that involves 5 different steps to analyze sentiment data. These steps are: [30]

**Step 1: Data collection:** the first step of sentiment analysis consists of collecting data from user generated content contained in blogs, forums, and social networks. These data are disorganized, expressed in different ways by using different vocabularies, slangs, context of writing etc. Manual analysis is almost impossible. Therefore, text analytics and natural language processing are used to extract and classify;

**Step 2: Text preparation:** consists in cleaning the extracted data before analysis. Non-textual contents and contents that are irrelevant for the analysis are identified and eliminated;

**Step 3: Sentiment detection:** the extracted sentences of the reviews and opinions are examined. Sentences with subjective expressions (opinions, beliefs and views) are retained and sentences with objective communication (facts, factual information) are discarded;

**Step 4: Sentiment classification:** in this step, subjective sentences are classified in positive, negative, good, bad; like, dislike, but classification can be made by using multiple points;

**Step 5: Presentation of output:** the main objective of sentiment analysis is to convert unstructured text into meaningful information. When the analysis is finished, the text results are displayed on graphs like pie chart, bar chart and line graphs. Also time can be analyzed and can be graphically displayed constructing a sentiment time line with the chosen value (frequency, percentages, and averages) over time.

## *2.5.* Approaches/techniques for sentiment analysis

There are a number of different approaches that have been used in an attempt to solve the problem of sentiment classification. One of the most widely used methods involves classifying a single word or phrase with sentiment, and then calculating an overall sentiment rating for a target document using some weighting [31]. The most commonly applied techniques for sentiment analysis are listed below.

- Machine Learning Techniques
- Ontology Based Techniques
- Lexicon Based Techniques
- Natural Language Processing Techniques
- Linguistic Techniques
- Lexicon Based Approach and ETC.

From our point of view or study natural language processing techniques selected for the Amharic sentiment analysis in order to assign polarity of comments that is given by the viewers of KANA TV program. The related approaches describes here:

## 2.5.1. Natural language processing technique

There are different language analysis techniques that fall under the sunshade of natural language processing, of which the most common are: part of speech (POS) tagging, co-reference resolution and full syntactic parse tree. POS tagging is the process of labeling word occurrences with its world class; for example, whether a word is occurring as an adjective, noun, or a verb. Effective tagging requires knowledge of not just the word but also its context, such as position within the sentences and surrounding words. Hidden Markov model is a common technique which is used in POS tagging [32]. To avoid constantly referring a subject by name, natural language usually contains alternative words that can be used when referring to a previously mentioned subject. Co-reference resolution is used for automating the process of connecting such references. Creating parse tree for natural languages is another central area of study in NLP. Parsing is related to POS tagging as determining the sentence structure requires knowledge of which sense words are being used.

## 2.5.2. Lexicon Based Approach

Opinion words are employed in many sentiment classification tasks. Positive opinion words are used to express some desired state, while negative opinion words are used to express some undesired states. There are also opinion phrases and idioms which together are called opinion lexicon. Opinion words are also known as polar words, opinion-bearing words and sentiment words. The lexicon based techniques to Sentiment analysis is unsupervised learning because it does not require prior training in order to classify the data.

To compile or collect the opinion words list, three main approaches have been investigated: *Manual approach, Dictionary based approach* and *Corpus-based approach* [33]*.*

**Manual Approach**

This approach is just a process of hand picking sentiment words from different sources with the goal of populating a lexicon with polar words. This manual approach is very time consuming and it is usually combined with automated approaches as the final check because automated methods make mistakes. The opinion words lexicon used in [34] are manually handpicked based on a reading of several thousand messages. In the work of J. Yi et al. [35] , they collected a sentiment lexicon of 3000 English sentiment terms manually from different sources.

**Dictionary Based Approaches**

[36] Presented the main strategy of the dictionary-based approach. A small set of opinion words is collected manually with known orientations. Then, this set is grown by searching in the well-known corpora WordNet [37] or thesaurus [38] for their synonyms and antonyms. The newly found words are added to the seed list then the next iteration starts. The iterative process stops when no new words are found. After the process is completed, manual inspection can be carried out to remove or correct errors.

The dictionary based approach has a major disadvantage which is the inability to find opinion words with domain and context specific orientations. Qi and He [39] used dictionary-based approach to identify sentiment sentences in contextual advertising. They proposed an advertising strategy to improve ad relevance and user experience.

They used syntactic parsing and sentiment dictionary and proposed a rule based approach to tackle topic word extraction and consumers' attitude identification in advertising keyword extraction. They worked on web forums from automotvieforums.com. Their results demonstrated the effectiveness of the proposed approach on advertising keyword extraction and ad selection.

**Corpus based approaches**

The Corpus-based approach helps to solve the problem of finding opinion words with context specific orientations.

Its methods depend on syntactic patterns or patterns that occur together along with a seed list of opinion words to find other opinion words in a large corpus. One of these methods was represented by Hatzivassiloglou and McKeown [40] . They started with a list of seed opinion adjectives, and used them along with a set of linguistic constraints to identify additional adjective opinion words and their orientations. The constraints are for connectives like AND, OR, BUT, EITHER-OR….; the conjunction AND for example says that conjoined adjectives usually have the same orientation. This idea is called sentiment consistency, which is not always consistent practically. There are also adversative expressions such as but, however which are indicated as opinion changes. In order to determine if two conjoined adjectives are of the same or different orientations, learning is applied to a large corpus. Then, the links between adjectives form a graph and clustering is performed on the graph to produce two sets of words: positive and negative.

The corpus-based approach is performed using statistical approach or semantic approach as illustrated in the following subsections:

## Statistical approach

Finding co-occurrence patterns or seed opinion words can be done using statistical techniques.

This could be done by deriving posterior polarities using the co-occurrence of adjectives in a corpus, as proposed by Fahrni and Klenner [41]. It is possible to use the entire set of indexed documents on the web as the corpus for the dictionary construction. This overcomes the problem of the unavailability of some words if the used corpus is not large enough [42] .

## Semantic approach

The Semantic approach gives sentiment values directly and relies on different principles for computing the similarity between words. This principle gives similar sentiment values to semantically close words. WordNet for example provides different kinds of semantic relationships between words used to calculate sentiment polarities.

WordNet could be used too for obtaining a list of sentiment words by iteratively expanding the initial set with synonyms and antonyms and then determining the sentiment polarity for an unknown word by the relative count of positive and negative synonyms of this word [43] .



Figure 3: General Sentiment classification techniques [44]

## 2.6. RELATED WORKS

In the field of natural language processing Sentiment analysis is one of the most widely held research areas. Many Researchers have been conducted in the area of sentiment analysis. Most researches of sentiment analysis are directed for English language. Some of sentiment analyses for other languages. Related researches done for different language such as English, French, Chinese, Indian and other using different techniques and approaches are reviewed. Different authors used different techniques such as machine learning, ontology based approaches, lexicon-based approaches and others.

Sandeep Balijepalli [45] used machine learning technique to categorize opinionated English documents taken from political blogs based on their sentiments and determine the polarity strength of the sentiments. Contents collected from the political domain are made to pass through pattern matching (Nave Bayes filter, bag of words and part of speech tagging) for obtaining the sentiment oriented sentences which are later to be indexed. The index helps to avoid the delays in fetching the data.

The framework proposed by Sandeep gets contents from the database of blogs, pass the sentences to the sentence chunkier for stripping unrelated data, then the sentence is passed through filters for filtering out objective sentences and classifying subjective sentences, index opinionated sentences, divide results by bloggers party and finally sort them by their polarity strength. In the above approach, filter analysis is done by making sentences to pass through the pattern recognizer first for checking the sentences if they follow the custom developed subjective pattern. If the sentence matches the pattern, it is indexed otherwise it is passed through Nave Bayes (unigram, bigram) for further analysis where this filter depends on the training dataset. If the sentence is not indexed at this filter, it again passes through the part of speech tagging and if the sentence is found to be subjective, it is indexed otherwise it is considered as objective sentence and it is skipped. Then, other sentence undergoes the entire procedure. The experimental result shows that the system performs well with unigram approach.

Alistair Kennedy and Diana Inkpen [46] proposed a method that counts positive and negative terms but also takes contextual valence shifters such as negations and intensifiers into account.

Two approaches are compared in their work. The first approach simply counts positive and negative terms where the review is positive if the review contains more positive than negative terms. Review is negative if it contains more negative than positive terms. A review is neutral if it contains equal number of positive and negative terms. The term counting method can be easily modified to use valence shifters. The second method counts positive and negative terms, but takes contextual valence shifters into account. Their approaches are classified as basic (uses the first approach) and improved one (uses the second approach).

The main lexicon used in this work was the General Inquirer (GI) though they added extra terms from other sources. As the authors stated, their motivation to use this approach was to see the effect of incorporating contextual valence shifters to the basic method of sentiment classification. The data sets they for experimental purpose are taken from two sources. The first data set is taken from www.epinions.com [47]. Epinions.com is a general consumer review site. The data set taken contains 70 positive and 70 negative reviews. The reviews were collected from a variety of different products, including air conditioners, sewing machine, vacuums cleaners, TVs, cookware, beer and wine. The second data set is a movie reviews that contains 2000 reviews, 1000 positive and 1000 negative taken from other movie review sources.

The experimental result indicated that the proposed approaches perform well as indicated in the following. The basic approach using GI lexicon gives an accuracy of 0.679 for product reviews and 0.595 for movie reviews. The improved method using GI lexicon gives an accuracy of 0.686 for product reviews and 0.627 for movie reviews. The experimental results of adding extra terms from other resources to GI are also given in their work and some improvements are shown. In most cases the method of classification performs better when classifying product reviews than movie reviews. This is because movie reviews are known to be more difficult to classify than other reviews such as product reviews [48] .

Lili Zhao and Chunping Li [49] used ontology based opinion mining for movie reviews with the goal of improving feature level opinion mining by employing ontology. The use of this approach was motivated by the role of ontology in conceptualizing domain specific information.

The main components of the proposed approach are: text collection (movie reviews), preprocessing, feature identification, polarity identification and sentiment analysis with the support of ontology development. Like others the polarity identification fully relies on a lexicon of tagged positive and negative sentiment terms which are used to quantify positive/negative sentiment. For this purpose SntiWN [50]was used as it provides a readily interpretable positive and negative polarity values for a set of 'affective' terms.

 The target of the ontology development is to define common terminologies in the area, and give the definition of the relationship among the terminologies. Iterative approach is used for developing the ontology following two steps. The first step is selecting the relevant sentences including concepts and the second step is extracting the concepts from those sentences. Criteria used for selecting the sentences are: the sentences that contain conjunction word and sentences that contain at least one concept seed. At the initial state, manually labeled feature are used as seeds. Randomly selected 1400 movie reviews from internet movie database (IMDB) [51] were used as dataset where half of them are positive and the other half are negative.  The experimental results indicate that the accuracy is satisfying, and proves that it is reasonable to compute the polarity score by the proposed method where the main factor is found to be the ontology structure. Even though this work is ontology based, it depends on the lexicon of opinion terms for assigning weights to the sentiment terms. Ontology is basically needed for feature extraction.

Xiaoying Xu et al. [52] on their work titled "categorizing term's subjectivity and polarity manually for opinion mining", proposed principles and guidelines to create a large-scale Chinese sentiment lexicon for opinion mining manually. Two experiments are conducted in their work: the first experiment is conducted to investigate the reliability of manual subjectivity labeling of the terms. The second experiment is conducted to see the effectiveness of the lexicon in judging the polarity of subjective sentences.

In this paper, it is indicated that for establishing the first and large scale human tagging Chinese sentiment lexicon, the agreement of different annotators and the reliability in sentence polarity judging system are key issues. As a result annotation principles and guidelines are needed to be established. The principles in tagging the terms subjectivity and polarity established by the authors are: the terms should be opinion mining oriented, the lexicon built will be used only in the subjective sentence in opinion mining and the word will indicate polarity in the subjective sentence. A clear guideline in annotating the sentiment word and qualified annotators are needed for realizing the principles of building the lexicon. During the tagging campaign and lexicon building, the main resource used was HowNet [53]. HowNet is an on-line common-sense knowledge base unveiling inter-conceptual relations and inter attribute relations of concepts as connoting in lexicons of the Chinese and their English equivalent. According to their analysis the answer they gave to the question "what kind of word can be selected in the lexicon?" is 'if a term has subjective meaning either in concept meaning or in emotion meaning overtone, and it can indicate the polarity in subjective sentence, it must be selected in our lexicon.

The experimental results for the first experiment show that the polarity of word sense can be reliability annotated in despite of the polarity ambiguous in words is common Chinese, and also because of its large-scale it could be very useful fundamental resource in opinion mining and other related fields. To evaluate the lexicon in real applications (the second experiment), they built simple sentence sentiment recognition system that contains text pre-processing, polarity words detection and weight assign, link construction and polarity propagation. The text pre-processing is used to segment words and tag POS. In the polarity words detection, every polarity word is checked whether it is polarity word defined in the sentiment lexicon and get the corresponding sentiment polarity if found. Every polarity word and modifier word get the initial weight defined in the sentiment lexicon. If the polarity word is linked to a modifier word, the polarity value should be multiplied by a coefficient in the polarity propagation step.

The results of the second experiment show that using the sentiment lexicon it can achieve an accuracy of more than 70%.

Sigrid Maurel et al. [54] used a combined approach (combination of symbolic and statistical) for the classification of opinionated texts in the French language. The symbolic approach includes systems for extracting information adapted to the corpora based on the rules of syntactic and semantic analyzer. This approach analyzes texts sentence by sentence and extracts relationships that convey feelings. Statistical method is based on machine learning techniques. It process text in a single step and assigns a global opinion at the whole text at the end. The hybrid approach is used in their work to increase the quality of the results. As they indicated the experimental results show that combination of statistical and symbolic (hybrid) approaches gives more accurate results than either method used separately. In the ground of NLP Sentiment analysis as we give clarity about this field is the one and the most widely held research areas. Many Researchers conducted in the area of SA. Most researches of sentiment analysis are focused for English language as we indicated in the portion of related works above. Few Related researches done for Amharic languages:

Salama [55] design sentiment mining model for opinionated Amharic texts in movie domain in which reviews manually collected from fans by distributing questionnaires. The researcher Used Lexicon based approach and the proposed model have the following components: preprocessing, sentiment words detection, weight assignment and propagation, polarity classification, polarity strength representation and sentiment lexica. After the reviews are preprocessed, each term is checked for existence in the sentiment lexica at the sentiment words detection component. The detected sentiment terms are assigned weight and the values of sentiment terms that are linked to contextual valence shifters are propagated in the weight assign and polarity propagation component. Based on the weights of the sentiment values, the reviews are classified into predefined categories: positive, negative or neutral. Finally, the polarity strength of the reviews is rated. The sentiment lexica are built manually from different sources based on the principles and guidelines. This lexicon is a collection of predefined sentiment words, they are words that express an opinion toward an object such as good that express positive opinion and bad which express a negative opinion, which are manually collected and stored in the dictionary. After having the lexicon, the review document is preprocessed and every valid term in the review is checked whether it is a sentiment word or not.

This is done by a detection mechanism where the whole lexicon is scanned for every term. If the term exists in the dictionary, then the term is a polarity word (positive or negative). Terms in the dictionary are tagged in the lexicon with a computer interpretable value as "+" for positive opinion terms and "-" for negative opinion terms. Based on this procedure, given a review document, if there are more positive terms than negative then it is considered to be positive. If there are more negative than positive terms then it is considered to be negative. If there are equal numbers of positive and negative terms then it is neutral.

Final experimental results found in movie review in [56] are shown in the Table 1 below.

Table 1: Experimental results found in research

| System | Class | Precision | Recall | F-measure |
|---|---|---|---|---|
| General purpose Amharic sentiment terms(Basic system) | Positive | 0.929 | 0.823 | 0.867 |
| Negative | 0.6 | 0.573 | 0.589 | |
| Basic + domain lexicon | Positive | 0.937 | 0.943 | 0.939 |
| Negative | 0.62 | 0.78 | 0.69 | |
| Both lexica + contextual valance shifter terms | Positive | 0.943 | 0.949 | 0.945 |
| Negative | 0.666 | 0.842 | 0.743 | |

The author showed result obtained in using single general purpose dictionary and using two sentiment lexica: the general purpose lexicon and the domain specific lexicon. And finally, the result of the experiment conducted using the two lexica and considering the contextual valence shifter terms.

Mengistu [57] also has done a research on Sentiment analysis for classifying Amharic opinionated text in to positive, negative or neutral by using ML approach in ERTV, Fana broadcasting and diretube.com domains.

The author employed three machine learning classification techniques (Naive Bayes, Multinomial Naïve Bayes and Support vector machines) using n-grams presence, n-grams frequency and n-grams-TF-IDF features selection methods. The experiments are conducted using 576 Amharic opinionated texts collected from ERTA, Fana Broadcasting and diretube.com manually. The Experiment indicates that uni-grams –term frequency feature selection methods perform the best for all algorithms (Support Vector machine, Naïve Bayes and multinomial Naïve Bayes). Based on their relative performance of classification, support vector machine registers with 78.8% accuracy outperform, Naïve Bayes with 77.6% and multinomial Naive Bayes with 74.7%.as shown from the result obtained SVM performed better than NB and MNB algorithm.

Tilahun has done the study with the title of opinion mining from Amharic blog using rule based approach. He used feature level opinion mining and summarization techniques are evaluated on 484 reviews manually collected from hotel, university and hospital domain.

Getachew works on the title of opinion mining from Amharic entertainment text using machine learning approaches (Naïve Bayes, Decision Tree and Maximum Entropy).the experiment conducted using 616 Amharic optioned texts. The study obtained 90.9 %, 83.1% and 89.6% using Naïve, Bayes, Decision Tree and Maximum entropy algorithms respectively. However, the study did not control negation, because the study uses uni-gram as a feature for classification. The result only shows positive and negative polarity but it did not include incremental positive and decremental negative. In a summary, one of the research work related to this study is Selama's work [55].He has done Sentiment analysis on Amharic opinionated text using rule based approach on movie review. To classify the opinionated text, the researcher has built a lexicon that contains information about Amharic word senses. All the words which are not in the list of the dictionary are ignored then the class will be unclassified therefore there is a need to change another method to solve this problem.

Summary of some of the related works in terms of their objective/goal, Methods used, Data source and result found in remark are summarized below in Table 2.

Table 2: Overview of some previous work

| Author | Objectives/goals | Methods/ techniques | Data Resource /domain | Result |
|---|---|---|---|---|
| (Selama, 2010) | Design and develop a sentiment mining model for opinionated Amharic documents. | Lexicon/dictionary based | Movie review | Better result found in using both basic lexica and domain lexica with contextual valance shifter terms |
| (Mengistu,2013) | Explore the possibility of applying Sentiment analysis and build a model. | Supervised machine learning (NB,MNB,SVM) and n-grams presence, n grams frequency and n-grams-TF-IDF techniques for feature selection | entertainme nt reviews in ERTA and Fana Broadcasti ng Corporate program | Support Vector machine achieves the best result using uni grams term frequency. |
| (Tulu Tilahun, *2013)* | Develop feature level opinion mining and summarization model for Amharic language. | rule based approach | Amharic blog | Two experiments are conducted. |
| (Abreha m Getache w ,2014) | Appling opinion mining to create classification model for the Amharic entertainment reviews. | Machine Learning Approach (Naïve Bayes, Decision Tree and Maximum Entropy) | Entertainm ent review | By combining the two methods they are able to improve the results over either of the method alone. |

## 2.7. Summary

Here in this chapter reviewed different research attempts to solve the problem of sentiment analysis for different languages.in this chapter mainly showed the classification of sentiment analysis in two parts which are machine learning approach and lexicon based approaches. And also the review showed that Machine Learning Techniques, Ontology Based Techniques, Lexicon Based Techniques, Natural Language Processing Techniques, Linguistic Techniques, Lexicon Based are the commonly used approaches to deal with sentiment analysis. To accumulate or gather the opinion words list, three main approaches have been investigated in this chapter: which are Manual approaches, Dictionary based approach and Corpus-based approach. Each approaches given detail explanation. In the manual approaches a process of hand picking sentiment words from different sources with the goal of populating a lexicon with polar words and clearly defined as very time consuming and it is usually combined with automated approaches as the final check because automated methods make mistakes. In the dictionary based approaches showed the advantages like to identify sentiment sentences in contextual advertising. The last lexicon based approaches is the Corpus-based approach which helps to solve the problem of finding opinion words with context specific orientations as defined in this chapter.

# CHAPTER THREE

## 3. DESIGN AND IMPLEMENTATION

In this research focused on detecting the overall polarity of kana TV program viewer comments and to compute the polarity (positive, negative or neutral) with the incremental, and decrement of words. Now here in this research we are going to use an approach based on dictionaries and rule based (hybrid) Approaches. A dictionary is no more than a list of words that share a category. For example, you can have a dictionary for positive and negative expressions. In our approach we need to process the input text, splitting, stemming or lemmatizing, and extracting information from it. Tokenization and stemming for the source material written in Amharic should be done before we compare each lexicon with positive and negative dictionaries. Stemming should be done to disambiguate the root word that should be compared with the positive and the negative dictionaries. The design of the dictionaries highly depends on the concrete topic where you want to perform the sentiment analysis.

### 3.1.    System architecture

The general architecture of the proposed system of sentiment analysis for kana TV broadcast is shown in figure 3.1. The system contains different components based on the processes required. These components are: pre-processing (Tokenization, Normalization, and Stemming), sentiment sentence/document analysis and sentiment scores. If the cumulative result is greater than zero the sentence/document is positive and if the cumulative result is less than zero the text is negative and if it is zero called neutral. The sentiment lexicon is also part of the general systems architecture. The input to the system is a corpus of documents in text format. The documents in this corpus are converted to text and pre-processed using a variety of linguistic tools such as stemming, tokenization…etc. The system may also utilize a set of lexicons and linguistic resources. The main component of the system is the sentence/document analysis module, which utilizes the linguistic resources to annotate the pre-processed documents with sentiment annotations. The annotations may be attached to whole documents (for document-based sentiment), to individual sentences (for sentence-based sentiment).

These annotations are the output of the system and they may be presented to the user using a variety of visualization tools. The General architecture of sentiment analysis is shown below in Figure 4.

Figure 4: Architecture of Sentiment Analysis

The proposed system has the following components which are listed below:

### 3.1.1.    Defining a structure for the text

Defining the text simply as a list of words. Or define a more elaborated structure carrying every possible attribute of a processed text (word lemmas, word forms, multiple tagging, inflections...) [58] .

***Structures to implement the polarity will be:***

- . Each input text to a list of sentences
- . Each input sentence to a list of tokens
- . Each input token to a tuple of three elements:
    - *A word form* (the exact word that appeared in the text)
    - *A word lemma* (a generalized version of the word) and
    - A list of associated tags

### 3.1.2.    Pre-processing the text

Once we have decided the Structures to implement the polarity of your processed text, we can pre-process this structured text. With pre-process we mean some common first steps in NLP such as: Tokenize, Split into sentences…etc.

- ▪ We will use the NLTK library for these tasks:

```
import nltk
class Splitter(object):
    def __init__(self):
        self.nltk_splitter = nltk.data.load('tokenizers/punkt/english.pickle')
        self.nltk_tokenizer = nltk.tokenize.TreebankWordTokenizer()
    def split(self, text):
        """

        Input format: a paragraph of text
        Output format: a list of lists of words.
```

```python
            """
            sentences = self.nltk_splitter.tokenize(text)
            tokenized_sentences = [self.nltk_tokenizer.tokenize(sent) for sent in sentences]
            return tokenized_sentences
class POSTagger(object):
    def__init__(self):
        pass
    def pos_tag(self, sentences):
        """
        Input format: list of lists of words
        Output format: list of lists of tagged tokens. Each tagged tokens has a
        form, a lemma, and a list of tags
        """
        pos = [nltk.pos_tag(sentence) for sentence in sentences]
        #adapt format
        pos = [[(word, word, [postag]) for (word, postag) in sentence] for sentence in pos]
        return pos
```

- Now, using this two simple wrapper classes, can perform a basic text pre-processing, where the input is the text as a string and the output is a collection of sentences, each of which is again a collection of tokens. By the moment, our tokens are quite simple. Since we are using NLTK, and it does not lemmatize words, our forms and lemmas will be always identical. At this point of the process, the only tag associated to each word is its own POS Tag provided by NLTK.

```python
text = """"""
splitter = Splitter()

postagger = POSTagger()
splitted_sentences = splitter.split(text)
print splitted_sentences
pos_tagged_sentences = postagger.pos_tag(splitted_sentences)

print pos_tagged_sentences
```

- Defining a dictionary of positive, negative, Increments, and decremental expressions.

A dictionary is no more than a list of words that share a category. For example, you can have a dictionary for positive expressions. The design of the dictionaries highly depends on the concrete topic where you want to perform the sentiment analysis.

The next step is to recognize positive, negative, Increments and decrement expressions. To achieve this, we are going to use dictionaries, i.e. simple files containing expressions that will be searched in our text.

Sample data will be listed below:

<u>positive.yml</u>

ጥሩ: [positive]

ማር: [positive]

ደስ: [positive]

ወጻብ: [positive]

<u>negative.yml</u>

መጥፎ: [negative]

አለመርካት: [negative]

ወጻድ: [negative]

አሳዛኝ: [negative]

አልመክርም: [negative]

<u>inc.yml</u>

በጣም: [inc]

እጅግ: [inc]

<u>dec.yml</u>

ጥቂት: [dec]

ያነስ: [dec]

አጥረት: [dec]

- Tagging the text with dictionaries

The following code defines a class that we will use to tag our pre-processed text with our just

defined dictionaries

```
class DictionaryTagger(object):
    def __init__(self, dictionary_paths):
        files = [open(path, 'r') for path in dictionary_paths]
        dictionaries = [yaml.load(dict_file) for dict_file in files]
        map(lambda x: x.close(), files)
        self.dictionary = {}
        self.max_key_size = 0
        for curr_dict in dictionaries:
            for key in curr_dict:
                if key in self.dictionary:
                    self.dictionary[key].extend(curr_dict[key])
                else:
                    self.dictionary[key] = curr_dict[key]
                    self.max_key_size = max(self.max_key_size, len(key))
    def tag(self, postagged_sentences):
        return [self.tag_sentence(sentence) for sentence in postagged_sentences]
    def tag_sentence(self, sentence, tag_with_lemmas=False):
        """
        """
        tag_sentence = []
        N = len(sentence)
        if self.max_key_size == 0:
            self.max_key_size = N
        i = 0
        while (i <N):
            j = min(i + self.max_key_size, N) #avoid overflow
```

```python
                tagged = False
                while (j >i):
                    expression_form = ' '.join([word[0] for word in sentence[i:j]]).lower()
                    expression_lemma = ' '.join([word[1] for word in sentence[i:j]]).lower()
                    if tag_with_lemmas:
                        literal = expression_lemma
                    else:
                        literal = expression_form
                    if literal in self.dictionary:
                        is_single_token = j - i == 1
                        original_position = i
                        i = j
                        taggings = [tag for tag in self.dictionary[literal]]
                        tagged_expression = (expression_form, expression_lemma, taggings)
                        if is_single_token: #
                            original_token_tagging = sentence[original_position][2]
                            tagged_expression[2].extend(original_token_tagging)
                        tag_sentence.append(tagged_expression)
                        tagged = True
                    else:
                        j = j - 1
            if not tagged:
                tag_sentence.append(sentence[i])
                i += 1
    return tag_sentence
```

- To compute the sentiment score of a sentence we are using recursive function
  sentence_score:

```python
def sentence_score(sentence_tokens, previous_token, acum_score):
    if not sentence_tokens:
        return acum_score
    else:
```

```
    current_token = sentence_tokens[0]
  tags = current_token[2]
  token_score = sum([value_of(tag) for tag in tags])
  if previous_token is not None:
     previous_tags = previous_token[2]
     if 'inc' in previous_tags:
        token_score *= 2.0
     elif 'dec' in previous_tags:
        token_score /= 2.0
 return sentence_score(sentence_tokens[1:], current_token, acum_score +
token_score)
```

- Next, we could adapt our sentiment_score function. We want it to flip the polarity of a sentiment word when is preceded by an inverter:

```
def sentence_score(sentence_tokens, previous_token, acum_score):
   if not sentence_tokens:
       return acum_score
  else:
       current_token = sentence_tokens[0]
       tags = current_token[2]
       token_score = sum([value_of(tag) for tag in tags])
       if previous_token is not None:
          previous_tags = previous_token[2]
          if 'inc' in previous_tags:
             token_score *= 2.0
          elif 'dec' in previous_tags:
             token_score /= 2.0
          elif 'inv' in previous_tags:
             token_score *= -1.0
     return sentence_score (sentence_tokens [1:], current_token, acum_score +
token_score)
```

### 3.1.3.     Tokenization

Tokenization refers to the process of splitting the text into a set of tokens (usually words). This process detects the boundaries of a written text. The Amharic language uses a number of punctuation marks which demarcate words in a stream of characters which include 'huletneTb' (፡), 'aratneTb'(፡ ፡), 'deribsereze' (፤), 'netelaserez'(፣), exclamation mark '!' and question mark'?' These punctuation marks don't have any relevance in opinion mining task and have to be removed. Review document divided into sentences and a sentence further divided into words. In our case the smallest piece or token is known as terms. The tokenization algorithm is shown in Algorithm 1 below.

1. Open the corpus

2. While not end of file is reached do

      For each character in the corpus

          If the character is Amharic word delimiters then

             Remove the character

          End if

       End for

3. End while

4. Close files

**Algorithm 1:** Tokenization Algorithm

### 3.1.4.     Normalization

Amharic writing system has homophone characters which mean characters with the same sound have different symbols for example; it is common that the character ስ and ሥ are used interchangeably as ስራ and ሥራ to mean "work". These different symbols must be considered as similar because they do not have effect on meaning .Such type of inconsistency in writing words will be handled by replacing characters of the same sound by a common symbol. Thus, for example, if the character was one of ሐ፣ሓ፣ሃ፣ኀ፣ኃ or ኸ (all of them with a similar sound h) then it was converted to ሀ.

By the same symbol, all orders of ሠ (with the sound s) were changed to their equivalent respective orders of ሰ; all orders of ፀ (with the sound tse) were changed to their equivalent respective orders of ጸ. Sample of such normalized Amharic words are shown in Table 4. Such inconsistency of characters may cause unnecessary increase in the number of document representative words that causes large data size processing. In the table below examples of the Amharic character redundancy where more than one symbol is used for a given sound is shown.

Table 3: Example of redundant Amharic characters with the same sound

| consonants | Other symbols with the same sound |
|---|---|
| ሀ(hä::) | ኀ ሐ ሓ and ኃ |
| ሰ(sä::) | ሠ |
| አ(ä::) | ኣ ዐ and ዓ |
| ጸ(tsä::) | ፀ |

This variation in Amharic language spellings would unnecessarily increase the number of words representing a document which could reduce the efficiency and accuracy of the classifiers. Amharic document processing for feature selection should therefore normalize word variants (spelling differences) caused by inconsistent usage of redundant characters. During the preprocessing activity of Amharic documents for the research, the different forms of a character that have the same sound are changed to one common form. The table below show that some of the redundant words having different word spellings.

Table 4: Example of redundant words which have different word spellings

| words in English | Words in Amharic | Redundant words in Amharic in the given English words |
|---|---|---|
| Pray | ጸሎት | ፀሎት |
| Religion | ሃይማኖት | ሀይማኖት፤ ሐይማኖት፤ ኃይማኖት |
| Holy water | ፀበል | ጸበል |
| Name | ስም | ሥም |
| Sun | ፀሐይ | ጸሎት |
| Life | ሐይወት | ሀይወት |

The algorithm for normalization is shown in Algorithm 2 below.

1.  *Open the corpus*
2.  *While not end of corpus file do*

      *For each character in the corpus*

          *If the character is ሐ፣ኅ or any order then*

             *Changed to ህ*

          *Else if it is ሠ or any order of it then*

             *Changed to ስ*

          *Else if it is or any other order of it then*

             *Changed it to θ*

          *End if*

      *End for*

3.  *End while*

**Algorithm 2:** Normalization algorithm

## 3.1.5.    Stemming

Most of the time words that appear in a document have many morphological variations. The process of stemming is an attempt to reduce a word to its stem or root form. Most often such common variants happened due to suffixing and prefixing. Stemming will bring the different forms of the word into common forms .For example the word 'comfort' (አልተመለሰም, ተመልሰናል) is converted to a stem word "መለስ". Thus, terms of a document are represented by stem words rather than by the original words. This also reduces the number of different terms needed for representing a document and also saves storage space and processing time.

The algorithm for Stemming is shown in Algorithm 3 below.

1. *Open corpus and exception list*
2. *While not end of corpus file is reached do*

   *For each term in the corpus*

   *If term starts with prefix*

   *If term not in exception list then*

   *Remove prefix*

   *End if*

   *End if*

   *If term ends with suffix*

   *If term not in exception list then*

   *Remove suffix*

   *End if*

   *End if*

   *End for*
3. *End while*
4. *Close files*

**Algorithm 3:** Stemmer Algorithm

## 3.1.6.    Amharic text transliteration

Next to pre-processing have done, the Amharic text need to be transliterated in order to make the text well-matched with the machine learning tools in my case Natural language processing toolkit is selected for the experiments. Transliteration is the representation of the characters of one language by matching characters of another language thus transliteration task was accomplished from Amharic to Latin characters. In this research Latin alphabets are used for transliteration. It enables easy, unambiguous and consistent communication of texts.

The transliteration of the Amharic documents was conducted by using SERA [59]. SERA (System for Ethiopic Representation in ASCII) are a convention for the transcription of Fidel (Ethiopic script) into the seven bit ASCII format.

Table 5: Translate from Amharic to Latin Scripts

| Amharic Words | Latin Script |
|---|---|
| አልመክርም | Ālimekirimi |
| አሳዛኝ | Āsazanyi |
| ጥሩ | t'iru |
| አስቀያሚ | āsik'eyamī |
| አስገራሚ | Āsigeramī |

In this research, the researcher transliterated Amharic texts into Latin using SERA representation python programming tool.

In general the steps to decide the polarity of the given text will be:

- Transliterating the Amharic text into an equivalent English text.
- Stemming will be done on every word of the input word
- Identify each disambiguated lexicons for each input word
- Compare each if there is a positive , negative, incremental  or decrement in the dictionary
- Calculate the scores, output the sentiment score!

# 4. Experiment and results

Here in this study we are using the hybrid approaches which are rule based and dictionary based approaches to identifies the comments either good (positive) or not good (negative).

## 4.1. Data collection and preparation

The data is collected from kana broadcasting viewers which are given comments from kana TV face book, website, and blogs and during focus group discussion and by distributing an open ended questionnaires. From all data collection method we have sample size 1022 comments given as an input for the system.

After the data is collected, preprocessing tasks were applied to construct the final data set (data that are used as input for). Data preparation tasks are usually performed multiple times depending on the quality and size of the initial dataset. A task mainly includes normalization; tokenization and stemming of the data were performed to come about with the final appropriate dataset for the selected algorithms. Finally the Amharic text (data set) converted or transliterated into Latin equivalent sentences).

## 4.2. Data analysis

The viewers of kana television are considered as over million, due to getting information from different social media and other methodology. Such as the website of kana, Facebook, Focusing group discussion, and questionnaire. The target groups, which means collected comments are university students of selam campus (ሠላም ካምፓስ) and youth from Bahir dar sefen selam kebele (ሰፈነ ሠላም) and yetebabrute condominium area (የተባበሩት ኮንዶሚኒየም አካባቢ). They are good for our study or research because youth spent much time from kana program. Specially, an open ended questionnaire distributed among this people (viewers/non-viewers of the program) and focusing group discussion have been done here in this area. We are used to collect comments from kana website page and Facebook by using video screen shots.

From this sample screen shot comments and likes over one million users are using the Facebook page of kana TV and they like the page by giving comments as shown below.





Figure 5: Sample face book comments by video screen shot

44

How much is positive, incremental positive, negative, decrement negative or neutral comments?          (Sample size determination formula by yamana)

$$n \quad = \quad \frac{N}{1+N\,(e)^{2}}$$

**Where n =** Sample Size                                        Level of tolerance=5%

**N=** Total population                                        Confidence=95%

**E=** level of tolerance (margin of error)

For our study the total population is around 11,270. From the above proved formula the total sample size is about 1022, which means 1022 comments represent the total population 11,270.

## 4.3.    Performance measurement

We used the following performance metrics to evaluate our experiment results and our classifiers.

a. *Confusion Matrix*

A confusion matrix is used as a form of visualizing the performance of a classifier. It is displayed in a table format in which the columns represent the actual values (true and false) and the rows represent the predicted values (positive, negative and neutral). It can easily be generalized for multi-class classifiers.

The table reports the results of a classifier in terms of the number of true positives (Tp), false positives (Fp), false negatives (Fn), true negatives (Tn), false neutral (Fneu), true neutral (Tneu).

Table 6: Confusion matrix

| Labeled | Truth | |
|---|---|---|
| | Tp | Fp |
| | Fn | Tn |
| | Fneu | Tneu |

Positive had Tp documents correctly classified and Fp and Fneu Documents wrongly classified under it, and class negative Tn documents correctly classified under it Fp and Fneu documents wrongly classified under it, while class neutral had Tneu documents correctly classified under it Fp and Fn wrongly classified under it. Then, in relation to class positive: Tp Fp Fneu Tn Fn Fneu Precision is the ratio of number of documents correctly classified under class positive to the number of all documents classified under class positive:

### b. Accuracy, Precision, Recall, F-Score

We use the following performance metrics to evaluate our experiment result and our F-classifiers accuracy, precision, recall, f-score.

Accuracy is defined as the percentage of positive predictions that are correct.

$$Accuracy = \frac{tp+tn}{tp+fp+tn+fn}$$

Precision is defined as the percentage of positive predictions that are correct.

$$Precision = \frac{tp}{tp+fp}$$

Recall is defined as the percentage of actual positives that are labeled as positive.

$$Recall = \frac{tp}{tp+fn}$$

F-score is the harmonic mean of precision and recall.

$$\text{F-score} \quad = \quad 2_x \quad \frac{\text{Precision x recall}}{\text{Precision + recall}}$$

## 4.4. Results with experiments

The result of this study based on four target area which are university (selam campus), Sefen selam kebel, Yetebaberut condominium area and Facebook page of kana. During this study the data collected by questionnaires, interview, focusing group and collected comments from Facebook page by video screen shots. The total number of data set (comments) by using yamana sample size formula $(n=N/1+N\ (e)^2)$ is 1022 (one thousands and twenty two).

## Experiment:

Target are 1: *Selam campus*

The comments collected from *Selam campus* by using focusing group and an open ended questionnaire method from total population (1350) the sample size which denotes the total population is figure out in number 308 by using the sample formula of yamana.

$$\text{n (sample size)} = \underline{\text{N (total population)}} \qquad \text{n} = \underline{\quad 1350 \quad} = 308$$
$$1+\text{N (e)}^2 \qquad\qquad 1+1350(0.05)^2$$

From 308 manual selected comments from *Selam campus* 216 positive, 74 negative and the rest 18 comments are neutral. From this manual comments then the proposed system identified from 216 positive comments 203 true positive (tp), 8 false negative (fn), 5 false neutral (fneu) and from 74 negative comments 61 true negative,11 false positive and 2 false neutral and from 18 neutral comments 14 true neutral, 3 false negative and 1 false positive. The false positive will be 12, false negative will be 11 and false neutral will be comments. fp=11+1=12 and fn=8+3=11 and fneu=2+5=7.

Polarity in percent    Tp=93.98%        Tn=82.43%        Tneu=77.8%

       Fp=20.46%        Fn = 20.44%        Fneu = 5%

From the above data the accuracy, precision, recall and f-score are:

$$Accuracy = \frac{tp+tn+tneu}{tp+fp+tn+fn+fneu+tneu}$$

$$= \frac{203+61+14}{203+12+61+11+7+14} = \frac{278}{308} = \mathbf{0.90}$$

$$Precision = \frac{tp}{tp+fp}$$

$$= \frac{203}{203+12} = \frac{203}{215} = \mathbf{0.94}$$

$$Recall = \frac{tp}{tp+fn}$$

$$= \frac{203}{203+11} = \frac{203}{214} = \mathbf{0.95}$$

$$F\text{-}score = 2*\frac{\overline{\phantom{i}}cision*Recall}{Precision+Recall}$$

$$= \frac{0.94*0.95}{0.94+0.95} = \frac{0.8}{1.89} = \mathbf{0.47*2= 0.94}$$

Target are 2: *Sefen selam kebele*

The comments collected from *Sefen selam kebele* by using distributing an open ended questionnaire from total population (240) the sample size which denotes the total population is figure out in number 150. From 150 manual selected comments 140 positive, 10 negative and there is no neutral comments given by the viewers. From this manual comments, the proposed system identified from 140 positive comments 128 true positive (tp), 12 false negative (fn), from 10 negative comments 7 true negative, 3 false positive.

| Polarity in percent | Tp=91.43% | Tn=70 % | Tneu=0% |
|---|---|---|---|
| | Fp=30% | Fn = 8.57% | Fneu = 0% |

From the above data the accuracy, precision, recall and f-score is: Accuracy, Precision, Recall and F-score respectively are **0.90, 0.98, 0.91 and 0.94**.

Target are 3: *Yetebaberut condominium area*

The comments collected from *Yetebaberut condominium area* by using distributing an open ended questionnaire from the total population (280) the sample size which denotes the total population comments is figure out in number 164. From 164 manual selected comments 130 positive, 25 negative and the rest 9 comments are neutral when manually collected. From this manual comments then the proposed system identified from 130 positive comments 94 true positive (tp), 29 false negative (fn), 7 false neutral (fneu) and from 25 negative comments 17 true negative,3 false positive and 5 false neutral and from 9 neutral comments 4 true neutral, 2 false negative and 3 false positive. The false positive will be 6; false negative 31 and false neutral will be 12 comments. fp=3+3=6 and fn=29+2=31 and fneu=7+5=12

| Polarity in percent | Tp=72.30% | Tn=68 % | Tneu=44.44% |
|---|---|---|---|
| | Fp=45.33 % | Fn = 44.52% | Fneu = 25.38% |

From the above data the accuracy, precision, recall and f-score respectively are **0.70, 0.94, 0.75 and 0.83.**

Target are 4: *Facebook page*

The comments collected from *Facebook* page from the total viewer (9,500) which are viewer of the KANA TV program the sample size is in number 400. From 400 manual selected comments 259 positive, 90 negative and the rest 21 comments are neutral when manually collected. From this manual comments then the proposed system identified from 259 positive comments 245 true positive (Tp), 10 false negative (Fn), 4 false neutral (Fneu) and from 90 negative comments 75 true negative (Tn), 14 false positive (Fp) and 1 false neutral (Fneu) and from 21 neutral comments 17 true neutral (Tneu), 3 false negative (Fn) and 1 false positive(Fp).

Polarity in percent   Tp=94.6%        Tn=83.33%         Tneu=80.95%

Fp=20.32%       Fn = 18.15%          Fneu = 2.65%

From the above data the accuracy, precision, recall and f-score are respectively **0.91,** 0.94, 0.95 and 0.945

Table 7: Experimental result of the system

| No | Target area | Sample size | Accuracy | Precession | recall | f-score |
|----|-------------|-------------|----------|------------|--------|---------|
| 1 | *Selam campus* | 308 | 0.90 | 0.94 | 0.95 | 0.94 |
| 2 | *Sefen selam* | 150 | 0.90 | 0.98 | 0.91 | 0.94 |
| 3 | *Yetebaberut* | 164 | 0.70 | 0.94 | 0.75 | 0.83 |
| 4 | *Facebook page* | 400 | **0.91** | 0.94 | 0.95 | 0.945 |
| **Average** | | | **0.85** | **0.95** | **0.89** | **0.91** |

The researches come up with four target area which are selam campus (ሠላም ካምፓስ), sefen selam kebele (ሰፈነ ሠላም ቀበሌ), yetebaberut condominium area (የተባበሩት ኮንዶምኒየም) and website/Facebook page. For all target area 1022 (one thousand and twenty two) comments are collected about the program and from each target area there would be accuracy, precession, recall and f-score. Generally the accuracy of this study is 85%

## CHAPTER FIVE

## 5.  Conclusion and future work

### 5.1.  Conclusion

Sentiment analysis in Amharic using NLP tasks like Amharic text tokenization, Amharic stemming is good as to generalize the sentiments of some product. A dictionary is no more than a list of words that share a category. For example, you can have a dictionary for positive expressions.  Therefore, developing corpus for the Amharic sentiment analysis will be done more as the paper giving direction. In this research applied sentence and document-level (Hybrid) Sentiment analysis in the viewers/non-viewers comment review domain. For Amharic language, there is no standardized corpus for the purpose of Sentiment analysis. People usually use positive words in negative reviews, but the word is preceded by "not" (or some other negative word), such as "not great". Sometimes people give objective text to express their opinion but the classifier did not identify those facts from opinions.

From this paper the experiment is conducted in four target area which are selam campus (ሠላም ካምፓስ), sefen selam kebele (ሰፈነ ሠላም ቀበሌ), and yetebaberut condominium area (የተባበሩት ኮንዶምኒየም). Mainly With the methodology of website/Facebook page, focusing group discussion and an open ended questionnaires. For all target area 1022 (one thousand and twenty two) comments (sample size) are collected about the program and from each target area there would be accuracy, precession, recall and f-score.

From selam campus target area the sample size for the population of 1350 (one thousands and three hundred fifty) is 308 (three hundred eight) using this figure the Accuracy, precession, recall and f-score respectively are 0.90, 0.94, 0.95and 0.94. In sefen selam kebele with sample size 150 the Accuracy, precession, recall and f-score respectively are 0.90, 0.98, 0.91 and 0.94. In *Yetebaberut condominium area* with sample size 164 the Accuracy, precession, recall and f-score respectively are 0.70, 0.94, 0.75 and 0.83. and with the target area *Facebook page* sample size 400 the Accuracy, precession, recall and f-score respectively are 0.91, 0.94, 0.95 and 0.945.

Generally the accuracy of this study is 85%. In order to answer the research questions, all the articulated objectives were reached in all chapters.

## 5.2. Future work

To accomplish good results in Sentiment analysis and enrich upcoming researches in this area, the following recommendations/Future works are suggested to the researcher to conduct the research.

➢ In the future should advance the performance of the sentiment and improve the available Amharic sentiment lexicon dictionaries. This may be done by increasing the size of the lexicon, by considering phrase-level Amharic sentiment terms and by improving the quality of the lexicon considering more different domains.

➢ Developing corpus for the Amharic sentiment analysis should be done more in the future.

➢ Handling sentiment documents with idiomatic expressions from different domains can also be another focus of future research works.

➢ In this paper is not considered the total concept of the sentence and document so it is open for the researchers to come up with the concept.

➢ In our study the domain is focused on KANA TV show for the next it should be consider the other popular TV show (EBS, Nahoo, Walta, LTV,…etc.) even Ethiopian broadcasting.

# REFERENCES

[1]  A. H. ,. K. Walaa Medhat, "Sentiment analysis algorithms and applications: A survey," *Ain Shams Engineering,* pp. 1093-1113, 2014.

[2]  U. S. TEJASWINEE WAKDE, "Rule Based Approach for Currency Interpretation in Indian Languages," in *Recent Advances in Computer Science*, INDIA, 2014.

[3]  B. Liu, Sentiment Analysis and Opinion Mining, Morgan & Claypool Publishers, 2012.

[4]  "WIKIPEDIA (The Free Encyclopedia)," April 2016. [Online]. Available: https://en.wikipedia.org/wiki.

[5]  "Wikipedia,the free encyclopedia," [Online]. Available: https://en.wikipedia.org/wiki.

[6]  B. P. a. L. Lee, Opinion Mining and Sentiment Analysis, USA, 2008.

[7]  B. P. a. L. Lee, Opinion Mining and Sentiment Analysis, 2008.

[8]  "an introduction to sentiment analysis opinion mining in meaningcloud," 13 october 2015. [Online]. Available: https://www.meaningcloud.com/blog1.

[9]  M. Bonzanini, "Mining Twitter Data with Python (Part 6 – Sentiment Analysis Basics)," 17 may 2015. [Online]. Available: https://marcobonzanini.com.

[10] J.-H. C. Anna Stavrianou, Opinion Mining Issues and Agreement, Atelier FODOP, 2008.

[11] F. X. &. X. CHENG, "Opinion Mining," 3 Dec 2007. [Online].

[12] "Opinion Miner - Online sentiment analysis," [Online]. Available: http://www.slideshare.net/igmelig/opinion. [Accessed April 2010].

[13] ". R. Huey Y., ""Chinese sentiment analysis using maximum entropy"," in *proceedings of the workshop on sentiment analysis where AI meets psychology(SAAIP)*, Thailand, 2011, pp. 81-84.

[14] Liu, B., Sentiment Analysis and opinion mining, Morgan and Claypool, 2012.

[15] "Opinion Miner - Online sentiment analysis," [Online]. Available: http://www.slideshare.net/igmelig/opinionminer. [Accessed April 2010].

[16] A. Katrekar, "Big Data Analytics," [Online]. Available: www.globallogic.com.

[17] S. Gebremeskel, Sentiment Mining Model for Opinionated Amharic Texts, Addis Ababa: Unpublished, 2010.

[18] L. L. a. S. V. Bo Pang, "Thumbs up? Sentiment Classification using Machine Learning Techniques," in *Proceedings of the Conference on Empirical Methods in Natural Language Processing (EMNLP)*, 2002, pp. 79-86.

[19] P. D. Turney, "Thumbs Up or Thumbs Down? Semantic Orientation Applied to Unsupervised Classification of Reviews," in *NewBrunswick,N.J.*, 2002.

[20] A. B. S. B. S. K. Sandip Mali, "Sentiment of Sentence in Tweets: A Review," *IOSR Journal of Computer Engineering,* vol. 17, no. 6, pp. 157-162, 2015.

[21] K. M. Razieh Asgarnezhada, "A Comparative Classification of Approaches and Applications in Opinion Mining," *International Academic Journal of Science and Engineering,* vol. 2, no. 5, p. 2015, 2015.

[22] S. &. C. M. Das, "Yahoo! For Amazon: Extracting Market Sentiment from Stock Message Boards," in *Paper presented in the 8th Asia Pacific Finance Association Annual Conference*, Bankok, Thailand, 2001.

[23] J. W. a. T. W. Ellen Riloff, "Learning Subjective Nouns Using Extraction Pattern Bootstrapping," In CoNLL03, 2003, pp. 25-32.

[24] V. H. a. J.Wiebe, Effects of Adjective Orientation and Gradability on Sentence Subjectivity, COLING, 2000.

[25] S. K. a. E. Hovy, Determing the Sentiment of Opinions, COLING, 2004.

[26] A. A. S. C.-O. a. E. K. R. M. Gamon, Pulse: Mining Customer opinions from Free Text, IDA, 2005.

[27] J. R. E. Wiebe, "Creating Subjective and Objective Sentence Classifers from Unannotated Text," in *Proceedings of International Conference on Intelligent Text Processing and Computational Linguistics*, 2005.

[28] L. Z. a. C. Li, Ontology Based Opinion Mining for Movie reviews, Springer-Verlag Berlin Heidelberg, 2009.

[29] X. X. e. al, Categorizing Terms' Subjectivity and Polarity Manually for Opinion Mining in Chinese, 2009.

[30] F. F. G. G. Alessia D'Andrea, "Approaches, Tools and Applications for Sentiment Analysis

Implementation," *International Journal of Computer Applications,* vol. 125 , no. 3, p. 0975 – 8887, 2015.

[31] A. O'Neill, Sentiment Mining for Natural Language Documents, Australian, 2009.

[32] N. J. a. B. Liu, "Opinion spam and analysis," in *Proceedings of the Conference on Web Search and Web Data Mining (WSDM)*, Palo Alto, California, USA, 2009.

[33] B. Liu, Sentiment Analysis and Subjectivity, Handbook of Natural Language Processing, Second Edition, Chemical Rubber Company (CRC) Press, 2010.

[34] S. R. D. a. M. Y. Chen, Yahoo! for Amazon: Sentiment extraction from small talk on the Web, 2007.

[35] T. N. R. B. a. W. N. J. Yi, "Sentiment analyzer: Extracting sentiments about a given topic using natural language processing techniques," in *Proceedings of the IEEE International Conference on Data Mining (ICDM)*, 2003.

[36] L. B. Hu Minging, "Mining and summarizing customer reviews," in *In: Proceedings of ACM SIGKDD international conference on Knowledge Discovery and Data Mining (KDD'04)*, 2004.

[37] B. R. F. C. G. D. M. K. Miller G, WordNet: an on-line lexical database, Oxford Univ. Press, 1990.

[38] D. C. D. B. Mohammad S, "Generating high-coveragesemantic orientation lexicons from overly marked words anda thesaurus," in *In: Proceedings of the conference on Empirical Methods in Natural Language Processing (EMNLP'09)*, 2009.

[39] H. X. Z. F. S. Y. B. J. C. Qiu Guang, DASA: dissatisfaction-oriented advertising based onsentiment analysis, 2010.

[40] M. K. Hatzivassiloglou V, "Predicting the semanticorientation of adjectives," in *In: Proceedings of annual meeting ofthe Association for Computational Linguistics (ACL'97)*, 1997.

[41] K. M. Fahrni A, "Old wine or warm beer: target-specific sentiment analysis of adjectives," in *In: Proceedings of the symposium on affective language in human and machine, AISB*, 2008.

[42] T. P, "Thumbs up or thumbs down?: semantic orientationapplied to unsupervised classification of reviews," in *In: Proceedingsof annual meeting of the Association for ComputationalLinguistics (ACL'02)*, 2002.

[43] H. E. Kim S, "Determining the sentiment of opinions," in *In: Proceedings of interntional*

*conference on Computational Linguistics (COLING'04); 2004*, 2004.

[44] A. H. K. Walaa Medhat, "Sentiment analysis algorithms and applications:A survey," *Ain Shams Engineering Journal,* vol. 5, p. 1093–1113, 2014.

[45] S. Balijepalli, A Modular Domain Independent Sentiment Analysis System, Maryland, 2007.

[46] A. K. a. D. Inkpen, "Sentiment Classification of Movie and Product Reviews Using Contextual Valence Shifters," *Computational Intelligence,* vol. 22, 2006.

[47] "Unbiased reviews by real people," a shopping company. [Online]. [Accessed october 2009].

[48] F. J. X.-Y. Z. a. L. Z. L. Zhuang, "Movie review mining and summarization," in *in Proceedings of the ACM SIGIR Conference on Information and Knowledge Management (CIKM)*, 2006.

[49] A. S. F. Esuli, "Sentiwordnet: A publicly available lexical resource for opinion mining," in *In: Proceedings of 5th Conference on Language Resources and Evaluation, LREC*, 2006.

[50] A. S. F. Esuli, "Sentiwordnet: A publicly available lexical resource for opinion," in *In Proceedings of 5th Conference on Language Resources and Evaluation, LREC* , 2006.

[51] "The Internet Movie Database," 17 october 1990. [Online]. Available: http://www.imdb.com.

[52] X. X. e. al, Categorizing Terms' Subjectivity and Polarity Manually for Opinion Mining in Chinese, IEEE, 2009.

[53] D. Q. Dong Zhendong, HowNet and the Computation of Meaning, World Scientific Publishing Co.Pte.Ltd, 2006.

[54] P. C. a. L. D. Sigrid Maurel, A Hybrid Method for Sentiment Analysis, France: Statistical Analysis Software (SAS) press, 2008.

[55] S. Gebremeskel, "sentiment analysis for opinionated Amharic text," in *Addis ababa university*, Addis Ababa, 2010.

[56] L. Pan, Sentiment Analysis in Chinese, Brandeis University, 2012.

[57] MarkosKassa, "Implementing an open source Amharic resource grammar in GF," in *master's thesis*, University of Gothenburg, 2010.

[58] Basic Sentiment Analysis with Python, 2012.

[59] P. d. N. t. B. L. Varela, "Sentiment Analysis," 2012.

[60] B. Liu, Sentiment Analysis and opinion mining, Morgan and Claypool , 2012.

[61] A. S. F. Esuli, "Sentiwordnet: A publicly available lexical resource for opinion mining," in *In Proceedings of 5th Conference on Language Resources and Evaluation LREC*, 2006.

[62] "The Free Encyclopidea," in *wikipedia*.

[63] S. al., "Unit selection for Amharic using FESTVOX,5th ISCA speech synthesis workshop," in *Language technology research center*, 2004.

[64] Y. Afework, "Automatic Amharic text categorization," Addis Ababa university, Addis Ababa, 2007.

[65] "The Free Encyclopidea," Wikipedia, 2016. [Online]. [Accessed 2016].

[66] S. N. a. R. J. Richa Sharma, "OPINION MINING OF MOVIE REVIEWS AT DOCUMENT LEVEL," *International Journal on Information Theory (IJIT),* vol. 3, 2014.

[67] M. Belete, "Sentiment Analysis for Amharic opinionated text," Addis Ababa university, Addis Ababa , 2013.

[68] M. woldekirkos, አማርኛ ሰዋሰው, addis abeba: Berhanenaselam printing press, 1934.

[69] woldekirkos, Mersehaizen, አማርኛ ሰዋሰው, addis abeba: Berhanenaselam printing press, 1934.

[70] Y. Afework, "Automatic Amharic text categorization," Addis Ababa university, Addis Ababa, 2007.

[71] M. woldekirkos, አማርኛ ሰዋሰው, addis abeba: Berhanenaselam printing press, 1934.

[72] BayeYemam, የአማርኛ፣ሰዋስው, addis abeba: ት.መ.ሟ.ማ.ድ, 1987.

Appendix 1: Questionnaire for KANA TV Viewer

<div align="center">

**ባሕር ዳር ዩኒቨርስቲ ኮምፒውቲንግ ፋኩልቲ**

**የኮምፒውተር ሳይንስ ትምህርት ክፍል**

*በቃና ቴሌቪዥን ተመልካቾች የሆኑና ያልሆኑ የሚሞላ መጠይቅ*

</div>

<div align="center">

**የመጠይቁ ዓላማ፤**

</div>

በባሕር ዳር ዩኒቨርስቲ በኮምፒውተር ሳይንስ ትምህርት ለማስተርስ ዲግሪ ምርምር ማሟያ የሚሆን በአማርኛ ቋንቋ ላይ መሠረት ያደረገ የዘመኑ ቴክኖሎጂ በሚፈቅደው መልኩ ኮምፒውተራይዝድ የሆነ የቃና ቴሌቪዥን ተመልካቾች አስተያየት መለያ ሲስተም ለመስራት ግብዓት የሚሆን ነው፡፡

<div align="center">

ስለሚሰጡኝ የመጠይቅ ምላሽ ከልብ አመሰግናለሁ

</div>

**ጥያቄ፡-** የቃና ቴሌቪዥን ተመልካች/ተከታታይ ነዎት ?

<div align="center">

አዎ ☐       አይደለሁም ☐

</div>

*መልስዎ አዎ ከሆነ ስለ አቀራረቡ ያልዎትን ማንኛውም አስተያየት ይግለፁልኝ*

_____
_____
_____
_____
_____::

*መልስዎ አይደለም ከሆነ የማይከታተሉበትን ምክንያት ይጻፉ*

_____
_____
_____
_____
_____::

<div align="center">

**የመጠይቁ አዘጋጅ**

*ዮሐንስ መኮንን*

በባሕር ዳር ዩኒቨርስቲ የኮምፒውተር ሳይንስ ትምህርት ክፍል

የማስተርስ ዲግሪ ተማሪ

</div>

Appendix 2: Sample Response of Viewers

ባሕር ዳር ዩኒቨርስቲ ኮምፒውቲንግ ፋኩልቲ

የኮምፒውተር ሳይንስ ትምህርት ክፍል

በቃና ቴሌቪዥን ተመልካቾች የሆኑና ያልሆኑ የሚሞላ መጠይቅ

## የመጠይቁ ዓላማ፤

በባሕር ዳር ዩኒቨርስቲ በኮምፒውተር ሳይንስ ትምህርት ለማስተርስ ዲግሪ ምርምር ማሟያ የሚሆን በአማርኛ ቋንቋ ላይ መሠረት ያደረገ የዘመኑ ቴክኖሎጂ በሚፈቅደው መልኩ ኮምፒውተራይዝድ የሆነ የቃና ቴሌቪዥን ተመልካቾች አስተያየት መለያ ሲስተም ለመስራት ግብዓት የሚሆን ነው፡፡

ስለሚሰጡኝ የመጠይቅ ምላሽ ክልብ አመሰግናለሁ

ጥያቄ፡- የቃና ቴሌቪዥን ተመልካች/ተከታታይ ነዎት ?

አዎ [ X ]     አይደለሁም [ ]

መልስዎ አዎ ከሆነ ስለ አቀራረቡ ያልዎትን ማንኛውም አስተያየት ይግለፁልኝ

_የቃና ቴሌቪዥን የልዩ ዝግጅት ............ ከፍ ሲል ....... ............_
_................... .................. ................._
_........ ......................... ................._
_.......... ......................... ................_
_............... ............................ ............._
_.......... ................. ............_
_____ ::_


መልስዎ አይደለም ከሆነ የማይከታተሉበትን ምክንያት ይጻፉ

_____
_____
_____
_____
_____
_____
_____ ::

59

Appendix 3: Sample Response of Non-Viewers

ባሕር ዳር ዩኒቨርስቲ ኮምፒውቲንግ ፋኩልቲ

የኮምፒውተር ሳይንስ ትምህርት ክፍል

በቃና ቴሌቪዥን ተመልካቾች የሆኑና ያልሆኑ የሚሞላ መጠይቅ

### የመጠይቁ ዓላማ፤

በባሕር ዳር ዩኒቨርስቲ በኮምፒውተር ሳይንስ ትምህርት ለማስተርስ ዲግሪ ምርምር ማማያ የሚሆን በአማርኛ ቋንቋ ላይ መሠረት ያደረገ የዘመኑ ቴክኖሎጂ በሚፈቅድለው መልኩ ኮምፒውተራይዝድ የሆነ የቃና ቴሌቪዥን ተመልካቾች አስተያየት መለያ ሲስተም ለመስራት ግብዓት የሚሆን ነው።

ስለሚሰጡኝ የመጠይቅ ምላሽ ከልብ አመሰግናለሁ

ጥያቄ፦ የቃና ቴሌቪዥን ተመልካች/ተከታታይ ነዎት ?

አዎ ☐       አይደለሁም ☒ *x*

መልስዎ አዎ ከሆነ ስለ አቀራረቡ ያልዎትን ማንኛውም አስተያየት ይግለፁልኝ

_____
_____
_____
_____
_____
_____
_____
_____ ።

መልስዎ አይደለም ከሆነ የማይከታተሉበትን ምክንያት ይፃፉ

የቃና ቴሌቪዥን ተመላካች አይደለሁም። ነገር ግን �ba ዝግጅቱ አቀራረብ እንዳለሁ ተመልካቻዬም ከዚ2ዜ8700ን ሥራ ፉታ በመለሁት እያለሁ። የፊልም አቀራረቡ ገአበ ወ,ጋፉ33 የመለፉ አይደለም። ገ211ጅ33 የማ2ቀ23 ነዉ። ተ0ሌ ፤ የወእገ0ኢፉና የጠፕ ሥራተኛን የለ0ሩ 2ዜ የሚ2ዛለ ገቶም ነዉ የ0ሐ0ፉ። ለእዝ0ቀበፉ የቃ0 ቴክ7ዥን ዝ0ፅ7ና አቀbLበ አዮ0ፉ7ፕ0።

_____
_____ ።

60

Appendix 4: Sample of KANA TV viewer comments

1. በአቀራረቡ የሚገርም ነው

2. እጅግ በጣም አሪፍ ዝግጅት ነው የሚቀርበው አዘጋጆች በርቱ

3. በነገራችን ላይ በቃና የቴሌቪዥን የሚቀርቡ ፊልሞች አይመቹኝም

4. ጥቁር ፍቅር ደስ የሚል ፊልም ነው

5. ፊልሙ በጣም አሪፍ ነው ግን ከሀገራችን ባሕል ጋር የሚቃረን ፍልም የሚቀርብበት ነው፡፡

6. አሪፍ ነው ግን አይኮርጅ

7. በጣም የሚደነቅ የአገራችንን የፊልም ደረጃ የሚያሳድግ ፊልም የሚቀርብበት ነው፡፡ በርቱ ተበራቱ እንላለን

8. ኩዚ ጉኒ ተመችቶናል፤ በማለቁ ደግሞ አዝኛለሁ

9. በጣም የሚገርም ፊልም የሚቀርብበት የቴሌቪዥን መስኮት ነው

10. የተማሪን ጊዜ የሚሻማ በመሆኑ የቃና ተቃዋሚ ነኝ፡፡

11. # ታግ ተመችቶናል፤

12. ኢትዮጵያዊያን የቤት እመቤቶችን ከባሎቻቸው ጋር ያጣላ ወይም ያፋታ የዝግጅት አቀራረብ ነው ያለው፤

13. ትዳርን የሚያፋታ ነው፡፡

14. የኢትዮጵያዊያንን አዕምሮ ለማደንዘዝ የታቀደ ቴሌቪዥን ነው፡፡

15. በቃና እንቅልፍ አጣን

16. ፋጡማ ጉል መሳጭ ልጅ ናት

17. የተቀማ ሕይወት አንደኛ ምርጥ ፊልም ነው

18. ከኢትዮጵያ ባሕል ጋር ጥሩ ያልሆነ ወይም የተቃረነየዝግጅት አቀራረብ በመኖሩ አይመቸኝም፡፡

19. በጣም አሪፍ ነው አቦ ይቀጥል ብለናል፤

Appendix 5: Rules for Stemming

# RULE 1 - Take input as it is

def stem(input1):

    print(input1)

    collection = [input1]

# RULE 2 - Take out the right most suffix - From input 1

    input2 = re.match("(.+)(iwu|wu|wi|awī|na|mi|ma|li|ne|ache)",input1)

    if input2:

        print(input2.group(1)+'-'+input2.group(2))

        input2 = input2.group(1);

        collection.append(input2)

    else:

        input2 = input1

# RULE 3 - Take out the inner most suffix

    input3 = re.match('(.+)(ochi|bache|wache)',input2)

    input3 = re.match('(.+)(chi|ku|ki|ache|wal)',input2) if not input3 else input3

    if input3:

        print(input3.group(1)+'-'+input3.group(2))

        input3 = input3.group(1)

        collection.append(input3)

    else:

        input3 = input2

# RULE 4 - Take out the most left prefix - From input 1

    input4 = re.match('(yete|inide|inidī|āli)(.+)',input1)

    input4 = re.match('(ye|yi|masi|le|ke|inid|be|sile)(.+)',input1) if not input4 else input4

```python
    if input4:

            print(input4.group(1)+'-'+input4.group(2))

            input4 = input4.group(2)

            collection.append(input4)

    else:

            input4 = input1
```

# RULE 5 - Take out the right most suffix - From input 4

```python
    input5 = re.match('(.+)(iwu|wu|w|awī|na|mi|ma|li|ne|che)',input4)

    if input5:

            print(input5.group(1)+'-'+input5.group(2))

            input5 = input5.group(1)

            collection.append(input5)

    else:

            input5 = input4
```

# RULE 6 - Take out the inner most suffix - From input 4

```python
    input6 = re.match('(.+)(ochi|bache|wache)',input5)

    input6 = re.match('(.+)(chi|ku|ki|che|wal)',input5) if not input6 else input6

    if input6:

            print(input6.group(1)+'-'+input6.group(2))

            input6 = input6.group(1)

            collection.append(input6)

    else:

            input6 = input5
```

# RULE 7 - Take out the inner most prefix - From input 1

```python
    input7 = re.match('(te|mī|mi|me|mayit|ma|bale|yit|āya|ā|āyi)(.+)',input4)
```

```python
        if input7:

                print(input7.group(1)+'-'+input7.group(2))

                input7 = input7.group(2)

                collection.append(input7)

        else:

                input7 = input4
```

# RULE 8 - Take out the right most suffix - From input 7

```python
        input8 = re.match('(.+)(iwu|wu|w|awī|na|mi|ma|li|ne)',input7)

        if input8:

                print(input8.group(1)+'-'+input8.group(2));

                input8 = input8.group(1)

                collection.append(input8)

        else:

                input8 = input4
```

# RULE 9 - Take out the innermost suffix - From input 8

```python
        input9 = re.match('(.+)([^iīaeou])'?',input8)

        if input9:

                print(input9.group(1)+'-'+input9.group(2));

                input9 = input9.group(1)

                collection.append(input9)

        else:

                input9 = input4

        print(collection)

        return collection
```
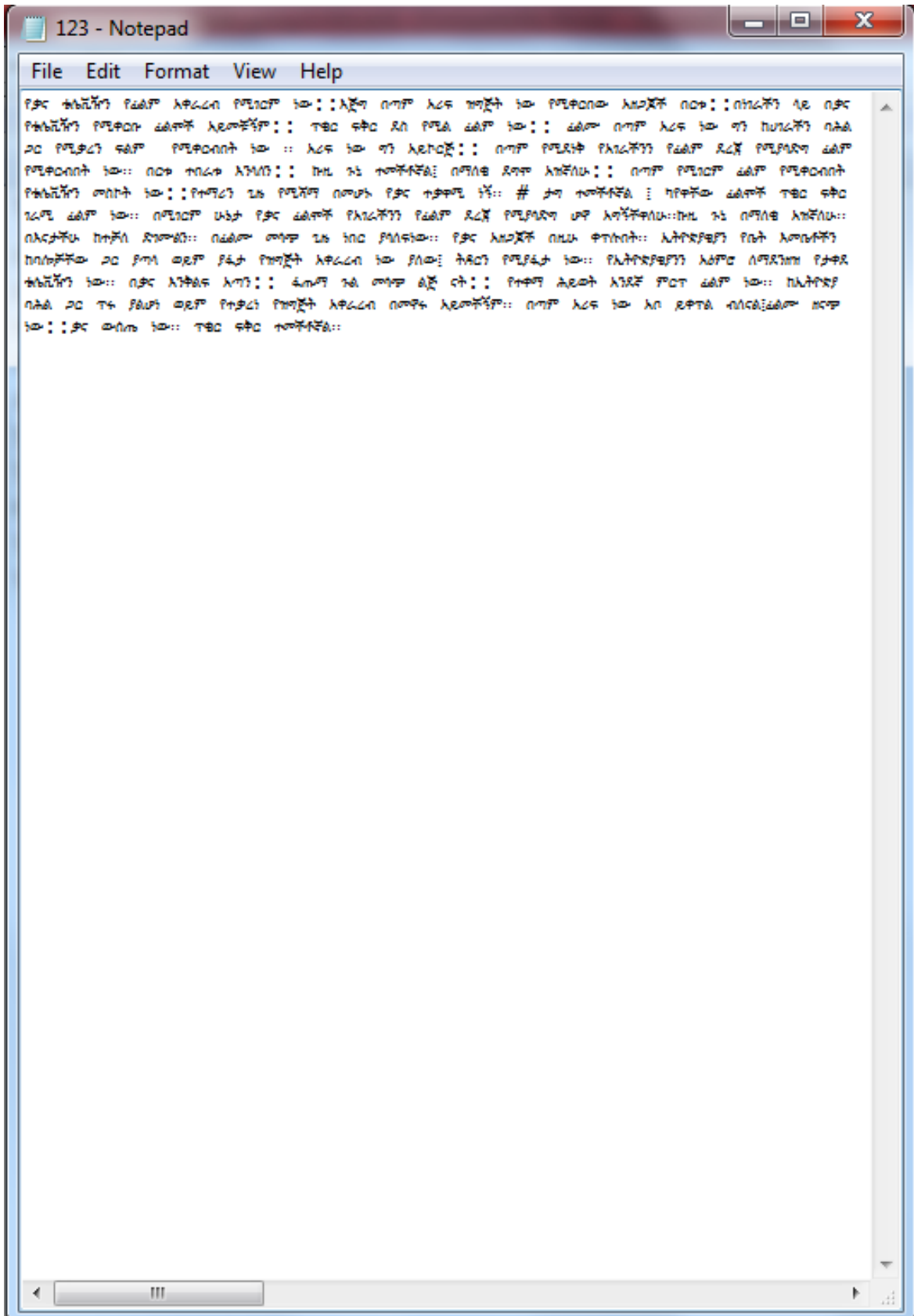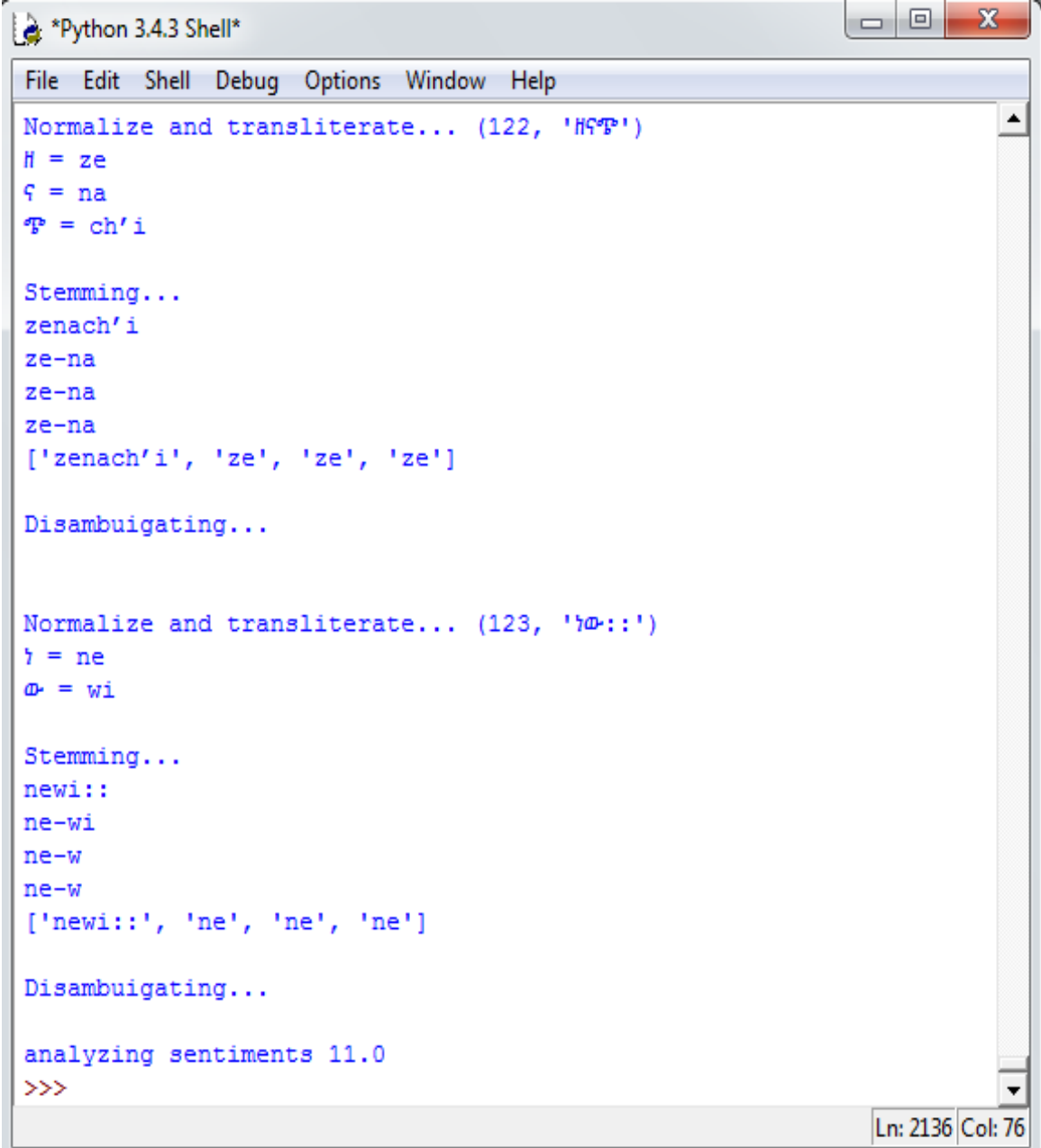
Appendix 6: Sample input text

123 - Notepad

File    Edit    Format    View    Help

የቃር ቴክኒሻን የልፊም አፈራረስ የሚገርም ነው፤፤እጅግ በጣም እርፍ ዝግጁት ነው የሚቀበለው አገልጋች ቦዑ፤፤በገሪቺን ላይ በቃር የቴክኒሻን የሚቀብ ሊልማች አደመጭኛም፤፤ ፐቦ ፍቅር ደስ የሚል ሊልም ነው፤፤ ሊልም በጣም እርፍ ነው ግን ከነገሪቺን በአል ጋር የሚቃረን ፍልም   የሚቀስሰት ነው ፡፡ እርፍ ነው ግን እደኮጅ፤፤ በጣም የሚደፍ የእገሪቺን የልፊም ደረጃ የሚያስን ሊልም የሚቀስሰት ነው፡፡ ቦዑ ተበረቱ እንህባ፤፤ ከዚ ነኪ ተመጭፍኛ፤ በግበቆ ደጣ እነኜሰሁ፤፤ በጣም የሚገርም ሊልም የሚቀስሰት የቴክኒሻን መስከት ነው፤፤የተግረን ጊዜ የሚሻግ ስሙሆኪ የቃር ተቃመጊ ነኝ፡፡ # ታግ ተመጭፍ ፤ ከቀቀጭው ሊልማች ፐቦ ፍቅር ገረሚ ሊልም ነው፡፡ በሚገርም ሁኔታ የቃር ሊልማች የእገሪቺን የልፊም ደረጃ የሚያስን ሆፍ እግኛጭቀሰሁ፡፡ከዚ ነኪ በግበቆ እነኜሰሁ፡፡ በእረታችሁ ከቀቦስ ደነማቆስ፡፡ በሊልም መጭፍ ጊዜ ከበር የእሰኖነው፡፡ የቃር እለጋጭ በዚሁ ቀተስስት፡፡ እትቀጵቀፍን የቤት እመቡተቺን ከሰቦቸጭው ጋር ያግባ ወደም የፊቃ የመግጁት እፈረረስ ነው ያሰበ፤ ትፈርን የሚፈታ ነው፡፡ የእትቀጵቀፍን እፈማቀ ለግደንጣ የታቀረ ቴክኒሻን ነው፡፡ በቃር እንቀስስ እማን፤፤ ፈጠማ ነል መፍቀ ብጅ ርት፤፤ የተቀማ ሐደመት እንደጆ ፍርኯት ሊልም ነው፡፡ ከእትቀጵ በአል ጋር ፐፉ ያሰበነ ወደም የተፈረን የመግጁት እፈረረስ በመፉ እደመቆኛም፡፡ በጣም እርፍ ነው እስ ደቀፕል ብስረበ|ሊልም ዘሞፍ ነው፤ ፤ ፐር ወሰሰ ነው፡፡ ፐቦ ፍቅር ተመጭፍኛ፡፡

Appendix 7: Sample Output

```
*Python 3.4.3 Shell*                                   [_] [□] [X]

File  Edit  Shell  Debug  Options  Window  Help

Normalize and transliterate... (122, 'ዘናጪ')
ዘ = ze
ና = na
ጪ = ch'i

Stemming...
zenach'i
ze-na
ze-na
ze-na
['zenach'i', 'ze', 'ze', 'ze']

Disambuigating...


Normalize and transliterate... (123, 'ነው::')
ነ = ne
ው = wi

Stemming...
newi::
ne-wi
ne-w
ne-w
['newi::', 'ne', 'ne', 'ne']

Disambuigating...

analyzing sentiments 11.0
>>>
                                              Ln: 2136 Col: 76
```

66