

A DISSERTATION REPORT ON

**SENTIMENT ANALYSIS USING ORIGINAL AND
REVERSED REVIEWS**

SUBMITTED TO THE SAVITRIBAI PHULE PUNE UNIVERSITY, PUNE
IN PARTIAL FULFILLMENT OF THE REQUIREMENTS
FOR THE AWARD OF THE DEGREE OF

MASTER OF ENGINEERING (Computer Engineering)

BY

Kaushik S. Hande

Exam No: 6989

Under The Guidance of

Prof. A. G. Phakatkar



**DEPARTMENT OF COMPUTER ENGINEERING
PUNE INSTITUTE OF COMPUTER TECHNOLOGY
Sr. No. 27, Pune Satara Road, Dhankawadi
Pune 411043**

PUNE INSTITUTE OF COMPUTER TECHNOLOGY



DEPARTMENT OF COMPUTER ENGINEERING

CERTIFICATE

This is to certify that the dissertation entitled

“SENTIMENT ANALYSIS USING ORIGINAL AND REVERSED REVIEWS”

Submitted by

Kaushik S. Hande

Exam No: 6989

is a bonafide work carried out by him under the supervision of Prof. A. G. Phakatkar and it is submitted towards the partial fulfillment of the requirement of Savitribai Phule Pune University, for the award of the degree of Master of Engineering (Computer Engineering)

Prof. A. G. Phakatkar	Dr. Rajesh Ingle	Dr. P. T. Kulkarni
Internal Guide	H.O.D	Principal
Dept. of Computer Engg.	Dept. of Computer Engg.	
PICT, Pune-43	PICT, Pune-43	PICT, Pune-43

ACKNOWLEDGEMENT

*It gives me great pleasure in presenting the dissertation report for my dissertation on “**SENTIMENT ANALYSIS USING ORIGINAL AND REVERSED REVIEWS**”.*

*I would like to take this opportunity to thank my internal guide **Dr. A. S. Ghotkar** for giving me all the help and guidance I needed. I am really grateful to them for their kind support throughout the analysis and design phase. Their valuable criticism and suggestions were very helpful.*

*I am grateful to **Dr. Rajesh B. Ingle**, Head of Computer Engineering Department, Pune Institute of Computer Technology for his indispensable support, priceless suggestions and for most valuable time lent as and when required.*

*In the end my special thanks to **Mr. Hrushikesh Kasar** for providing various resources such as well-equipped lab with all needed software platforms, continuous Internet connection, for my dissertation work.*

Kaushik S. Hande

(M.E. Computer Engg.)

List of Figures

6.1	Task Network	22
6.2	Timeline Chart	23
7.1	Use Case : System	29
7.2	Activity Diagram	30
7.3	Level 0 DFD	31
7.4	Level 1 DFD	31
7.5	Level 2 DFD	32
7.6	Sequence Diagram	32
8.1	Proposed System Architecture	36
8.2	Class Diagram	37
8.3	Component Diagram	38
9.1	Sample matrix creation with word counts	42
9.2	Sample matrix creation with tf-idf	42
9.3	Review classification with negative and positive percentages	43
9.4	Three classifiers results(Accuracy, recall, precision, f-measure with confusion matrix)	43
9.5	Three classifiers results(Accuracy, recall, precision, f-measure with confusion matrix)	44
11.1	Classifiers accuracy score comparison	48
11.2	Variations in classifier accuracies as the no. of features increases	48
A.1	IJIRCCE Certificate	57

A.2	cPGCON Certificate	58
A.3	cPGCON Review Sheet	59

List of Tables

4.1	Literature Survey	12
4.2	Literature Survey	12
6.1	Risk Table	18
6.2	Risk Probability Definitions	19
6.3	Risk Impact Definitions	19
6.4	Risk 1	19
6.5	Risk 2	20
6.6	Risk 3	20
9.1	Polarity Training Dataset	41
11.1	Unigram Features Results	47
11.2	Bigram Features Results	47
11.3	Trigram Features Results	47
B.1	Dissertation Task Set	61

ABSTRACT

Bag of words is used for modeling in machine learning algorithms. However, BOW is not able to handle negation well because of its fundamental deficiencies . Many ways are used to handle the problem of negation which results into polarity shift . They require either knowledge about language constructs or extra human interventions which eventually increases the complexity. In this paper, a data expansion technique, called dual sentiment analysis (DSA), is used to address the polarity shift problem due to negation in sentiment classification. Original and reversed training reviews are used for learning in a sentiment classifier and prediction is done on test reviews.

INDEX

1	SYNOPSIS	1
1.1	Dissertation Title	2
1.2	Internal Guide	2
1.3	Problem Statement	2
1.4	Objectives	2
1.5	Hypothesis	2
1.6	Relevant Mathematics Associated with Dissertation	2
1.6.1	Mathematical Model	2
1.6.2	Metrics for Performance Evaluation	3
2	TECHNICAL KEYWORDS	5
2.1	Area of Dissertation	6
2.2	ACM Keywords	6
3	INTRODUCTION	7
3.1	Dissertation Idea	8
3.2	Motivation of Dissertation	9
4	LITERATURE SURVEY	10
4.1	Sentiment Analysis and Polarity Shift	11
4.2	Gap Identification Through Literature Survey	12
5	PROBLEM DEFINITION AND SCOPE	13
5.1	Goals	14
5.2	Objectives	14

5.3	Statement of Scope	14
5.4	Software Context	14
5.4.1	Scikit-learn	14
5.4.2	NumPY	15
5.4.3	Natural language toolkit	15
5.4.4	Matplotlib	15
6	DISSERTATION PLAN	16
6.1	Purpose of the Document	17
6.2	Technical Constraints	17
6.3	Dissertation Estimates	17
6.3.1	Reconciled Estimates	17
6.4	Risk Management	17
6.4.1	Risk Identification	18
6.4.2	Risk Analysis	18
6.4.3	Overview of Risk Mitigation, Monitoring and Management	19
6.5	Staff Organization	20
6.5.1	Team Structure	20
6.5.2	Management Reporting and Communication	20
6.5.3	Dissertation Task Set	21
6.5.4	Task Network	21
6.5.5	Timeline Chart	22
7	SOFTWARE REQUIREMENT SPECIFICATION	24
7.1	Introduction	25
7.2	Purpose and Scope of the Document	25
7.3	Overview of Responsibilities of Developer	25
7.4	Product Overview	25
7.5	Hardware Resources Used	26
7.5.1	Software Requirements	26
7.5.2	Hardware Requirements	26
7.6	Functionality	27

7.7	Input	27
7.8	Output	27
7.9	Major Constraints	27
7.10	Applications	27
7.11	Usage Scenario	28
7.11.1	User Profiles	28
7.11.2	Use Case Views	28
7.12	Model and Description	29
7.12.1	Activity Diagram	29
7.13	Functional Model and Description	31
7.13.1	Data Flow Diagram	31
7.13.2	Sequence Diagram	32
7.14	Non-Functional Requirements	33
7.14.1	Availability	33
7.14.2	Scalability	33
7.14.3	Performance	33
7.14.4	Usability	33
7.14.5	Reliability	33
7.14.6	Maintainability and Changeability	33
8	DETAILED DESIGN DOCUMENT	34
8.1	Introduction	35
8.2	Architectural Design	36
8.3	Class Design	36
8.4	Component Design	38
9	IMPLEMENTATION DETAILS	39
9.1	Introduction	40
9.2	Algorithm	40
9.2.1	Document Classification	40
9.3	Dataset	41
9.4	Snapshots	42

10 TEST SPECIFICATION	45
11 DATA TABLES AND DISCUSSIONS	46
12 CONCLUSION	49
13 FUTURE ENHANCEMENTS	51
14 REFERENCES	53
A PAPERS PUBLISHED	56
A.1 Paper Title	57
A.1.1 IJIRCCE Certification	57
A.2 Paper Title	57
A.2.1 cPGCON Certificate	58
A.2.2 cPGCON Review	59
B DISSERTATION PLANNER	60

CHAPTER 1

SYNOPSIS

DISSERTATION TITLE

SENTIMENT ANALYSIS USING ORIGINAL AND REVERSED REVIEWS

INTERNAL GUIDE

Prof. A. G. Phakatkar

PROBLEM STATEMENT

“ To make use of the original and reversed review samples in pairs for training a statistical classifier and make predictions. ”

OBJECTIVES

- To obtain reversed reviews from each corresponding original reviews.
- To train the classifiers using these reviews.
- To obtain the predictions of labels(positive review or negative review) for test data.

HYPOTHESIS

Polarity shift causes accuracy of classifier to decrease. We assume that original review and corresponding opposite review can be used together to increase the accuracy of review class label prediction and to avoid the problem caused due to polarity shift.

RELEVANT MATHEMATICS ASSOCIATED WITH DISSERTATION

Mathematical Model

$$S = \{s, e, I, O, fmain | \phi\}$$

where,

s = start state

e = end state

I = Inputs to the system

$$I = \{x, x', y, y', D, D'\}$$

where,

x = original sample

x' = reversed sample

$y \in \{0, 1\}$ = The class label of the original sample

$y' = 1 - y$ = The class label of the reversed sample

$D = (x_i, y_i)_{i=1}^n$ = original training set

$D' = (x'_i, y'_i)_{i=1}^n$ = The reversed training set

O = Output

$$O = \{p(x), p(x'), p(x, x')\}$$

where

$p(x)$ = Prediction for the original sample

$p(x')$ = Prediction for the reversed sample

$p(x, x')$ = Dual prediction based on a pair of sample

$$f_{main} = \{f_{reverse}, f_{classifier}\}$$

$f_{reverse}$ = function for reversing the corresponding each review

$f_{classifier}$ = classifier for the prediction of class of review

Metrics for Performance Evaluation

Several statistical measures are used for performance evaluation -

- Accuracy-is the proximity of measurement results to the true value.

$$\frac{TP + TN}{TP + TN + FP + FN} \quad (1.1)$$

- Sensitivity- measures the proportion of positives that are correctly identified

$$\frac{TP}{TP + FN} \quad (1.2)$$

- Specificity- measures the proportion of negatives that are correctly identified

$$\frac{TN}{TN + FP} \quad (1.3)$$

- Positive predictive value- are the proportions of positive results in statistics and diagnostic tests

$$\frac{TP}{TP + FP} \quad (1.4)$$

- Negative predictive value- are the proportions of negative results in statistics and diagnostic tests

$$\frac{TN}{TN + FN} \quad (1.5)$$

CHAPTER 2

TECHNICAL KEYWORDS

AREA OF DISSERTATION

Natural language processing, machine learning, sentiment analysis, opinion mining.

ACM KEYWORDS

A Information Systems

A.1 Information Retrievals

A.1.1 Retrieval tasks and goals

A.1.1.1 Sentiment analysis

A.1.1.2 Clustering and classification

B Computing methodologies

B.1 Machine learning

B.1.1 Supervised learning by classification

B.1.1.1 Multinomial Naive Bayes

B.1.1.2 Random Forest

B.1.1.3 Support Vector Machines

CHAPTER 3

INTRODUCTION

DISSERTATION IDEA

Sentiment is an attitude, thought, or judgement prompted by feeling. Sentiment analysis is also known as opinion mining, it involves studying of peoples sentiments towards certain entities. Internet is a resourceful place with respect to sentiment information. From a perspective of a user, people are able to express their views through various social media, such as forums, micro-blogs, or online social networking sites [4].

With the advent of Web 2.0 techniques, users started preferring to share their opinions on the Web. These user-generated and sentiment-rich data are valuable to many applications like credibility analysis of news sites on the Web, recommendation system, business and government intelligence etc. At the same time, it brings urgent need for detecting overall sentiment inclinations of documents generated by users, which can be treated as a classification problem. Sentiment analysis includes several subtasks which have seen a great deal of attention in recent years:

1. To detect whether a given document is subjective or objective.
2. To Identify whether given subjective document express a positive opinion or a negative opinion.
3. To determine the sentiment strength of a document, such as strongly negative, weakly negative, neutral, weakly positive and strongly positive.

In this work we are focusing on second subtask.

Besides individuals on social media marketers also need to monitor all media for information related to their brands whether its for public relations activities, fraud violations, or competitive intelligence. Thus, aside from individuals, sentiment analysis is also the need of companies which are anxious to understand how their products and services are perceived by the public.

The dominating text representation method in both supervised and semi supervised sentiment classification is known as the bag-of-words (BOW) model [3], which is difficult to meet the requirements for understanding the review text and dealing with complex linguistic structures such as negation. For example, the BOW representations of two opposite reviews “It works well” and “It doesn’t work well”

are considered to be very similar by most statistical learning algorithms [2]. The two sentiment opposite texts are considered to be very similar by the BOW representation. This is exactly why standard machine learning algorithms often fail under the circumstance of polarity shift due to negation in the sentences of the review text [7] [1].

Several approaches have been proposed in the literature to address the polarity shift problem. They require either knowledge about language constructs or extra human interventions which eventually increases the complexity in classification of sentiment. Such high-level dependency on external resources makes the systems difficult to be widely used in practice. There were also some efforts to address the polarity shift problem with the absence of extra annotations and linguistic knowledge. However, results are still far from satisfactory.

MOTIVATION OF DISSERTATION

Polarity shift is a kind of linguistic phenomenon which can reverse the sentiment polarity of the text. Negation is the most important type of polarity shift [1]. For example, by adding a negation word *don't* to a positive text *I like this book* in front of the word *like*, the sentiment of the *this book* in front of the word *like*, the sentiment of the text will be reversed from positive to negative. However, the two sentiment-opposite texts are considered to be very similar by the BOW representation. This is the main reason why standard machine learning algorithms often fail under the circumstance of polarity shift.

CHAPTER 4

LITERATURE SURVEY

We studied the related work on sentiment analysis and polarity shift.

SENTIMENT ANALYSIS AND POLARITY SHIFT

According to the levels of granularity, tasks in sentiment analysis can be divided into four categorizations: document-level, sentence-level, phrase-level, and aspect-level sentiment analysis.

For document and sentence-level sentiment classification, there are two main types of methods in the literature: term-counting and machine learning methods [2] [3] [4] [14] . In term-counting methods, the overall orientation of a text is obtained by summing up the orientation scores of content words in the text, based on manually-collected or external lexical resources [6] [9]. In machine learning methods, sentiment classification is regarded as a statistical classification problem, where a text is represented by a bag-of-words; then, the supervised machine learning algorithms are applied as classifier [3]. Accordingly, the way to handle polarity shift also differs in the two types of methods.

The term-counting methods can be easily modified to include polarity shift. One common way is to directly reverse the sentiment of polarity-shifted words, and then sum up the sentiment score word by word [17]. Compared with term counting methods, the machine learning methods are more widely discussed in the sentiment classification literatures. However, it is relatively hard to integrate the polarity shift information into the BOW model in such methods. For example, Das and Chen [2] proposed a method by simply attaching NOT to words in the scope of negation, so that in the text I dont like book, the word like becomes a new word like-NOT. Yet Pang et al. [3] reported that this method only has slightly negligible effects on improving the sentiment classification accuracy.

GAP IDENTIFICATION THROUGH LITERATURE SURVEY

The following table shows the literature survey about different techniques of sentiment analysis used for classification.

Table 4.1: Literature Survey

No.	Reference	Techniques	Description
1	Dual Sentiment Analysis: Considering Two sides of one review	Support vector machine (SVM), Naive bayes, Logistic Regression	Dual training and Dual Prediction technique is used.
2	Thumbs up?Sentiment Classification using Machine learning algorithms	Learning algorithms and n-gram model	Classify the dataset using different machine.
3	Classification of sentiment reviews using N-gram machine learning approach	Support Vector Machine Naive Bayes	Converting text reviews into numeric matrices using countvectorizer and TF-IDF
4	Thumbs up or thumbs down? Semantic orientation applied to unsupervised classification of reviews	Unsupervised learning algorithm for classifying a review	A specific unsupervised learning technique based on the mutual information
5	Automatic Opinion polarity Classification of movie	Naive Bayes And Markov Model (MM)	Accessed overall opinion polarity(OvOp) concept using machine learning algorithm
6	Dual Training and dual prediction for polarity classification	SVM and Naive Bayes	Dual training and dual prediction (DTDP)

Table 4.2: Literature Survey

CHAPTER 5

PROBLEM DEFINITION AND SCOPE

GOALS

- Understanding existing sentiment analysis approaches.
- Study corpus based, lexical based and semantic based techniques.
- Understanding unigram, bigram, trigram and combination of them for modeling purpose.
- Training the model with naive bayes, support vector machine, maximum entropy.
- Applying this learned model to the test dataset.
- Evaluating the results generated by classifiers.

OBJECTIVES

Please refer Chapter 1, Section 1.7 on Page 2

STATEMENT OF SCOPE

- Preprocessing the reviews
- Classify reviews into two polarities.
- Evaluate the classification accuracy by each classifier.

SOFTWARE CONTEXT

Scikit-learn

Scikit-learn (formerly scikits.learn) is a free software machine learning library for the Python programming language. It features various classification, regression and clustering algorithms including support vector machines, random forests, gradient boosting, k-means and DBSCAN, and is designed to interoperate with the Python numerical and scientific libraries NumPy and SciPy.

NumPY

NumPy is the fundamental package for scientific computing with Python. It contains among other things:

- A powerful N-dimensional array object
- Sophisticated (broadcasting) functions
- Tools for integrating C/C++ and Fortran code
- Useful linear algebra, Fourier transform, and random number capabilities

Besides its obvious scientific uses, NumPy can also be used as an efficient multi-dimensional container of generic data. Arbitrary data-types can be defined. This allows NumPy to seamlessly and speedily integrate with a wide variety of databases.

Natural language toolkit

NLTK is a leading platform for building Python programs to work with human language data. It provides easy-to-use interfaces to over 50 corpora and lexical resources such as WordNet, along with a suite of text processing libraries for classification, tokenization, stemming, tagging, parsing, and semantic reasoning, wrappers for industrial-strength NLP libraries, and an active discussion forum.

Matplotlib

Matplotlib is a Python 2D plotting library which produces publication quality figures in a variety of hardcopy formats and interactive environments across platforms. Matplotlib can be used in Python scripts, the Python and IPython shell, the jupyter notebook, web application servers, and four graphical user interface toolkits. You can generate plots, histograms, power spectra, bar charts, errorcharts, scatterplots, etc., with just a few lines of code.

CHAPTER 6

DISSERTATION PLAN

PURPOSE OF THE DOCUMENT

This document specifies and estimates various risks associated with this project and states how they are handled. It also states the project plan in terms of task and their dependency.

TECHNICAL CONSTRAINTS

- To build a classification module that distributes data and execution among spark executors.
- To fetch candidate's tweets in a CSV format and store it in Alluxio.

DISSERTATION ESTIMATES

Reconciled Estimates

Cost Estimates

No cost is required for tools and software as open source softwares are used.

Time Estimates

Calendar time required: 11 months.

Dissertation Resources

- People : Single Person
- Software resources used are mentioned in Chapter 7, Section 7.5.1 on Page 26
- Hardware resources used are mentioned in Chapter 7, Section 7.5.2 on Page 26

RISK MANAGEMENT

This section discusses dissertation risks and the approach to managing them.

Risk Identification

For risks identification, review of scope document, requirement specifications and schedule is done. Answers to questionnaire revealed some risks. Following risk identification questionnaire has been referred.

- Are requirements fully understood by the software engineering team and its customers?
- Have customers been involved fully in the definition of requirements?
- Do end-users have realistic expectations?
- Does the software engineering team have the right mix of skills?
- Are project requirements stable?
- Is the number of people on the project team adequate to do the job?
- Do all customer/user constituencies agree on the importance of the project and on the requirements for the system/product to be built?

Risk Analysis

The risks for the dissertation are analyzed within the constraints of time and quality. Risk can be as follows:

- Out of memory error, while creating model and training the model.
- Review text contains unrecognized characters.
- Out of memory error, while predicting on the test dataset.

Please refer Table 6.1, 6.2 and 6.3 for detail description.

Table 6.1: Risk Table

ID	Risk Description	Probability	Impact		
			Schedule	Quality	Overall
1	Out of Memory	High	Low	High	High
2	Unrecognized characters	Low	Medium	High	Medium
3	Out of Memory	Low	Medium	High	High

Table 6.2: Risk Probability Definitions

Probability	Value	Description
High	Probability of the occurrence is	>75%
Medium	Probability of the occurrence is	26% - 74%
Low	Probability of the occurrence is	25%

Table 6.3: Risk Impact Definitions

Impact	Value	Description
Very High	>10%	Schedule impact or Unacceptable quality
High	5%-10%	Schedule impact or Some parts of the project have low quality
Low	<5%	Schedule impact or Barely noticeable degradation in quality Low Impact on schedule or Quality can be incorporated

Overview of Risk Mitigation, Monitoring and Management

Please refer Table 6.4, 6.5 and 6.6 for detail description.

Table 6.4: Risk 1

Risk ID	1
Risk Description	Out of memory error, when training Model
Category	Configuration
Source	Software Requirement Specification Document
Probability	High
Impact	High
Response	Mitigate
Strategy	Changing number of features resolves this issue.
Risk Status	Occurred and Resolved

Table 6.5: Risk 2

Risk ID	2
Risk Description	Unreconized characters
Category	Configuration
Source	Software Requirement Specification Document
Probability	Low
Impact	Low
Response	Mitigate
Strategy	Convert all characters into unicode format
Risk Status	Occurred and Resolved

Table 6.6: Risk 3

Risk ID	3
Risk Description	Out of memory error, when predicting
Category	Development Environment
Source	Software Requirement Specification Document
Probability	Low
Impact	Low
Response	Mitigate
Strategy	Using sparse matrix to represent text
Risk Status	Occurred and Resolved

STAFF ORGANIZATION

Team Structure

- Internal Guide : Prof. A. G. Phakatkar
- Student : Kaushik S. Hande

Management Reporting and Communication

The progress of dissertation is reported once in a month.

Dissertation Task Set

Major tasks in the Dissertation stages are -

Task 1 : Requirement

1. Define problem statement
2. Identify scope, requirements
3. Related mathematical model

Task 2 : Design

1. Identifying of key objects, functional relation
2. UML diagrams and functional dependency graph
3. System design

Task 3 : Implementation

1. Import dataset
2. Classifier implementation
3. Comparison using graphs

Task 4 : Testing

1. Unit testing
2. Integration testing
3. System testing

Task 5 : Integration and Maintenance

1. Integration
2. Maintenance

Please refer figure 6.1 Task Network.

Task Network

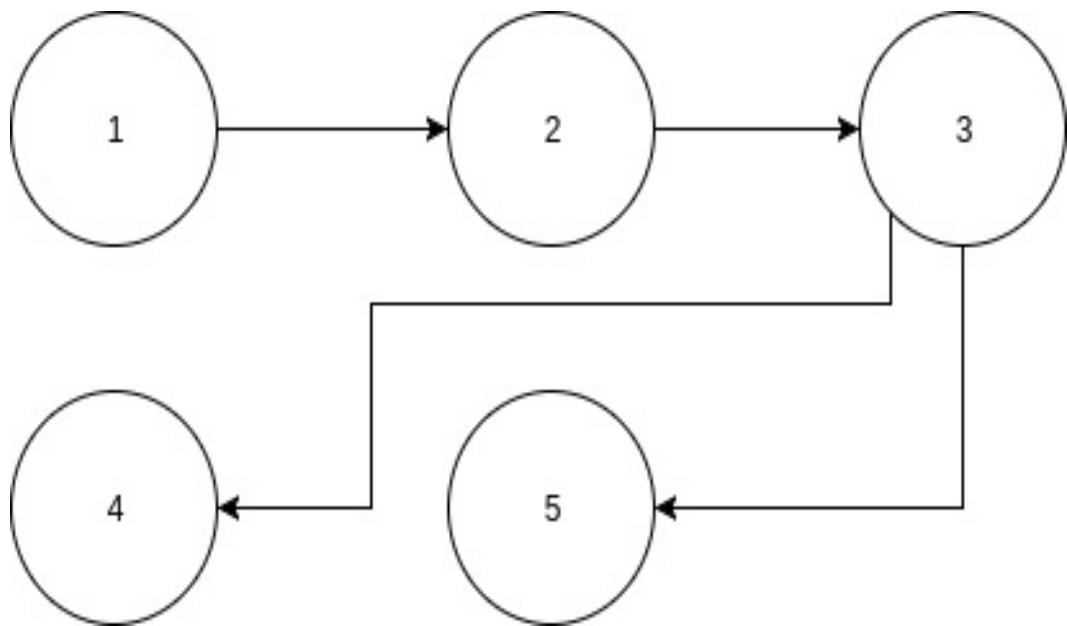


Figure 6.1: Task Network

Timeline Chart

Please refer Annexure B, Table B.1 on Page 66 for all Dissertation Tasks.

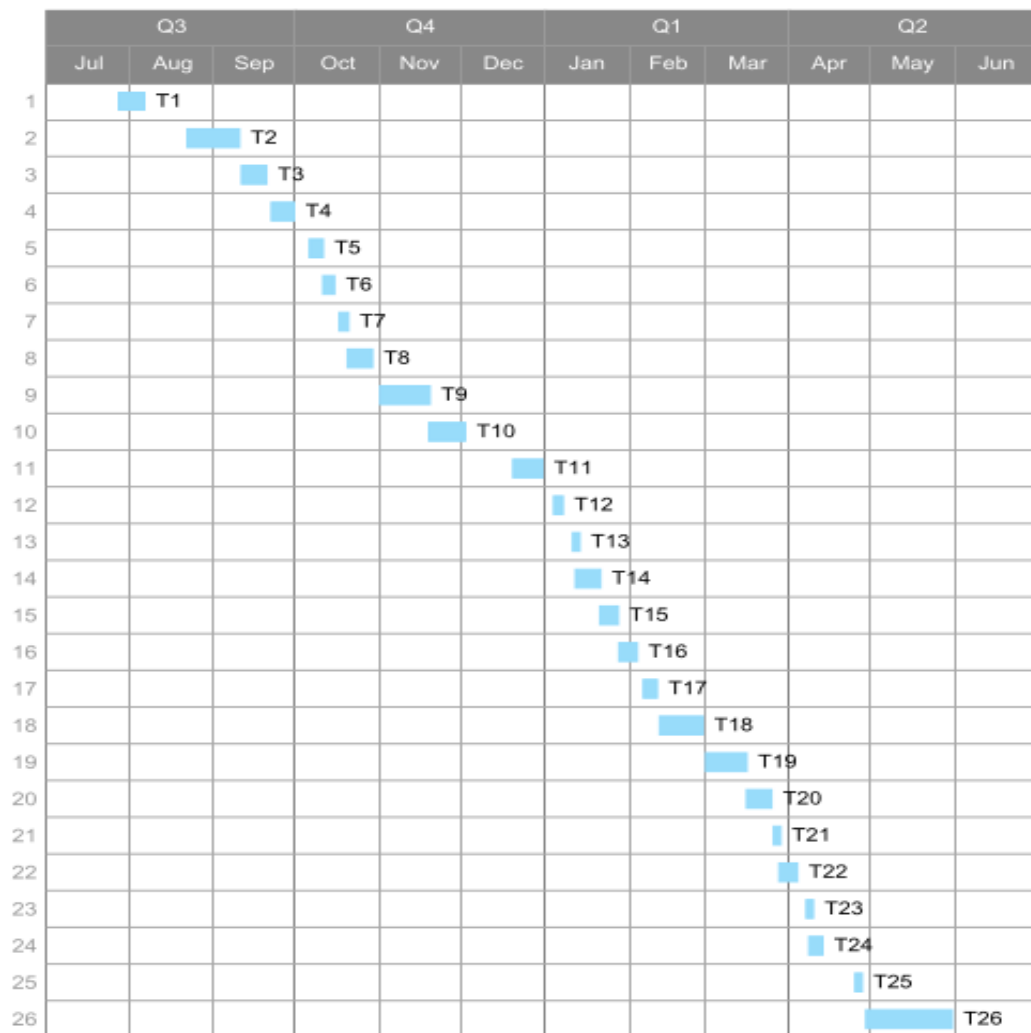


Figure 6.2: Timeline Chart

CHAPTER 7

SOFTWARE REQUIREMENT

SPECIFICATION

INTRODUCTION

The aim of this document is to specify the software requirements for classification of movie reviews.

PURPOSE AND SCOPE OF THE DOCUMENT

The purpose of the document is to enlist various software requirements to build the system. This document has functional and non-functional requirements for the software being developed.

OVERVIEW OF RESPONSIBILITIES OF DEVELOPER

The responsibilities of a developer includes gathering of information about the classification libraries, that can be used to design and develop the system to categorize movie reviews. The developers responsibilities include:

- Planning for dissertation (Scheduling)
- Designing of system (High Level Design Document)
- Coding of system (Implementation)
- Testing of system (Test Cases)

PRODUCT OVERVIEW

System builds classifier models for classification of reviews. Different functionality of the system are :

- Review loader : It loads the reviews into python environment.
- Stopwords remover : It removes the stopwords like the words is, a, an, the, was, were etc.
- Vectorizer : It converts the reviews into matrix of rows and columns where columns represents the words and rows represents each reviews. Presence of particular word in a review is shown by the values in the columns.
- Classifier : It classifies the reviews into positive and negative review.

HARDWARE RESOURCES USED

Software Requirements

- Python 2.7.6
- Rstudio Version 0.99.893
- R version 3.3.2
- Operating Systems:
 - Windows XP, 7, 8, 10
 - Linux(Any flavor)
 - Mac OS

Hardware Requirements

- Intel(R) Core(TM) i3 CPU @ 2.90GHz or later, width : 64 bits
- Memory : 4 GB DDR3 or more
- Capacity : 1697MHz or more
- Cores : 4 or more
- PCI Express Gigabit Ethernet Controller, Size: 100Mbit/s, Capacity: 1Gbit/s, Width: 64 bits
- Hard Disk : 500 GB (EXT4 Primary/Logical Partition)

FUNCTIONALITY

- Download movie reviews from IMDB dataset.
- Import the movie review dataset into python environment using csv package.
- Convert the text reviews into matrix form.
- Remove the stopwords from reviews.
- Show positive and negative polarity score for test reviews.
- Compare classifiers for accuracy of classification.

INPUT

- Dataset that consists of movie reviews and their corresponding labels.
- List of stopwords which play no role in classification.

OUTPUT

- Classification of each test review into positive or negative.
- Percentage of accuracy achieved in classification.
- Comparison of accuracies obtained by each classifier.

MAJOR CONSTRAINTS

- To store movie reviews as input in csv file format.
- To execute classifiers in configured environment.
- To train the model for polarity classification.

APPLICATIONS

- Businesses and organisations which require consumer opinions to do with products they market and services they produce.

- Individuals who make decisions to purchase products or services based upon word of mouth or on-line reviews, or to find public opinion, e.g. concerning politics or local issues.
- On-line advertising where in social media, an organisation may place an advertisement in response to a favourable review of a product, or a rival product could be advertised upon receipt of a bad review
- Opinion retrieval for general searches of opinions
- HR Analytics.

USAGE SCENARIO

A use case represents a particular functionality of a system. Hence, use case diagram is used to describe the relationships among the functionalities and their internal/external actors. This section provides various usage scenarios for the system to be developed.

User Profiles

Actors of the system are user and IMDB movie dataset system. IMDB dataset provides input to the system. System downloads the dataset from given url. It performs text analysis on each row for that stopwords removal, vectorizer and sentiment analysis is used. The predicted class of the reviews is accessed by the user.

Use Case Views

System

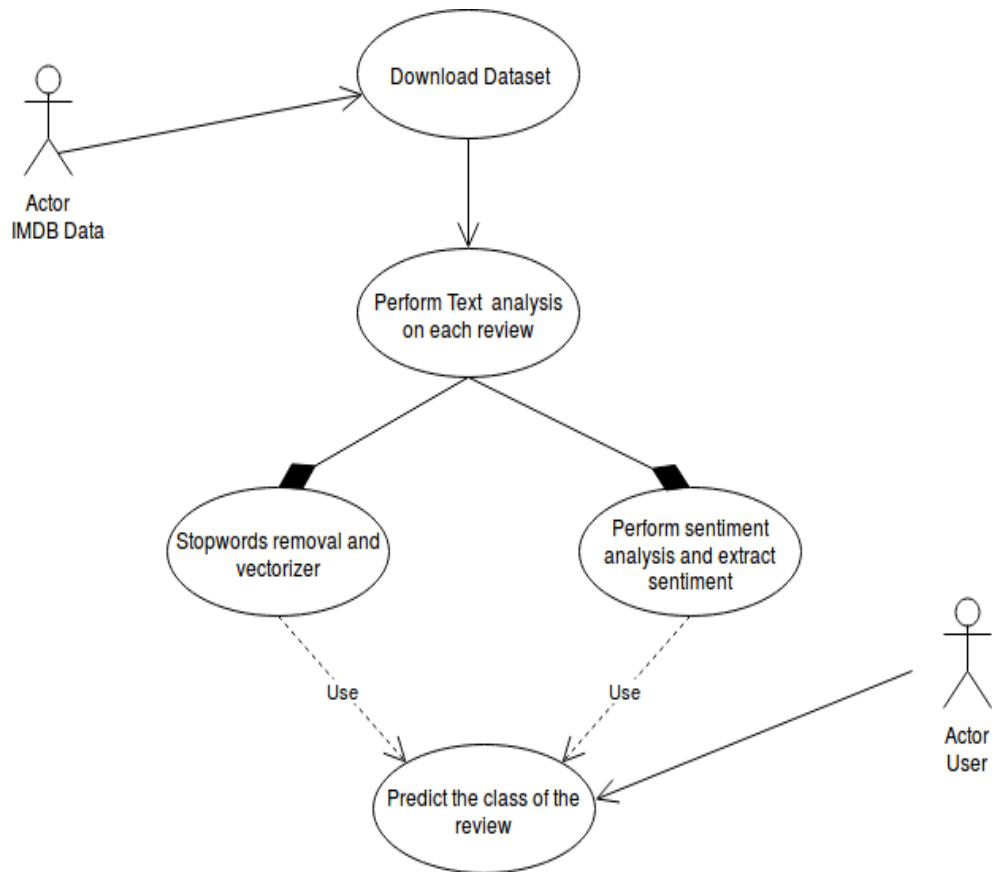


Figure 7.1: Use Case : System

MODEL AND DESCRIPTION

This section contains details about events and associated behaviour of the system which is shown using diagram below.

Activity Diagram

Activity diagram is a flow chart to represent the flow from one activity to another activity. The activity can be described as an operation of the system. The control flow is drawn from one operation to another. This flow can be sequential, branched or concurrent. The purpose of activity diagrams is to capture the dynamic behaviour of the system.

Description : As shown in figure 8.6, User downloads movie dataset. This dataset is added for processing using csv import library. Documents are preprocessed to remove stopwords and reviews are converted into matrix. Classifiers are trained after

feature extraction. This trained classifiers are used on test dataset to classify.

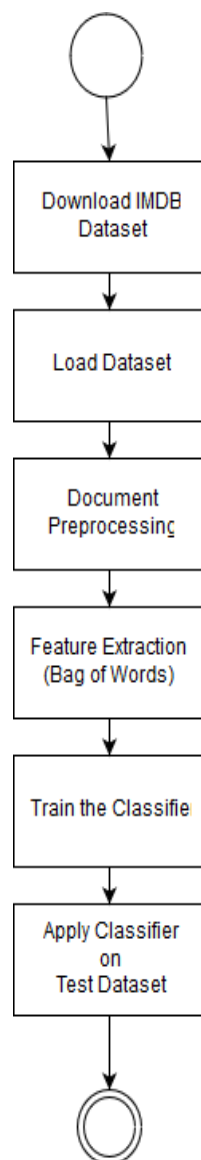


Figure 7.2: Activity Diagram

FUNCTIONAL MODEL AND DESCRIPTION

This section describes data flow diagrams (DFD) of the proposed system. There are three types of DFDs explained in the section. These diagrams explain the system in brief.

Data Flow Diagram

Level 0 Data Flow Diagram

In the level 0 DFD as shown in figure 7.7, Downloaded reviews are loaded in sentiment analysis system. The system analyzes and presents report through charts to user

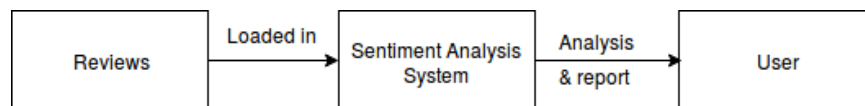


Figure 7.3: Level 0 DFD

Level 1 Data Flow Diagram

In the level 1 DFD as shown in figure 7.8, Reviews are loaded through csv file loader. Data is preprocessed and model is created. Analysis done by the system is given to the user.

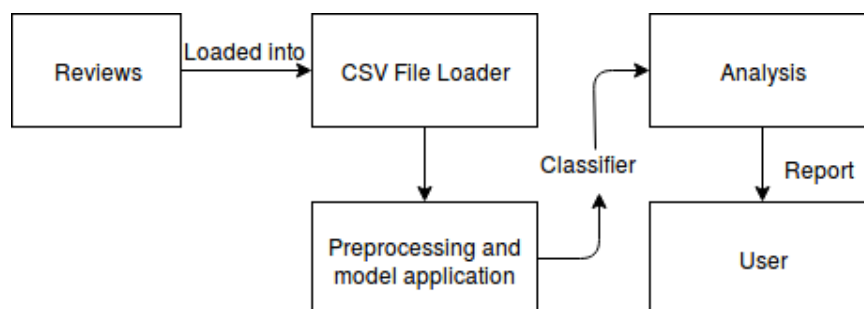


Figure 7.4: Level 1 DFD

Level 2 Data Flow Diagram

In the level 2 DFD as shown in figure 7.9, Reviews are preprocessed to give bag of words features. Classifiers are trained using this model. Classification is performed

on the test dataset. Detailed analysis of results obtained by classifiers is done. It is provided to the user.

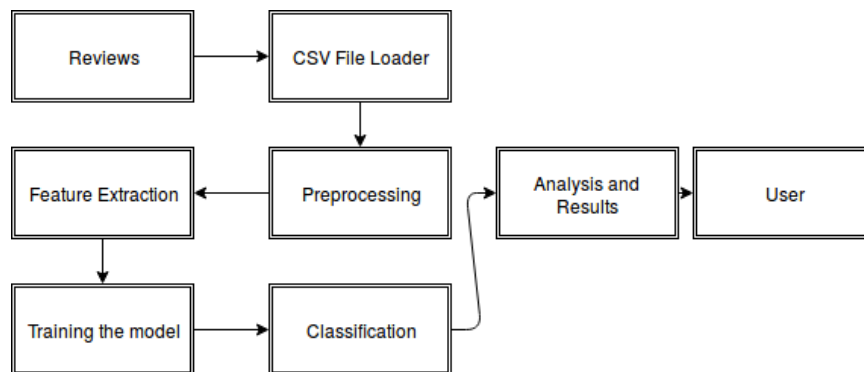


Figure 7.5: Level 2 DFD

Sequence Diagram

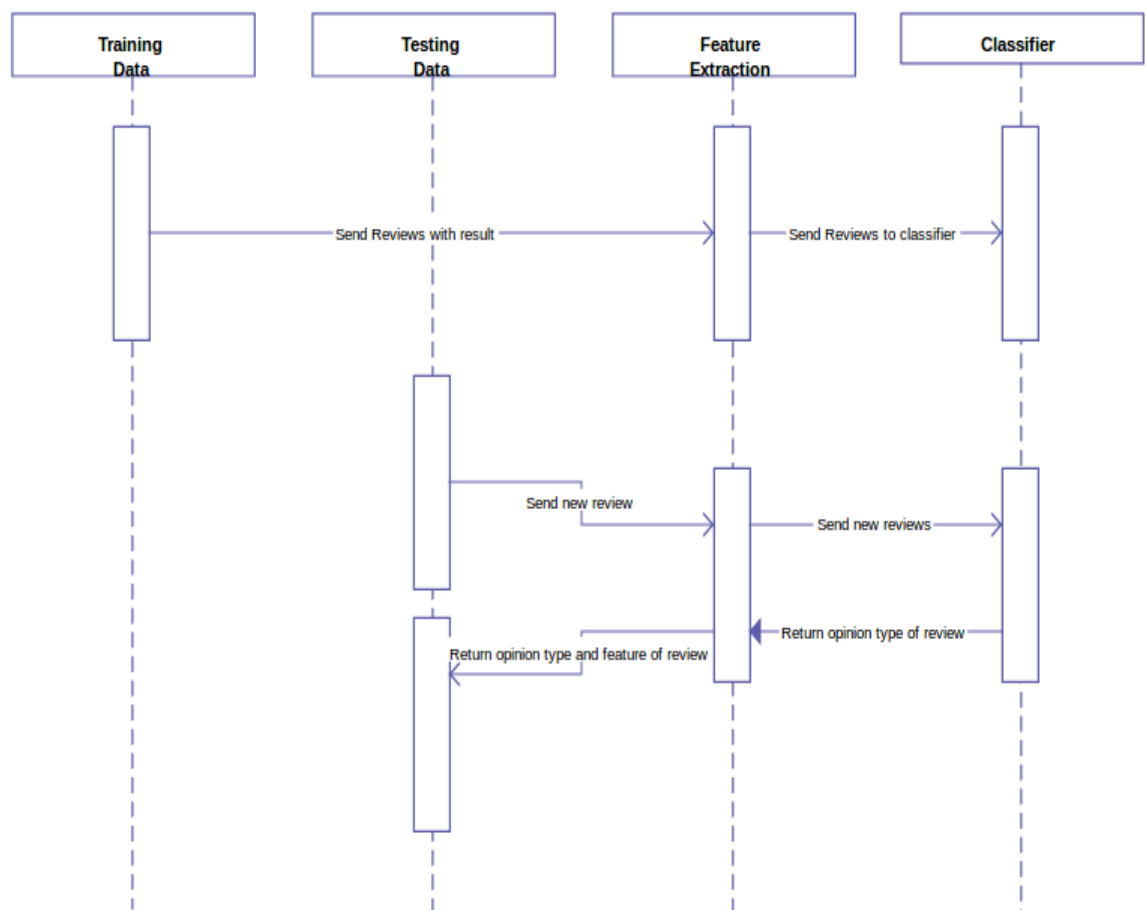


Figure 7.6: Sequence Diagram

NON-FUNCTIONAL REQUIREMENTS

Availability

Required libraries must be installed and loaded in the python environment with the required configurations. Dataset must be downloaded from specified url [18].

Scalability

The system should be scalable to classify reviews even if the training and test data are increased. System can comfortably handle reviews dataset upto 25000 reviews.

Performance

The system must be interactive and delays involved must be less. There should be no immediate delays for every action and response of the system. Training time increases as the training data increases. It takes 4 to 5 seconds in training the dataset. Training increases further when bigram and trigram models are used.

Usability

The system should be easy to handle and process requests efficiently. System's functions are designed to use with ease and provide results. Results are presented in the form of graphs and are easy to comprehend.

Reliability

The system should efficiently analyze movie reviews entirely and give correct classification result. It should be reliable to perform classification effectively on any review dataset.

Maintainability and Changeability

The system is made up of different independent modules that can be modified to correct faults, improve performance or other attributes, or adapt to a changed environment. System can be improved for new features and will be able to include new requirements.

CHAPTER 8

DETAILED DESIGN DOCUMENT

INTRODUCTION

This document specifies the design that is used to fetch movie reviews in CSV format, classifies individual movie review into polarity categories.

- **Polarity Categories**

- Positive
- Negative

Candidate's tweets are collected and classified based on these categories. They are profiled and compared based on emotional and polarity categories.

ARCHITECTURAL DESIGN

Figure 8.1 shows architectural design of proposed system. Following are important components in the system :

- **Movie reviews data :** It contains 1,000 positive and 1,000 negative movie reviews from IMDB.
- **Preprocessing :** It has stopwords removal and vectorizer.
- **Training data and Test Data :** Data is divided into training data and test data. Training data consists of 75 % of data and test data consists of 25 % of total data. Both are mutually exclusive.
- **Classifier :** Training data is given as a input to one of the classifier.
- **Prediction system :** It takes test data and applies trained model to it .
- **Sentiment :** It gives the end output of polarity into positive and negative classification.

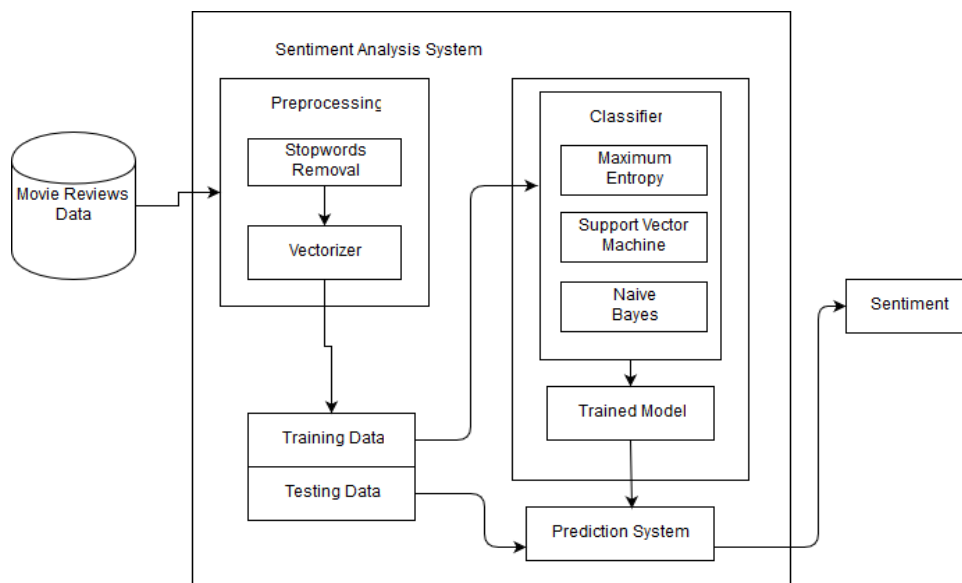


Figure 8.1: Proposed System Architecture

CLASS DESIGN

It is a static diagram that represents the static view of an application. It is not only used for visualizing, describing, and documenting different aspects of a system but

also for constructing executable code of the software application. It describes the attributes and operations of a class and also the constraints imposed on the system.

Description : In figure 9.2, modules and their relationships are shown. Document classifier used for classifying SentimentClassification class has ClassifyReview method which predicts the polarity of review.

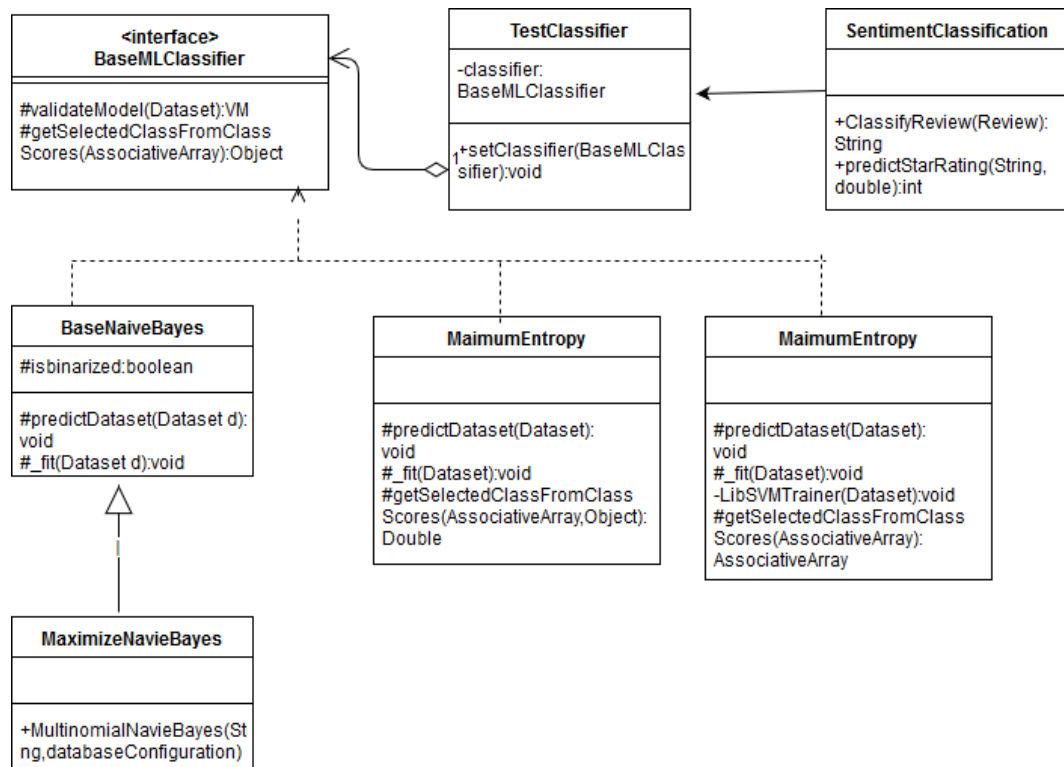


Figure 8.2: Class Diagram

COMPONENT DESIGN

It is used to model the physical aspects of a system. It is also used to visualize the organization and relationships among components in a system. It does not describe the functionality of the system but it describes the components used to make those functionalities.

Description : Figure 8.3 describes primary components of the system. Movie dataset is stored in a component which then gives to classifier system. Classifier system then preprocess the reviews and does the prediction of polarity and sends the output to the console subsystem where result evaluation and comparison takes place.

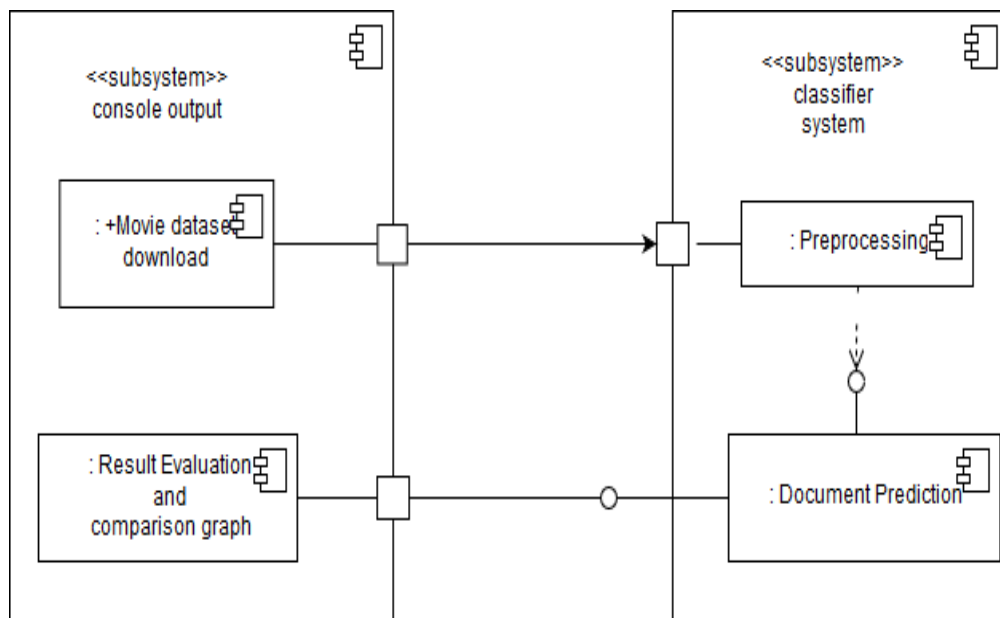


Figure 8.3: Component Diagram

CHAPTER 9

IMPLEMENTATION DETAILS

INTRODUCTION

This section describes implementation of the system, required libraries and dependencies needed for components of the system and use of implementation strategy.

ALGORITHM

Document Classification

Input : Movie review dataset in .csv file format

Output : Polarity classification

1. Initialize review_data list to empty
2. Initialize review_label list to empty
3. Read csv file
4. Put first column of each row in csv in review_data
5. Put second column of each row in csv in review_label
6. Specify ngram range i.e. unigram, bigram or trigram
7. Fit the review_data into matrix form using countvectorizer which stores counts of each word in a review.
8. Transform review_data to tfidf using tfidftransformer.
9. Split review_data and review_label into review_train, label_train, review_test and label_test
10. Give review_train and label_train as input to the classifier and store learned model in classifier variable.
11. Apply review_test on classifier model with predict method and take predicted labels output in predicted variable
12. Compare label_test with predicted variable label values
13. Print the confusion matrix

DATASET

Pang and Lee's Movie Review Data was one of the first widely-available sentiment analysis datasets. It contains 1,000 positive and 1,000 negative movie reviews from IMDB, so it is now considered too small for serious research and development purposes. However, the text is similar to movies reviews on IMDB today.

The file movie-pang02.zip contains a copy of Pang and Lee's Movie Review Data in a csv format that can be imported directly in python. It has two categories: Pos (reviews that express a positive or favorable sentiment) and Neg (reviews that express a negative or unfavorable sentiment). For this assignment, we will assume that all reviews are either positive or negative; there are no neutral reviews.

For document classification, a training and testing dataset is required. Training records for polarity categories are mentioned in Table 9.1

Table 9.1: Polarity Training Dataset

Polarity	Training Records
Positive	1000
Negative	1000
Total	2000

SNAPSHOTS

```
raphael@raphael-Inspiron-N4010: ~/Documents/6107/Dissertation/pychennai-sentiment-al
ction module into which all the refactored classes and functions are moved. Also
note that the interface of the new CV iterators are different from that of this
module. This module will be removed in 0.20.
"This module will be removed in 0.20.", DeprecationWarning)
/usr/lib/python2.7/dist-packages/scipy/sparse/compressed.py:130: VisibleDeprecat
ionWarning: `rank` is deprecated; use the `ndim` attribute or function instead.
To find the rank of a matrix see `numpy.linalg.matrix_rank`.
    if np.rank(self.data) != 1 or np.rank(self.indices) != 1 or np.rank(self.indpt
r) != 1:
/usr/lib/python2.7/dist-packages/scipy/sparse/coo.py:200: VisibleDeprecationWarn
ing: `rank` is deprecated; use the `ndim` attribute or function instead. To find
the rank of a matrix see `numpy.linalg.matrix_rank`.
    if np.rank(self.data) != 1 or np.rank(self.row) != 1 or np.rank(self.col) != 1
:
(0, 3)      2
(0, 0)      1
(0, 2)      1
(0, 4)      1
(0, 5)      1
(1, 1)      1
(1, 6)      1
(1, 3)      1
(1, 4)      1
>
```

Figure 9.1: Sample matrix creation with word counts

```
raphael@raphael-Inspiron-N4010: ~/Documents/6107/Dissertation/pychennai-sentiment-al
ing: `rank` is deprecated; use the `ndim` attribute or function instead. To find
the rank of a matrix see `numpy.linalg.matrix_rank`.
    if np.rank(self.data) != 1 or np.rank(self.row) != 1 or np.rank(self.col) != 1
:
(0, 3)      2
(0, 0)      1
(0, 2)      1
(0, 4)      1
(0, 5)      1
(1, 1)      1
(1, 6)      1
(1, 3)      1
(1, 4)      1
>
(0, 3)      0.707106781187
(0, 0)      0.353553390593
(0, 2)      0.353553390593
(0, 4)      0.353553390593
(0, 5)      0.353553390593
(1, 1)      0.5
(1, 6)      0.5
(1, 3)      0.5
(1, 4)      0.5
>
```

Figure 9.2: Sample matrix creation with tf-idf

```

=====
A great movie. Lansbury and Tomlinson are perfect, the songs are wonderful, the
dances, with a particular mention for the "Portobello Ballet" are gorgeous. As f
or the animated section, the match between animals has become an instant classic
; the climax with the attack of the armatures is chilling and fascinating. I rec
ommend to see the restored 134 minutes version or at least the 112 minutes video
. Here in Italy we have only the 98 minutes version, although the film was prese
nted in its original release at the running of 117 minutes. If possible, watch a
lso the German videocassette: it was generated from the 98 minutes running but i
t's missing of every refer to World War II and of all the scenes between English
people and their Nazi invaders! positive positive
[ 0.00411253  0.99588747]
=====
Why oh why don't blockbuster movies simply stick to their selling point? Everyon
e in the cinema, young and old, was there to see talking animals make jokes, and
whilst they did that we were all happy... And then, as with Lost In Space, came
the two killer blows - plot and sentiment. Who really cared what happened to th
e tiger or whether Eddie Murphy made up with his daughter? Not me, that's for su
re. negative negative
[ 9.99933065e-01  6.69352950e-05]
=====
What happened? What we have here is basically a solid and plausible premise and
with a decent and talented cast, but somewhere the movie loses it. Actually, it
never really got going. There was a little excitement when we find out that Angi

```

Figure 9.3: Review classification with negative and positive percentages

```

kaushik@subhash-OptiPlex-9020: ~/Documents/6107/to_be_mailed
package will not be available.
warnings.warn("The twython library has not been installed. "
-----Unigram features-----
1467
491
[[195  56]
 [ 42 198]]
The accuracy score is 80.04%

-----
SINGLE FOLD RESULT (Naïve Bayes)
-----
('accuracy:', 0.8004073319755601)
('precision', 0.8011561845908501)
('recall', 0.8009462151394422)
('f-measure', 0.8003998805336165)
[[210  41]
 [ 62 178]]
The accuracy score is 79.02%

-----
SINGLE FOLD RESULT (SVM)
-----
('accuracy:', 0.790224032586558)

```

Figure 9.4: Three classifiers results(Accuracy, recall, precision, f-measure with confusion matrix)

```
kaushik@subhash-OptiPlex-9020: ~/Documents/6107/to_be_mailed
('f-measure', 0.8003998805336165)
[[210  41]
 [ 62 178]]
The accuracy score is 79.02%

-----
SINGLE FOLD RESULT (SVM)
-----
('accuracy:', 0.790224032586558)
('precision', 0.7924221058286328)
('recall', 0.789160026560425)
('f-measure', 0.7893292009814336)
[[226  25]
 [ 84 156]]
The accuracy score is 77.80%

-----
SINGLE FOLD RESULT (Maximum Entropy)
-----
('accuracy:', 0.7780040733197556)
('precision', 0.7954553555515951)
('recall', 0.775199203187251)
('f-measure', 0.7733983682006595)
```

Figure 9.5: Three classifiers results(Accuracy, recall, precision, f-measure with confusion matrix)

CHAPTER 10

TEST SPECIFICATION

CHAPTER 11

DATA TABLES AND DISCUSSIONS

Table 11.1: Unigram Features Results

Classifier	Accuracy	Precision	Recall	F-measure
Naive Bayes	80.04	80.11	80.09	80.03
Support Vector Machine	79.02	79.24	78.91	78.93
Maximum Entropy	77.80	79.54	77.51	77.33

Table 11.2: Bigram Features Results

Classifier	Accuracy	Precision	Recall	F-measure
Naive Bayes	81.26	81.43	81.34	81.25
Support Vector Machine	80.44	80.49	80.39	80.41
Maximum Entropy	80.85	82.13	80.62	80.57

Table 11.3: Trigram Features Results

Classifier	Accuracy	Precision	Recall	F-measure
Naive Bayes	80.85	80.90	80.90	80.85
Support Vector Machine	80.04	80.34	79.92	79.93
Maximum Entropy	77.59	79.52	77.30	77.09

The above Table 11.1, Table 11.2, Table 11.3 shows accuracy, precision, recall and f-measure of three different classifiers. Three different features are used as features. First table shows results with unigram features i.e. each words is considered separate features. Second table shows results with bigram features i.e. 2 words together are considered as features. Third table shows results with trigram features i.e. 3 words combine are considered as features.

Below fig. 11.1 shows comparison of accuracy scores shown in Table 11.1, Table 11.2, Table 11.3 .

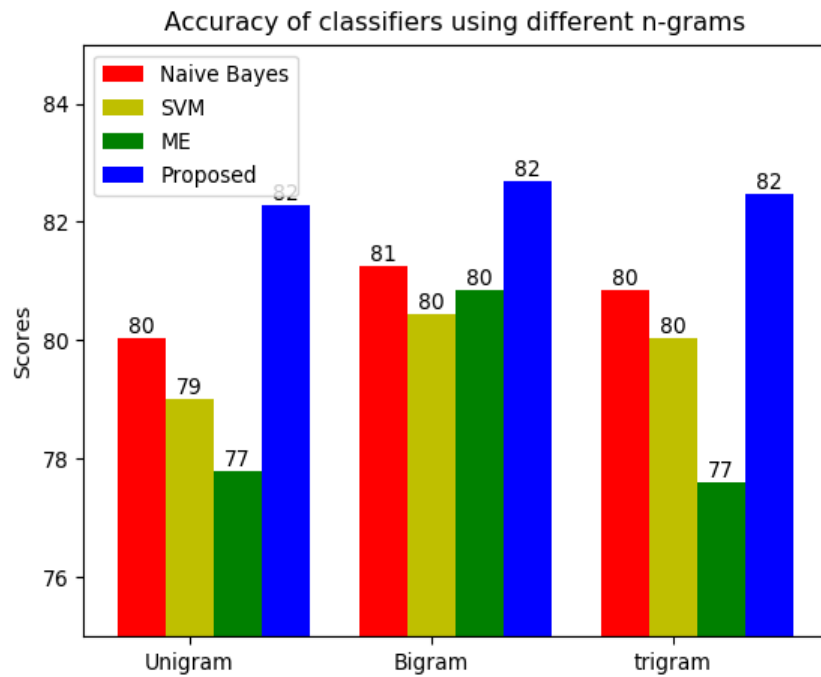


Figure 11.1: Classifiers accuracy score comparison

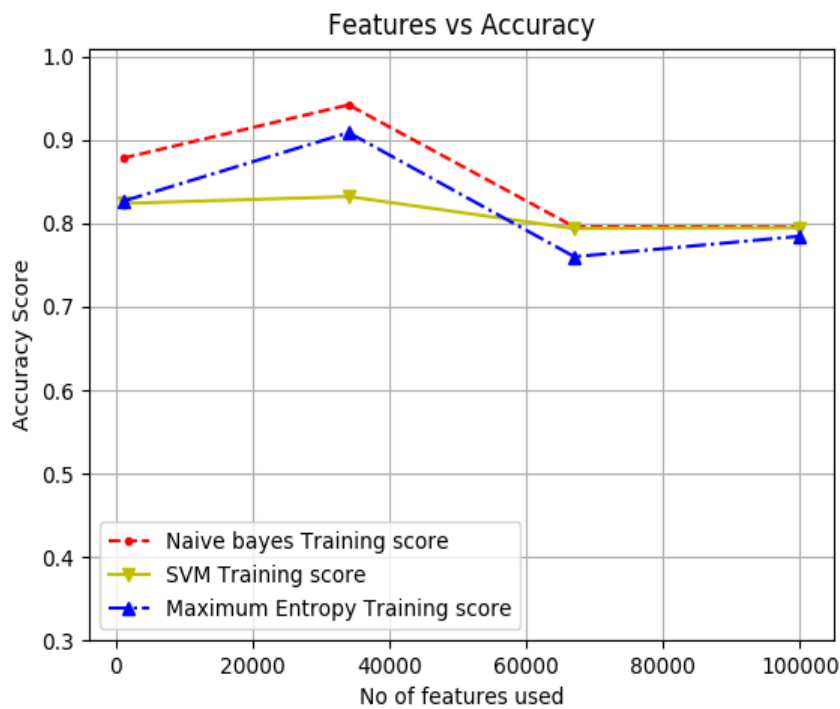


Figure 11.2: Variations in classifier accuracies as the no. of features increases

CHAPTER 12

CONCLUSION

In this work, the Dual sentiment analysis technique is used to address the polarity shift problem in sentiment classification. The idea behind this is to create opposite reviews of the original reviews and used them together in Bag of words model which will be the feature inputs to various machine learning algorithms such as Naive Bayes, SVM and Maximum Entropy. It performs better than the methods which were used to address the problem of polarity shift due to negation .

CHAPTER 13

FUTURE ENHANCEMENTS

The research can be extended to include the neutral reviews as well. Also antonym can be obtained using words that has opposite polarity in the corpus. The corpus-based pseudo-antonym dictionary is also good at obtaining more domain-relevant antonym words by learning from the corpus.

CHAPTER 14

REFERENCES

- [1] R. Xia, F. Xu, C. Zong, Q. Li, Y. Qi and T. Li, “ Dual sentiment analysis: Considering two sides of one review, ” in *IEEE transactions on knowledge and data engineering*, vol. 27, no. 8, pp. 2120 - 2133, Aug. 2015.
- [2] S. Das and M. Chen, “ Yahoo! for Amazon: Sentiment extraction from small talk on the web, ” *Management science* , Vol.53, Issue no.9, pp.1375-1388, 2007.
- [3] Pang, L. Lee, and S. Vaithyanathan, “ Thumbs up?: Sentiment classification using machine learning techniques, ” *Proceedings of the ACL-02 conference on Empirical methods in natural language processing*, pp. 79-86, 2002.
- [4] B. Pang and L. Lee, “ Opinion mining and sentiment analysis, ” *Foundations and Trends in Information Retrieval*, vol. 2, no. 1-2, pp. 1-135, 2008.
- [5] R. Xia, T. Wang, X. Hu, S. Li, and C. Zong, “ Dual Training and Dual Prediction for Polarity Classification,” *Proceedings of the Annual Meeting of the Association for Computational Linguistics (ACL - 02)*, pp. 521-525, 2013.
- [6] P. Turney, “ Thumbs up or thumbs down? Semantic orientation applied to unsupervised classification of reviews, ” *Proceedings of the Annual Meeting of the Association for Computational Linguistics (ACL)*, pp. 417-424, 2002.
- [7] M. Li and C. Huang, “ Sentiment classification considering negation and contrast transition, ” *Proceedings of the Pacific Asia Conference on Language, Information and Computation (PACLIC)*, pp. 307-316, 2009.
- [8] Li, S. Lee, Y. Chen, C. Huang and G. Zhou, “ Sentiment Classification and Polarity Shifting, ” *Proceedings of the International Conference on Computational Linguistics (COLING)*, pp. 635-643, 2010.
- [9] D. Turney and Michael L. Littman, “ Un-supervised learning of semantic orientation from a hundred-billion-word corpus. ” *Technical Report EGB-1094, National Research Council Canada*, arXiv preprint cs/0212012, 2002.

- [10] Yuan Wang, Zhaohui Li, Jie Liu, Zhicheng He, Yalou Huang and Dong Li, “ Word Vector Modeling for Sentiment Analysis of Product Reviews ” *Natural Language Processing and Chinese Computing 2014*, pp. 168-180, 2014.
- [11] Xia, Rui and Wang, Cheng and Dai, Xinyu and Li, Tao, “ Co-training for Semi-supervised Sentiment Classification Based on Dual-view Bags-of-words Representation ” *Association for Computational Linguistics (ACL 1)*, pp. 1054-1063, 2015.
- [12] Na, J.C., Sui, H., Khoo, C., Chan, S., and Zhou, Y., “ Effectiveness of Simple Linguistic Processing in Automatic Sentiment Classification of Product Reviews ” *Proceedings of the Eighth International ISKO Conference* pp. 49-54, 2004.
- [13] Rui Xia, Feng Xu, Jianfei Yu, Yong Qi and Erik Cambria, “ Polarity shift detection, elimination and ensemble: A three-stage model for document-level sentiment analysis ” *Information Processing & Management* 52, no. 1, pp. 36 - 45, 2016
- [14] Rui Xia, Chengqing Zong and Shoushan Li, “ Ensemble of feature sets and classification algorithms for sentiment classification ” *Information Sciences* 181, no. 6 pp. 1138-1152, 2011
- [15] Anderson Uilian Kauer and Viviane P. Moreira, “ Information retrieval for sentiment polarity prediction, ” *Expert Systems With Applications* 61, pp. 282 - 289, 2016
- [16] Yuming Lin, Jingwei Zhang, Xiaoling Wang and Aoying Zhou, “ An Information Theoretic Approach to Sentiment Polarity Classification ” *Proceedings of the 2nd joint WICOW/AIRWeb workshop on web quality, ACM Lyon France* , pp. 35 - 40, 2012
- [17] L. Polanyi and A. Zaenen, “Contextual lexical valence shifters,” in *Proc. AAAI Spring Symp. Exploring Attitude Affect Text*, pp. 110, 2004.
- [18] <http://boston.lti.cs.cmu.edu/classes/95-865-K/HW/HW3/movie-pang02.zip>

ANNEXURE A

PAPERS PUBLISHED

PAPER TITLE

Sentiment Analysis using Machine Learning Algorithms: A Survey

IJIRCCE Certification



Figure A.1: IJIRCCE Certificate

PAPER TITLE

Sentiment Analysis using Original and Reversed Reviews

cPGCON Certificate



Figure A.2: cPGCON Certificate

cPGCON Review

Sandip Foundation's Sandip Institute of Technology and Research Centre
Department of Computer Engineering
in association with Board of Studies Computer Engineering
Savitribai Phule Pune University, Pune

Evaluation Sheet for cPGCON 2017

Track ID: T1005743 Session No: 2 Date: 24/03/2017

Name of the Student: Kaushik S. Hande College Code: 005 Day (1/2): 1 - Friday

Title of Paper: Sentiment analysis using original and reversed Reviews

Name of Session Chair: Prof Dr. D.D. Patil, Prof Dr. Suoil Wankhede

Evaluation by Reviewer			
Sr. No	Title	Max. Marks	Marks by Reviewer
1	Relevance of Title, Abstract & Keyword	5	4
2	Adequate Literature Survey	10	7
3	Implementation Details a) Software Requirement Specification (05) 02 b) Mathematical Model & Design (05) 03 c) Implementation Status (10) 05	20	10
4	Algorithms-Measures & Metrics a) Performance Measures Used (05) 02 b) Result Tables (10) 04 c) Comparison with similar systems (05) 01 d) Efficiency Calculation (05) 02 e) Outcome & Success definition of work (05) 02	30	11
Total		65	32

Evaluation by Session Chair		Max. Marks	Marks by Session Chair
5	Concluding Remarks (Result Discussion, Conclusion & Future Scope)	5	03
6	References (Journal/Conference/Recent)	5	03
7	Format, Organization of Paper, Clarity & Linguistic Quality	5	04
8	Organization of Content & Presentation Skills	5	03
9	Question & Answers	5	03
10	Contribution of Authors: To society at large /Technology/Research(Publication or Patents) /Interdisciplinary approach/ Use of modern architecture & tools	10	08
Total		35	24
Grand Total		Out of 100	56

Remarks/ suggestions: _____

Internal Name & sign. <u>Prof. R. D. Pawar</u>	External Name & sign <u>Prof Dr. D.D. Patil</u>
--	---

Recommendation:

1. Recommended for Best paper: Yes/No 2. Recommended for Journal: Yes/No

3. Whether commercially feasible: Yes/No

Figure A.3: cPGCON Review Sheet

ANNEXURE B

DISSERTATION PLANNER

Table B.1: Dissertation Task Set

Task Title	Dissertation Task
T1	Study of Domain - Machine Learning and Natural Language Processing
T2	Identification of problem in existing systems
T3	Review of Literature
T4	Building Mathematical Model
T5	Report On Scheme of Implementation
T6	Identification of Prerequisites and Installation
T7	Configuring python and python package installer pip in the system
T8	Study of various machine learning algorithms and its implementation in python
T9	Studying libraries in python required for implementation
T10	Downloading and extracting reviews from IMDB movie datasets
T11	Removing stopwords, punctuation marks, numbers etc.
T12	Report Preparation
T13	Dissertation Project Stage I Presentation
T14	Document Preprocessing
T15	Creating Bag of words model from movie reviews.
T16	Splitting the dataset into training and test dataset
T17	Train machine learning classifiers using bag of words model.
T18	Create unigram, bigram, trigram variations of model
T19	Train machine learning classifiers using this model.
T20	Cpgcon Paper Presentation
T21	Predictive Model Construction
T22	Model Testing
T23	Experimental results, Analysis and Validation of results
T24	Project Review with Demonstration
T25	Report Validation and Submission, Report Submission