

Classification of Opinion on movie reviews by using classifiers with 3-gram TF-IDF and SVD Features

A Thesis

submitted in partial fulfillment of the requirements for the award of the degree of

Master of Engineering

in

Computer Science and Engineering Department

by

Shveta

(Roll No.: 801632046)



THAPAR INSTITUTE
OF ENGINEERING & TECHNOLOGY
(Deemed to be University)

Thapar Institute of Engineering and Technology

Patiala-147001

June 2018

Certificate

I hereby certify that the work which is being presented in the thesis entitled, " *Classification of Opinion on movie reviews by using classifiers with 3-gram TF-IDF and SVD Features* ", in partial fulfillment of the requirements for the award of degree of **Master of Engineering** and submitted in Computer Science and Engineering Department of Thapar Institute of Engineering and Technology, Patiala, is an authentic record of my own work carried out during the period under the supervision of **Assistant Prof. Ajay Kumar Loura** and refers other researcher's work which are duly listed in the reference section.

The matter presented in this thesis has not been submitted elsewhere for the award of any other degree or diploma from any institution.

Shveta

Shveta

Candidate

This is to certify that the above statement made by the candidate is correct and true to the best of my knowledge.

A Kumar

(Dr. Ajay Kumar Loura)

Assistant Professor

CSED

Thapar Institute of Engineering
and Technology
Patiala

Acknowledgement

I have been waiting long for this moment to acknowledge all those who contributed in building this work. It is my pleasure to thank all of them here. First of all, I offer my sincere gratitude to my supervisor, Dr. Ajay Kumar Loura, for accepting to be my supervisor. Without his help and encouraging advice, this work would have never begun. I am deeply indebted to him for providing me wonderful research atmosphere and platform to explore my research to the fullest.

I would also extend my gratitude to Dr. Maninder Singh, Head, CSED for providing me the opportunity to conduct my research work. I would also like to thank the Director of the institute, Prof. Prakash Gopalan for his continuous support.

My special thanks goes to my friends for discussing thoughts and sharing all ups and downs with me during the course of this work. At the same time I would also like to thank all my colleagues for their continuous support.

Last but not the least I would like to thank my parents and family members, who made me capable of reaching this point of life and for giving me their kind support and love. I dedicate my work to them.

Shveta
Shveta
ME(CSE)
801632046

Abstract

Extraction of features plays an effective role in sentiment analysis or opinion mining about an issue, customer reviews and products etc. in which these are fed to machine learning approaches to get the sentiments classified. Existing techniques widely used TF-IDF feature extraction from the unigram lexicons of the sentiment documents, some used the term frequency score of the unigram words as features. In this work, unigram, two word clusters (bigram) and three word clusters (trigram) are generated after filtering the collected sentiment data. Data used are the web movie reviews collection of the users. Data is collected manually from various websites(www.imdb.com, bookmyshow.com, google user reviews) about the conflict of Bollywood movie Padmaavat in which three different sentiments were found. Many people have positive moods about the issue and releasing of the movie Padmaavat and some of them are against of the movie, which are taken as negative reviews. A very little quantity was showing neutral moods, which show sentiments both in favor and against the movie. Hence three different categories of reviews are marked and fed to the proposed opinion mining system. All three unigram, bigram and trigram word lexicons are used further to get the TF-IDF of all the reviews. After that singular value decomposition (SVD) features are generated. Four different machine learning classifiers named as a K-Nearest Neighbor, Support Vector Machine, Naive-Bayes and Decision Tree are used for the classification step in which results are compared. Experimental results show more accuracy in classification when proposed feature extraction techniques are used as compared to existing method. Among the classifiers, decision trees give better accuracy in classification of sentiments than all other used classifiers. Decision tree gives 0.9272% accuracy in classification for positive sentiments, 0.8901% accuracy for negative sentiments and 0.9629% accuracy for neutral sentiments.

Keywords: Opinion, uni-gram, bi-gram, tri-gram, TF-IDF, SVD, Classification, Classifiers, Decision Tree, KNN, Naive-Bayes, SVM.

TABLE OF CONTENTS

TITLE	PAGE NUMBER
CERTIFICATE.....	(i)
ACKNOWLEDGEMENT	(ii)
ABSTRACT.....	(iii)
TABLE OF CONTENTS.....	1
1. INTRODUCTION.....	4
1.1 OPINION MINING	4
1.1.1 TERM PRESENCE VERSES TERM FREQUENCY.....	4
1.1.2 TERM POSITION.....	5
1.1.3 N-GRAM FEATURES.....	5
1.1.4 PARTS OF SPEECH.....	5
1.1.5 ADJECTIVES ONLY	5
1.2 TEXT CLASSIFICATION.....	5
1.2.1 TEXT CLASSIFICATION PROCESS.....	6
1.3 BASE CLASSIFIERS.....	7
1.3.1 NAIVE BAYES (NB).....	7
1.3.2 SUPPORT VECTOR MACHINE (SVM).....	8
1.3.3 K NEAREST NEIGHBOR.....	8
1.3.4 DECISION TREE.....	8
1.4 KNN AND TF-IDF FRAMEWORK.....	9
1.4.1 FRAMEWORK STRUCTURE.....	9
1.4.2 KNN CLASSIFIER.....	10
1.5 LEARNING PROCESS.....	10
1.6 DETERMINATION OF WEIGHT MATRIX.....	10

1.6.1 TERM BINARY.....	11
1.6.2 TERM FREQUENCY.....	11
1.6.3 TERM FREQUENCY INVERSE DOCUMENT FREQUENCY(TF-IDF).....	12
1.7 CONFLICT OF PADMAAVAT MOVIE.....	13
2. LITERATURE SURVEY.....	14
3. PROBLEM FORMULATION.....	21
3.1 PROBLEM DEFINITION.....	21
3.2 OBJECTIVES.....	22
3.3 METHODOLOGY.....	23
4. CURRENT WORK.....	24
4.1 OPINION MINING.....	24
4.1.1 PREPROCESSING.....	25
4.1.2 FEATURE EXTRACTION.....	26
4.1.3 OPINION CLASSIFICATION.....	27
4.1.3.1 MACHINE LEARNING APPROACH	28
4.2 PROPOSED SYSTEM MODULE.....	29
4.2.1 TEXT PRE PROCESSING.....	29
4.2.2 TRANSFORMING CASES.....	30
4.2.3 TOKENIZING.....	30
4.2.4 FILTERING STOPWORDS.....	30
4.2.5 FEATURE EXTRACTION BASED ON TF-IDF OF UNIGRAM, BIGRAM, AND TRIGRAM	
4.2.6 SINGULAR VALUE DECOMPOSITION FEATURES.....	30
4.3 CLASSIFICATION.....	31
4.3.1 DECISION TREE CLASSIFIER.....	31
4.3.2 SVM CLASSIFIER.....	32

4.3.3 NAIVE-BAYES.....	32
4.3.4 K-NEAREST NEIGHBOR.....	33
5. RESULTS AND DISCUSSIONS.....	34
5.1 HARDWARE AND SOFTWARE REQUIREMENTS	34
5.2.1 CLASSIFICATION MATRIX.....	36
5.2.1.1 CLASSIFICATION ACCURACY.....	37
5.2.1.2 SENSITIVITY.....	37
5.2.1.3 SPECIFICITY.....	37
5.3 CLASSIFICATION RESULTS OF MOVIE REVIEWS USING UNIGRAM, BIGRAM AND TRIGRAM AND SVD OF 3-GRAM TF-IDF'S.....	38
5.4 CLASSIFICATION RESULTS USING TF-IDF FEATURES OF UNIGRAM AND SVD OF TF-IDF UNIGRAM.....	44
6. CONCLUSION AND FUTURE SCOPE.....	50
6.1 CONCLUSION.....	50
6.2 FUTURE SCOPE.....	51

REFERENCES

CHAPTER 1

INTRODUCTION

1.1 Opinion Mining

Opinion Mining can be described as a sub-request of computational semantics that spotlights on expelling people's decision from web. Late advancement of web urges customers to contribute and impart by methods for web diaries, chronicles, long range casual correspondence goals, et cetera. Each one of these stages gives an unorganized measure of huge information that we are captivated to analyze. Highlight designing is a phenomenal degree fundamental and critical undertaking for Opinion Mining. Changing over a pinch of substance to a section vector is the critical stroll in any data driven way to deal with oversee Opinion. Some routinely utilized portions as a bit of Opinion Mining and their evaluations have been talked about in the running with sections.

1.1.1 Term Presence vs Term Frequency

Term Frequency [27] has continually been viewed as a key in standard Information Retrieval and Text Classification tries. In any case, it is found that term presence is more essential to Sentiment examination than term frequency. This is twofold regarded part vectors in which the segments simply display whether a term happens (value 1) or not (value 0). It has besides been seen that the occurrence of rare words contains a large information of data than irregularly happening words.

1.1.2 Term Position

Most words when showing up specifically positions in the substance pass on more conclusion or weight age than words showing up somewhere else. This looks like IR where words showing up in subject Titles, Subtitles or Abstracts and whatnot are given more weight age than those showing up in the body. In various cases, despite the way that the substance contains positive words all through, the closeness of a negative assumption toward the end sentence expect the picking part in choosing the evaluation. So most of the things considered words appearing in an initial couple of

sentences and last couple of sentences in a substance are given more weight age than those showing up elsewhere.

1.1.3 N-gram Features

N-grams [27] are prepared for getting association with some degree and are comprehensively used as a piece of Natural Language Processing aspect. Regardless of whether higher demand n-grams are significant includes wrangle about. It has been represented by examiners that unigrams defeat bigrams while gathering film reviews by estimation extremity, yet unique researchers found that in a couple of settings, bigrams and trigrams perform better.

1.1.4 Parts of Speech

Parts of Speech data is most regularly abused in all NLP errands. A champion among the most basic reasons is that they give an unpleasant kind of word sense disambiguation.

1.1.5 Adjectives only

Adjectives have been utilized constantly as segments among all parts of talk. A powerful relationship between enlightening words and subjectivity has been found. Though each one of the parts of talk are basic people, most typically used descriptors to outline a huge bit of the decisions and a high exactness have been represented by each one of the works concentrating on modifiers for feature time.

1.2 Text classification

For the most part the data store as content like messages, pages, daily paper article, statistical surveying reports, protest letter from client and inside created reports. With respect to an online daily papers give news under different classifications like national, global, governmental issues, back, sports, amusement and such messages orders are likewise an essential piece of content mining. Text classification constructed on pro knowledge how to categorize the manuscript under the mention set of groups. Data mining characterization begins with preparing a set of report that is as of now mark with class. Content grouping has two flavors as single mark and multi name. A solitary mark record is having a place with just a single class and multi name archive might have a place with in excess of one class. Information store among most

content databases is semi-structure information in which they are neither totally unstructured nor totally organized. For instance an archive may consist couple of organized fields, for example, title, creators, production date, and class yet additionally contain some to a great extent unstructured content parts, for example, unique, substance.

1.2.1 Text Classification Process

The phases of text classification are talking about as following focuses.

a) Document Collection

In Document Collection gather the distinctive sorts of a report like .html, .doc, .pdf and so forth.

b) Pre-Processing

In this present content report is a reasonable word design. For instance, Perform Tokenization, stemming word, evacuating stop words, tokenization a record is dealt with as a string and after that segment into a rundown of tokens. In evacuating stop word expel the stop words for instance "the", "an", "and". In stemming words changes various words from into comparative canonical form.

c) Indexing

It is to furnish the file to each report with this effortlessly distinguish as each record.

1) Feature selection

In the situation of preprocessing and ordering the essential advance of text classification defines feature selection. The essential idea of feature choice is to pick a subset of features from the primary report. It is performed by upholding the words with most astounding score as per foreordained measure of the significance of word. There are different sorts of highlights which can be utilized for content information i.e. term frequency, TF-IDF, unigram, bigram or trigram highlights of the words

2) Classification:

In this segment, records are grouped into predefined classifications. The records can be arranged by directed, unsupervised strategies. At the point when the class name of each archive is realized that is called regulated characterization when the class mark

of records isn't realized that is called unsupervised classification. There are different sorts of classifiers accessible i.e. Decision trees, Naive-Bayes, knn, SVM and so forth.

3) Performance Evaluation:

This is the end phase of text classification, this is tentatively, as opposed to systematically. In this measure the execution. Numerous measures have been utilized like accuracy and recall.

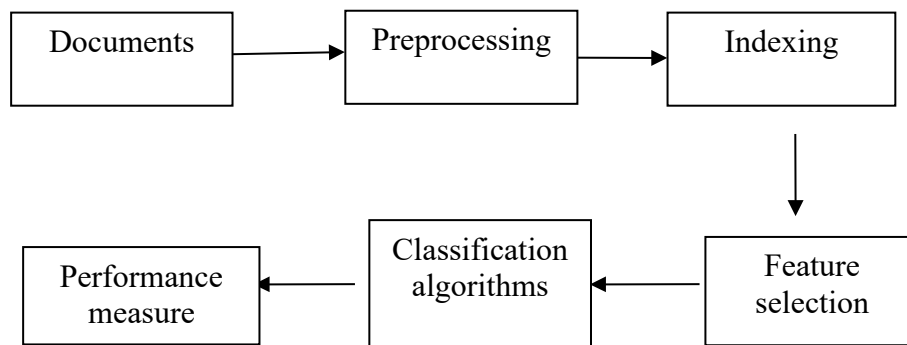


Figure 1.1: Document Classification Process

1.3 Base Classifiers

Base classifier is divided into following four categories.

1.3.1 Naive Bayes (NB)

The NB classifier [19] manages likelihood of a classifier in light of Bayes Theorem. Independent probability of features is explained by this classifier. The NB classifier creates an arrangement of classes i.e. collection of class name and likelihood for that class name. Bayes theorem for probability is characterized as beneath [19].

$$P(A/B) = P(A) P(B/A)/P(B) \dots(1.1)$$

$P(A/B)$ defines A's probability given that B is valid, $P(B/A)$ is probability of B given that A is valid (probability), $P(A)$ is earlier likelihood of the class and $P(B)$ defines the prior probability of indicator.

1.3.2 Support Vector Machine (SVM)

SVM [19] is a technique utilized for the order of both straight and nonlinear information. It gives high precision than rest arrangement systems.

As SVM develops hyper plane to create separation among two classes. It is exceptionally utilized for grouping systems. Maximum marginal hyperplane is accomplished by the greatest separation to the closest training data point of any class. MMH can be registered as takes after [19]

$$d(X^T) = \sum y_i a_i x_i x^T + b_o \dots\dots (1.2)$$

(Where Σ varies from $i=1$ to l)

where y_i defines the class label of support vector x_i ; x^T defines a test tuple; a_i and b_o are numeric parameters and l is a support vector's number.

1.3.3 K-Nearest Neighbor (KNN)

KNN classification methods [19] are instance based learning, utilized for pattern classification. KNN is a basic classifier that stores every accessible case and based on similarity classifies unlabeled new cases. Euclidean distance is used to calculate distance metric as given below [19].

$$\text{Dist}(X,Y) = \sqrt{\Sigma(X_i - Y_i)^2} \dots\dots (1.3)$$

(Where Σ varies from $i=1$ to n)

1.3.4 Decision Tree (DT)

Decision Tree is a machine learning model where nodes represent in tree like structure with various interior and exterior nodes, associated as branches. Interior nodes speak to a test on property and every branch means a result for every test. Root hub is the highest hub which speaks to properties; based on traits classification process starts. At long last leaf hub indicates a class name. There are two stages include in DT, i.e. pruning and developing. In developing stage, decision tree ends up being broad and in pruning stage, tremendous size tree diminishes its size. In order to register the estimation of DT, we need entropy and information pick up which depends on following figuring.

Entropy: $E(S) = \sum -p_i \log_2 p_i \dots \dots \dots (1.4)$

(Where Σ varies from $i=1$ to c)

Information Gain: $\text{Gain}(T, X) = \text{Entropy}(T) - \text{Entropy}(T, X) \dots \dots \dots (1.5)$

1.4 KNN and TF-IDF Framework

1.4.1 Framework structure

The structure comprises of a few discrete modules. Fig. 1.3 demonstrates the structure of the new idea of Framework. It has five fundamental modules: Documents module, GUI module, Preprocessing module, measuring module, TF-IDF module and KNN. GUI module empowers the customer to manage application and whole framework. Documents module is intended for record administration and choosing information assets. Record assets can be situated on the nearby PC or on the Internet.

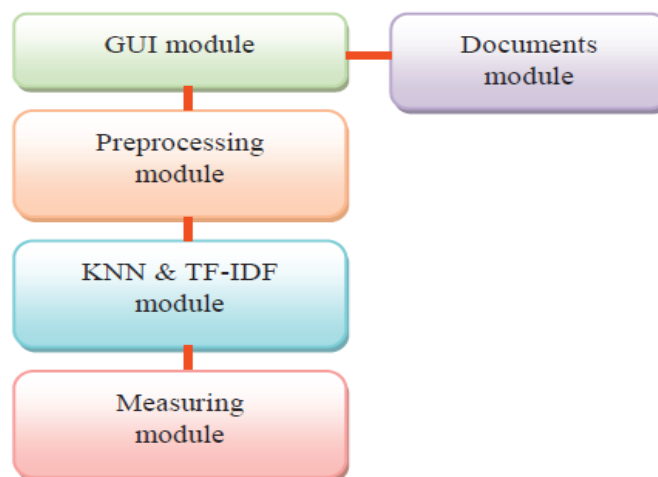


Figure 1.2 Framework structure [7]

In this module, a client can characterize report classes, the element which impacts the last after effects of content order. Preprocessing module examines the record attributes plans and changes them into an organization appropriate for characterization. The module alters reports to great content arrangement, paying little respect to the sort of textual style and character set. Reports consist of particular configuration module consequently expel control characters which may negatively affect the aftereffect of grouping. The fundamental module in structure is KNN and TF-IDF module. The

primary strategies are contained by module for characterization and assurance of record weight esteem. The consequences of characterization propagate to the endmost module for the estimation and introduction of measurable pointers. The endmost module permits showing grouping comes about and their straightforward measurable investigation.

1.4.2 KNN classifier

As effectively determined the calculations rely upon the machine learning. Report availability and Preprocessing is trailed by learning stage. The calculation decides the essential reports that will contrast and each new record. Computation checks where a report is organized by simply looking arrangement records that are most similar to it. The calculations consider that it is conceivable to classify documents in the Euclidean distance as focuses. Euclidean separation is the separation among two focuses in Euclidean space. The separation among two focuses in the plane with organizes $p=(x, y)$ and $q=(a, b)$ can be calculated [24].

$$d(p, q) = d(q, p) = \sqrt{(x-a)^2 + (y-b)^2} \dots\dots\dots(1.6)$$

1.5 Learning process

Learning process begins with parsing the fundamental content which looks through the words in records and structures a vector [20]. Parsing process evacuates all control characters, spaces between dabs, words, commas, and comparative characters. The framed vector speaks to a key protest that will be utilized for grouping of tried archives.

1.6 Determination of the weight matrix

To give text classification and looking through records it is important to set up the weight grid. The framework contains the estimations of relations between every exceptional words and archive. It is an underlying article in the calculation to figure individual significance (weight) of each looked report. Each report is spoken to as a vector in n-dimensional vector space. Imagine a framework A with measurements NxM, where N measurement is characterized by various one of kind words in an

example all things considered. M speaks to the number of records to be arranged. Weight network can be described as a social grid of word - archive. Measurement of the grid is equivalent to item, the quantity of various one of a kind words and the aggregate number of reports. Every framework component a_{ij} speaks to weight estimation of word I in record j . Weight network is shown in Fig1.4.

j - number of all documents

i - number of all unique words	$a_{(0,0)}$	$a_{(1,0)}$	$a_{(2,0)}$	$a_{(3,0)}$
		\						
			\					
				\				
					\			
						\		
							\	
								\
	$a_{(i,0)}$							

Figure 1.3 Weight matrix [7]

In deciding the weight values in the lattice, diverse measurements and techniques for figuring can be utilized.

1.7 Determination of weight matrix

a) Term binary

Binary strategy is to a great degree straightforward and simple to actualize. The technique checks if a specific word (term) shows up in archive. The qualities in grid can be 0 or 1. On the off chance that the word shows up in the record, at that point the weight esteem a_{ij} is set to 1, generally is set to 0.

b) Term frequency

Term Frequency has been used to fetch out information and it has been revealed how frequently an expression has occurred in the document j [14].

$$a_{ij} = f_{ij} = \text{frequency of term } i \text{ in document } j \dots\dots(1.7)$$

Some text classification algorithms performed standardize term recurrence by isolating with the recurrence of the most widely recognized term in the documentation.

$$a_{ij} = f_{ij} / \max_i \{f_{ij}\} \dots\dots\dots(1.8)$$

c) Term frequency – inverse document frequency (TF-IDF)

An extremely well-known research strategy in the field of natural language processing (NLP) which is utilized as a part of the usage of the calculation defined in this article. TF-IDF technique decides the relative frequency of words in a particular report through an inverse proportion of the word over the whole archive corpus. In deciding the esteem, the technique utilizes two components: TF - term recurrence of term I in archive j and IDF - inverse document frequency of term I. In our examination and testing the calculation of system this technique indicated great outcomes. TF-IDF can be computed as [15]:

$$a_{ij} = tf_{ij}idf_i = tf_{ij} \times \log_2 \left(\frac{N}{df_i} \right) \dots\dots\dots(1.9)$$

Where N is quantity of reports in gathering, a_{ij} is heaviness of term I in document j, tf_{ij} is the term recurrence of term I in document j and df_i is the document frequency of term I in the accumulation. Keeping in mind the end goal to acquire better outcomes with records of various length, an adjusted condition has been utilized as indicated [18]:

$$a_{ij} = tf_{ij}idf_i = \frac{f_{ij}}{\sqrt{\sum_{s=1}^N (tfidf(a_{sj}))}} \times \log_2 \left(\frac{N}{df_i} \right) \dots\dots\dots(1.10)$$

1.8 Conflict of Padmaavat Movie

In this work, Padmaavat film has been chosen for opinion mining. The genuine story behind the challenges by the Rajput Karni Sena started from conflict of expanded inflated egos. In an elite Facebook Live with Firstpost, the man who is driving the rally against the arrival of Padmaavat, benefactor of Karni Sena Lokendra Singh Kalvi said everything started with a somewhat inconsequential proclamation made by

Singh amid a communication with the media when shooting for the motion picture was going to start. As indicated by Kalvi, Singh, reacting to an inquiry whether he was assuming the part of a reprobate in the movie, said that he could go two indents past assuming the part of a scoundrel on the off chance that he is given two personal scenes with the lead female character in the film. Somebody sent that cut-out to Karni Sena in Rajasthan. Bhansali, in an announcement issued over the challenges, stated, "This film involved into such a significant number of debates due to some gossip. Gossip is this, in the movie a fantasy scene has been recorded between Rani Padmaavati and Alauddin Khilji. I have officially dismissed this claim and furthermore given a composed evidence of this. Today once more, I am repeating through this video that in our film there is no such scene between Rani Padmaavati and Alauddin Khilji, which would hurt the sentiments of anybody."

CHAPTER 2

LITERATURE SURVEY

Dong et al. [32] introduced comes about because of 6 diverse product domains it could be contended that, since these areas all identify with shopper hardware, the generalizability of our approach might be overstated; maybe it functions admirably for customer gadgets with specialized highlights however not maybe for different domains. At the season of composing they have finished a comparative assessment for inns utilizing Trip Advisor inn surveys. They likewise exhibited and display a fundamentally the same as picture both as far as the sort of inn cases that are created and our capacity to deliver valuable suggestions by consolidating comparability and opinion. On this premise, they can be to some degree certain that the methodologies they have portrayed in this work give a helpful new way to deal with case-based item proposal.

Brindha et al. [12] proposed a word weighting plan in light of inference through statistical functions. The word weighting mixes are connected to three prominently known and generally utilized audit informational indexes. The outcomes demonstrate that the proposed plans, particularly SWD(3)*IFO, beat the broadly utilized audit informational indexes in view of the gathered weighting and furthermore create the best precision of 97.4% (most elevated in the Cornell movie surveys). The proposed techniques are meant to uncover that the derived weighting plans are superior to the plans that don't think about the relationship between words and their communicated extremity. Notwithstanding the change of grouping exactness, our examination uncovers that the generally utilized stop-word expulsion process is not essential for feeling mining classification.

Kim et al. [9] exhibit a strategy to anticipate the movies execution of a film in light of client remarks made on its trailer and the related advertising properties accessible on informal communities and somewhere else on the Internet. To appraise client conclusion, word records were characterized in light of well-known feeling words, emojis, shortened forms, and casual words. They likewise gathered advertising properties including data on the main stars, chief, author, and their past works. The

movies record forecast framework was prepared utilizing data on discharged films as the ground truth. Through the investigation, they demonstrated that considering client remarks and showcasing properties together prompts better expectation exactness.

Armentano et al. [23] assessed distinctive tokenization methodologies, pre-handling techniques and calculations to assemble classification models that can decide the sentiment and extremity communicated in the short messages distributed in Twitter. They at that point utilized the best classifier to separate the primary opinion of Twitter clients in regards to an arrangement of target movies with a specific end goal to assist clients with deciding to see the motion picture or not. With respect to wanted arrangement for an order show for supposition mining, they found that the models in light of SVM acquire preferable execution over those in view of Bayesian Networks and Decision Trees. Likewise, the best models are those that utilize an element determination strategy, (for example, Information Gain).

Manek et al. [3] proposed a measurable strategy utilizing weight by Gini Index technique for feature selection in sentiment investigation while in the meantime enhancing the exactness of sentiment polarity prediction utilizing different vast movies informational collection. Our proposed system for slant examination utilizing the SVM classifier is contrasted and other component determination techniques on film audits and results have demonstrated that order by utilizing this proficient and novel strategy has enhanced the exactness.

Hua et al. [39] proposed a novel text summarization procedure to decide the best k most instructive sentences from online lodging audits. Most past investigations on survey rundowns have concentrated on text preprocessing, which ignores basic factors, for example, the validity of audit creators, survey time, audit value, and clashing sentiments. In addition to the data (i.e., catchphrases and key expressions) extricated from the surveys, this investigation additionally thought to be all the previously mentioned factors. Specifically, they utilized the survey creators' recorded information to ascertain their belief, the audit time to decide the estimation of a lodging survey, and the number of proposals to characterize the convenience of an inn audit. Sentence positions, sentence length, regardless of whether a sentence contained

pointer expressions and the greater part of the previously mentioned factors were then consolidated to acquire the sentence significance score.

Taylor et al. [11] executed a framework and handled the issues in the Lake District tourism industry, in the south of Chile. The input given by the framework clients demonstrated that our outline and perception diagrams, which were likewise proposed as a piece of our augmentation, are straightforward and give real bits of knowledge about feeling, demonstrating how helpful and great our device is. Their outline and models for angle based opinion can be utilized as a part of numerous conceivable applications in the tourism space. Advantages that may emerge involve the two travellers and specialist co-ops.

Kang et al. [22] proposed another sentiment analysis technique utilizing an outfit of content based concealed Markov models, the EnsembleTextHMM strategy. Rather than depending on removed sentiment dictionaries or predefined catchphrases, it utilizes marked preparing writings to reflect assorted examples of assumptions. For that reason, they assembled a classifier utilizing a content-based concealed Markov display through grouping the words, and they connected an outfit of such classifiers by changing the number of covered states and examining the training texts.

Yang and Lin [17] proposed new methodologies for examination at the sentence level and the documentation level. The fundamental commitment of this examination is that they proposed another learning-based opinion mining structure that thinks about numerous kinds in workmanship item/benefit surveys. They directed tests where they looked at the proposed approach (numerous compose examination) and the customary approach (single summed up type investigation). Their approach performed superior to the traditional approach. All the more significantly, they demonstrated that the weight table is reusable and it is just important to prepare another weight table for another sort. It can be expected that the positive/negative ramifications of different kinds of feeling (e.g., cheerful and irate) in the audit messages on a blog/discussion for workmanship items/administrations will stay steady for quite a while. This reusability could diminish the time required to mine information without performing the whole procedure once more.

Domeniconi et al. (2016) [13] formulated a base tf.idfec plot where idf is registered without considering reports having a place with the demonstrated classification: this anticipates giving a low weight to terms to a great extent show in it. The tf.idfec-based variation blends our base plan with the significance recurrence from tf.rf, along these lines likewise successfully boosting weights of terms which show up as often as possible in the classification under examination.

Ibrahim and Silva (2015) [28] infer that the proposed Term Frequency - Average Term Occurrences (TF-ATO) term-weighting plan (TWS) can be viewed as focused when contrasted with the broadly utilized TF-IDF. The proposed TWS gives higher viability in the two instances of static and dynamic report accumulations. Likewise, the record centroid vector can go about as a limit in standardization to separate between archives for better viability in recovering applicable reports. They watched a variety and diminishment in framework adequacy when utilizing dynamic rather than static record accumulations, in addition to there is an extra cost for each refresh to the gathering.

Bansal and Srivastava (2018) [5] utilize word2vec demonstrations to arrange more than 400,000 online purchaser audits for different worldwide cell phone brands gained from Amazon. They first discovered includes most items like angles by word2vec and demonstrate that word2vec can discover semantically comparable words. At that point, they utilize CBOW and skip-gram strategies with four distinctive order calculations: Naïve Bayes, SVM, Logistic relapse and Random Forest. Results demonstrate that CBOW performs well when contrasted with skip-gram, showing that information may comprise much of the time happening same words. Arbitrary Forest beats every one of the calculations when utilized with word2vec portrayals. Consequently, conveyed word vector portrayals can be productively utilized for the errand of feeling arrangement by fusing semantic word relations and contextual data.

Trstenjak et al. (2014) [7] show a structure for text classification in light of KNN algo and TF-IDF strategy. Principle inspiration for examination was to create idea of structures with accentuation on KNN and TF-IDF module. System with installed techniques gave great outcomes, affirmed our idea and starting desires. Assessment of system was performed on a few classifications of archives in online condition. Tests

should give replies about the nature of characterization and to figure out which factors affect execution of order. The system work was exceptionally steady and dependable. Amid testing the nature of grouping they have accomplished great outcomes paying little heed to the K consider esteem the KNN algorithm.

Kukkar and Mohana (2018) [4] utilized the hybrid approach of combining text mining, natural language processing and machine learning methods to distinguish which bug report as non-bug or bug. Just because the clamors of misclassification is lessened by separating the bug reports and improve the execution of programmed bug forecast. In this work the four join fields (journalist, rigorousness, precedence and constituent) with textual data like remarks, depiction; outline and so forth are added to testing and preparing dataset. The bigram and TF-IDF approach is utilized with Info pick up for the bug seriousness expectation. The bigram approach helped in lessening the sparest of dataset. To compute the exactness of proposed display, Precision, Recall and F-measure are utilized. It is watched that, the execution of KNN classifier is changed by the dataset.

Du et al. (2018) [40] propose the various leveled grouping classification approach in light of the relaxation methodology which eases the effect of the 'blocking' issue. It defers the indeterminate classification decision until the point when it can be grouped certainly, thus the mistake that has happened in the upper level will not be exchanged to the lower level. They additionally apply the Least Information Theory in term weighting and documentation drawing and it offers another fundamental data evaluate demonstrate by various likelihood circulations.

Akimushkin et al. (2018) [8] presented an approach by which the demonstration of text with complex systems is improved by considering the words relating to the nodes. This is finished with a comparability metric to think about two bits of content where the nearness of the most important words, as per arrange measurements, is considered. At the point when the information acquired with the comparability measurements were utilized as contribution to machine learning calculations, a high exactness was accomplished which achieved 98.75% for one of the book accumulations. Essentially, the exactness was impressively higher than for conventional techniques in view of tf-idf, being additionally higher than other system approaches that did not think about

the mark of the hubs. Additionally significant is that the execution was enhanced with dimensionality decrease with MDS, which is favorably inferable from the lower computational cost.

Kwon et al. (2018) [29] assess the pattern in deqi research and the related components to create deqi, they got the accompanying conclusions and feasible arrangements. To start with, the present examination set excessively accentuation on mediations identified with deqi, and future investigations should concentrate more on instruments. Second, muscle, force, profundity and weight were identified with the age of deqi. More examinations about neuronal instruments are prescribed to take in more about the generation and utilization of deqi. In conclusion, there is presently an opportunity to propel the technique of deqi look into and to examine the instrument of needle therapy treatment by exploiting the advance of related innovation.

Niemann et al. (2017) [16] presented patent paths, which they characterize as the arrangement of patent clusters over the span of time. They have built up a non-specific process comprising of five stages, and talked about six outline decisions identified with these means. The field of carbon filaments in a blend with bike innovation was shown the patent path strategy, prompting bits of knowledge into the advancement of the mechanical field after some time.

Mahmud et al. (2017) [31] built up a framework that can predict movie success in light of viewer sentiment analysis. Utilizing our sentiment classifier they accomplished stable outcomes in both of the prevalent datasets for film audits so they can state that our classifier is genuinely steady and it is not inclined to a particular dataset. The framework has its own shortcomings. It performs extremely well on films that have either a high or low evaluating yet the execution diminishes when they consider the hazy areas are meaning when they consider motion pictures that have a normal rating.

Liu et al. (2017) [33] built up a model in which the substance of the tag group got in each run are not precisely the same in the little piece of the labels, the circumstances that every news is labeled is changed, however the general change is little, the general consequence of labeling news isn't bigger influenced. Especially, the labels separated

from this model are construct in light of the news content, as well as on the qualities of the client's labeling conduct in the social labeling framework on the Internet, and supplement the label aggregate in view of the news content by breaking down countless remarks on the system of applicable news points. They completely think about the conduct of Internet clients and concerns and system social popular assessment, to accomplish the high estimation of news labels. Furthermore, through the investigation of the substance and recurrence everything being equal, this model likewise gives a technique for extricating the focal point of the news theme content.

CHAPTER 3

PROBLEM FORMULATION

3.1 Problem definition

The Sentiment is thought, attitude, judgment prompted by feelings. Sentiment analysis is also known as Opinion mining, studies people's opinions towards certain entities. Large amount of information is available on the internet on different products. These products can be mobile phones, Cameras, Movies Etc. All Information is available on net regarding these products. This information helps the customer to select the relevant product by comparing their reviews on internet. After getting the product, the customers give the reviews regarding the product on the internet. This will be helpful to other customers to get information about their products. Now a day's Sentiment Analysis become a big issue. Explosion of social media has provided many opportunities for citizens to publicly spread their opinions, but it has created serious problems when it comes to make sense of these opinions and because of viral nature of social media importance to get an understanding of citizens opinions has grown (when attention is very unevenly and frequently distributed) Some issues become important through word-of-mouth.

In many of such situations where investigation could assume a noteworthy part in its advancement. With the progression in the technological field and having web access to numerous individuals are giving their opinions in different web destinations and online journals. Client surveys are essential for different fields (e.g. Films, Products, and Services). Film surveys assume a crucial part in portraying its prosperity and disappointment. Individuals have now turned out to be unmistakable on what films to watch and what not to watch. Thus individuals would prefer not to sit idle on a film that has awful audits. Nowadays online review is important for a personal recommendation. The major part of aggregate the opinion of people is the information gathering process. Sentiment is an opinion, feeling, view that is expressed in the form of negative or positive. Classify and summarize the opinions expressed online regarding services, products. The opinion mining has numerous challenges. The foremost is that opinion word is not always considered positive or negative, in one

condition it may be positive and in another condition it may be negative. Second challenge is opinion can be in the form of a compound sentence or a simple sentence. Moreover, the compound sentence is more challenging. There isn't sufficient work done in opinion mining of compound sentences. Our investigators concentrated on movie audits. There are expansive measures of client created films survey are accessible on the web like IMDB, BookMyShow, google user reviews and so forth. There are numerous difficulties like at least one terrible feature of film does not make it overall bad similarly as one or more good features does not make it good overall. Therefore, opinion mining of film review is viewed as more difficult than opinion mining of other type of reviews.

In our work, Padmaavat movie reviews are collected in order to explore opinion mining. Reviews have been collected from online websites like www.imdb.com, bookmyshow.com and google users reviews.

Padmaavat, in view of a lyric composed by Sufi artist Malik Muhammad Jayasi in 1540, portrays how Rajput ruler Maharawal Ratan Singh and his wonderful and politically keen ruler Padmaavati endeavor to spare their hundreds of years old kingdom from a savage madman named Alauddin Khilji, the eager Sultan of Delhi set on having Padmaavati and overwhelming India. At the point when Khilji triumphs, Padmaavati picks self-immolation over oppression- a demonstration considered brave in now is the ideal time.

In the releasing dates of the movie a conflict arises and people were divided into two groups in which one type of people were in favor to release it and another were against to release. There were some neutral reviewers as well. Hence People have positive as well as negative views on this movie.

3.2 Objectives

Based on literature survey, following objectives are outlined:

- 1) To study various feature extraction algorithms i.e. TF-IDF(term frequency-inverse document frequency), unigram, bigram, tri-gram, SVD etc.
- 2) To collect the database of Padmaavat movie. Reviews are collected manually from various sources (www.imdb.com, bookmyshow.com, google user reviews) of internet.

- 3) To preprocess the database i.e. exclusion of special characters, numeric, article words like a, an, the etc.
- 4) To classify the data using machine learning Algorithms namely Decision trees, k-nearest neighbor, Naive Bayes and Support Vector Machine. To do training and testing of data with the classifiers of these algorithm.
- 5) Performance evaluation by comparing the results of different classification Algorithms and to select the most effective one in order to optimize the result.

3.3 Methodology

The steps carried out for the proposed framework are:

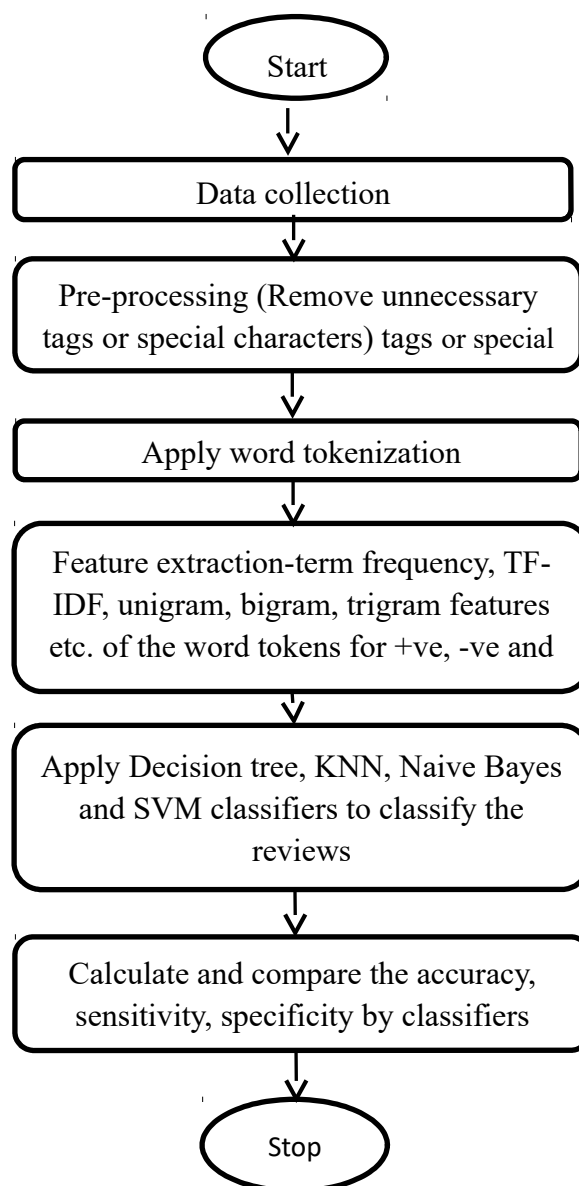


Figure 3.1: Steps of Methodology

CHAPTER 4

CURRENT WORK

4.1 Opinion Mining

Sentiment Analysis (SA) or Opinion Mining (OM) is a procedure to examine the discussions about a subject or an item, considering the framework which mechanizes this procedure. Among the undertakings to opinion mining, it can be said the influence investigation, subjectivity examination [35], feeling investigation, and the relevant extremity (negative or positive) about a record or remark. Here, there will be concentrated continuously to reviews of opinion mining wherever the extremity of the supposition about an item's element need to be analyzed. The particular opinion can

have five fundamental segments, i.e. $(o_j, f_{jk}, so_{ijkl}, h_i, t_l)$ where:

- 1) o_j is a main object for the opinion to be expressed. It o_j can be an event, a person, organization, product, or topic
- 2) f_{jk} is a feature of the object o_j ,
- 3) so_{ijkl} is the sentiment value of the opinion of the opinion holder h_i on a feature f_{jk} of object o_j at a time t_l ,
- 4) h_i is an opinion holder or source of the opinion,
- 5) t_l is the time when the opinion expressed.

Opinion mining is an arena of content mining [6]. The motivation would be to order remark to be negative or positive sentiment. In this way, opinion mining can be anticipated to a parallel content order issue while mulling over a few attributes of OM issues. It can, then, have the same building pieces of a content order framework.

The procedure begins with information accumulation. Numerous sources like websites, online networking and web reviews contain items opinions. When all is said in done a remark or a conclusion must be re-handled to outline content remark to a representation reasonable for the programmed order, then element ID and extraction takes after, at long last, remark's extremity grouping is executed as appeared in Fig 4.1

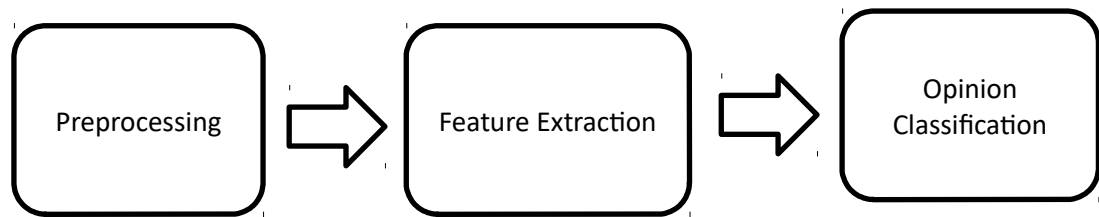


Fig. 4.1 Opinion mining process

Figure 4.1 shows the process of opinion mining, in this preprocessing is initiated and further it gets to the Feature Extraction and finalizing the whole process by getting to opinion classification.

4.1.1. Preprocessing

The content records is a proper information structure which comprise an assessment should be preprocessed and put away in a proper information structures for further handling. For maximum part, these assessments comprise rare syntactic components that might not be valuable for the following strides. Tokenize these suppositions, standardized after that cleaned. Some propelled preparing may be done on content feelings, to give some examples, standardization, gathering of equivalent words and spelling errors checking.

Normalization is acknowledged through outcome and expelling word additions that are associated with enunciation, to distinguish and consider distinctive events of the same term freely of their part with regards to the particular sentences in which they are utilized this could be vital undertaking if machine learning methodology is utilized later.

a) Grouping of equivalent words: The issue experienced is to communicate the same significance of a word by an extensive range of words, e.g. 'profit, "income" and "profit" have same significance. By considering the issue the words must be gathered to encourage highlight extraction. The works that bargain this issue it has been finding [42] that use a semi-directed strategy in view of Word Net considers the conclusion to aggregate these equivalent words.

b) Spelling correctors and check: There are numerous methodologies in addition to strategies for spelling checking which could be arranged into connection free and setting subordinate blunder redresses. The methodologies for connection autonomous

rectification execute words amendment freely of the setting by utilizing probabilistic systems or neuronal system ones. The methodologies for connection subordinate blunder remedy complete a revision as indicated by the appropriate data accessible utilizing semantic separation ones or machine learning.

4.1.2. Feature Extraction

It is determined by appliance domain, as options of “image” in image process are: intensity, distinction, luminosity. Though, merchandise opinion mining characteristics (e.g. reviews features) are: Keywords, or vocabulary. This could be a crucial stage in product opinion mining. This can be a crucial step in product opinion mining that may be categorized into four groups: lexicon-based, machine learning, dependency-relation-based and ontology approaches.

a) **Machine learning feature extraction approach:** Creation options area unit usually noun phrases or nouns. Machine learning method principally depends to chances on those noun phrases and nouns that supported their incidence frequencies. There is a supervised learning technique supported lexicalised Hidden Markov Models (L-HMMs) supported linguistic options [25].

b) **Ontology feature extraction approach:** Metaphysics is employed to abstract the options enclosed within opinion. There is domain metaphysics to spot the options of opinions conveyed by users; however with the fast upsurge as well as type of on-line merchandise, there is an affinity to can’t continually notice

1) **Lexicon primarily based feature extraction approach:** The product options supported a lexicon are extracted. A lexicon has been made that contained an inventory of opinion words and feature words and analyzed it to assign a polarity tag to all features. At that time informatics methods and apply math ways area unit to extract opinion words and feature words supported the outlined lexicon.

2) **Dependency relation primarily based feature extraction approaches:** In this, examine term dependency associations in sentences, and after that put on certain algorithms and rules to extract product options from the known dependency relations [41].

4.1.3. Opinion classification (OC)

OC worries with distinctive alignment of sentiments and categorizes to either negative or positive to spot this alignment, a machine learning approach, lexicon primarily based approach or a hybrid approach is used.

a) Machine Learning Approach: It is categorized in unsupervised and supervised learning strategies. Supervised learning systems require labeled coaching documents. On the other side, unsupervised strategies square measure used once it's tough to seek out these labeled coaching documents. The foremost used classifiers are Artificial Neural Network (ANN), Support Vector Machine (SVM) etc.

b) Lexicon-based Approach: A set of noted and precompiled sentiment terms determine this approach which associates sentimentality coordination and words is a sentiment lexicon. It differs considerably; when specific sentimentality wordlists outline sentimentality marks along variable arithmetic series whereas other outline single or additional sentimentality classes appreciate negative, neutral, positive, or a spread of emotions appreciate happiness and unhappiness.

In the existing work, suggestive square measure studied at the word level. Thus, the matter is expressed in sight of words. Task (1) given a Reviews paragraph, P , predicts its category, c , wherever c could be a regular opinion, comparative or suggestive. As stated previous, a feature is indicated via specific keywords or a lot of complicated phrases and syntax patterns.

1. Bag-of-Words options

As each speech in its raw kind is considered equally an order of words, training algorithms cannot use them; because of its peak expected input within sort of real-valued feature vectors within a set size instead of raw text documents of variable length.

2. Suggestive Clues

Provide the tokenized contented annotated with part-of-speech tags, generate options which provide the occurrence of patterns. The restricted and tokenized text is combined and traversed with its part-of-speech tag sequence to describe that a pattern matches or not.

3. Comparison Clues

Comparators clues are directly combined with consecutive category rules to the task of detect comparative sentences. These options are evaluated because comparative sentences area use a disjointed set from suggestive sentences.

4. Performative Verbs

In the direction of requirement under consideration the training scheme on however recommendations area unit given, additionally enclosed a lexicon of per formative verbs. We expect that together with an indicator of a specific statement comprising a per determinative verb is also indicative that the presenter is determining whether or not a product deserves a negative or positive plan. This feature set give out one feature, with its worth what number per formative verbs were noticed within the sample.

5. Sentiment Lexicon options

In sentiment lexicon options, external information is integrated within the current models. Specially, we have a tendency to use SenticNet, a conceptual object that maps a sentiment worth from -1.0 to 1.0 to a group of words or multiword terms.

6. Mixed Sequential Rules

A series of items is symbolized as $S = \langle i_1, i_2, \dots, i_n \rangle$, in which each item represent certain token. Latest approaches for contrast grouping have productively analyzed successive forms of part-of-speech tags adjoining definite keywords [25]. The researchers planned a sequence construction methodology that considers every sequence among 3-token radius from keywords and interchanges all words excluding the keyword with numerous part-of-speech tags. To higher make a case for the development, a way is hereby planned that generates many options from a similar symbolic series.

A diverse sequent instruction feature generator is outlined into 2 portions: word of consideration and surrounds. Abs initio all token series square measure infatuated sizes up to n that comprise a minimum of one word of concern. After that the feature is produced by uniting the exact illustration of word of attention, therefore the close tokens square measure portrayed by the several part-of-speech tags and the subsequent rules outline numerous categories of words of attention for producing sequent rule options.

4.2 Proposed system module

The three steps are considered in this method. These are:

4.2.1 Text Preprocessing

The first phase is the text preprocessing in the process of reviews classification. The unstructured data are received, at that time, text is effectively preprocessed. Collect the data from diverse sources, and then need to be cleaned. Firstly, clean the text data from useless and distortion information such as exclamations, irrelevant sentences, punctuation, dates, semicolon, quotations, etc.

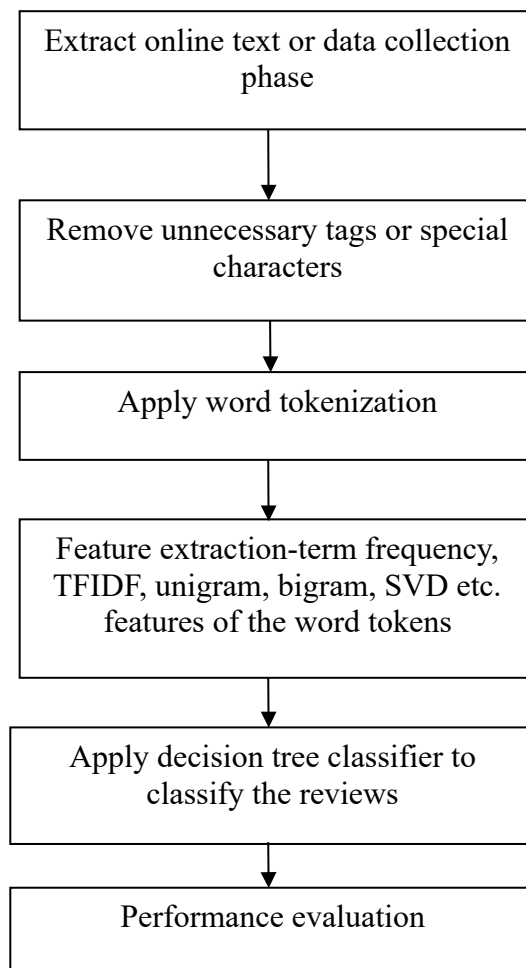


Figure 4.2: Flowchart for the Reviews classification process.

Figure 4.2 represents the complete steps in opinion classification on movie reviews with help of different machine learning classifiers with uni-bi-tri gram TF-IDF and SVD of uni-bi-tri gram TF-IDF features.

4.2.2 Transforming Cases

Transform all the (lower case as well as upper case) characters into lower case characters with the help of transform case operator. It is basically used to eradicate similar words that are different only in their case.

4.2.3 Tokenizing

The research accompanied on text mining contains sentences or words that must be separated word by word to increase the processing. So, separate all the words in sentences, and all the punctuations are disposed of since they cannot characterize any collection. It makes simpler calculations in the subsequent stages [1].

4.2.4 Filtering Stopwords

One of the techniques used meanwhile the principal studies were led on information retrieval is filtering stopwords. This system deletes needless things like words occurring too often or too rarely in sentences or text documents. This is also used to remove insignificant words or the words without any definite meanings such as a, an, or the. Similar to preceding stages, this phase is used to lessen computational complexity or processing time [1].

4.2.5 Feature Extraction based on TF-IDF of unigram, bigram and trigram

After generating the wordlist, the weight of every word is considered with respect to TF-IDF that is one of the best processes applied in text mining research. The word frequency means how many times a term is recurrent in text and IDF stands for Inverse Document Frequency, a procedure used to calculate the inverse probability of finding a word in text [1]. Equation 4.1 is a classic TF-IDF equation used to calculate weight:

$$W_{ij} = tf_{ij} * \log \frac{N}{df_i} \dots\dots\dots(4.1)$$

In this equation, N is the number of documents in the set of total documents, w_{ij} is the weight of weight of the word i in the document j, df_i is the number of documents containing the word I, tf_{ij} is the frequency of the word i in the document j.

4.2.6 Singular Value Decomposition (SVD) features of TF-IDF for bigram, unigram and trigram features

Calculation with SVD is done after TF-IDF Vector obtained after TF-IDF calculation will be used to calculate SVD. Given a term-document matrix $A=[a_1, a_2, \dots, a_n]$, SVD can be calculated using:

$$A_{mn} = U_{mn} \times S_{mn} \times V_{mn}^T \dots\dots\dots(4.2)$$

The result from the feature selection will be processed with the classifier.

4.3 Classification

4.3.1 Decision Tree classifier

A decision tree [2] is characterized as a tree and it is a classifier utilized for text classification where every hub can go about as a leaf or decision hub. In therapeutic science [30], decisions are to be taken quickly and minor postponement may prompt significant issues. Performing calculated examination [2] and basic leadership demonstrates suitable for such circumstances and framework may work better than anyone might have expected. For controlling such circumstances, decision trees have helped scientists a considerable measure where therapeutic specialists can improve examinations on premise of results they acquire and guidelines can be deducted on their premise. Decision trees are straightforward and standards can be effectively produced through them. They can take care of complex issues effortlessly. In any case, preparing through decision trees is extremely costly and is expensive. In addition, an audits archive class must be associated with one branch as it were. A solitary error in higher upper level can cause entire subtree invalid. There are issues of constant factors and over fitting in decision trees [37]. The fundamental reason for the decision tree calculations is to part the element space into one of a kind locales relating to the classes [36]. An obscure component vector is relegated to a class through a grouping of Yes/No decisions along a way of hubs of a decision tree. C4.5 is a calculation used to produce a decision tree and it is known as one of the fruitful decision tree grouping calculations.

4.3.2 SVM Classifier

It is support vector machine. It is a binary classifier. It groups the contribution to paired frame. At the point when input is sure then it speaks to as 1 else it speak to as 0. It depends on a negative and positive dataset which is extraordinary in other characterization. SVM is a supervised learning model. It is utilized for investigation the information. The SVM takes an arrangement of information to performing straight order. SVM can productively perform nonlinear characterization utilizing portion traps. A fundamental property of SVM is minimization the arrangement mistake and augments the edge. So it is additionally called most extreme edge classifier. There are different uses of SVM like it is useful in content and hypertext arrangement. Order of pictures can likewise perform utilizing SVM. A Manually written character can be perceived by utilizing SVM. SVM is additionally utilized as a part of medicinal science to arrange diverse sicknesses.

4.3.3 Naïve-Bayes

Thomas Bayes introduced this classifier method. This method learning from data and predict class in which each class have probability [10]. Bayes theorem is shown in Eq. (4.3)

Naive Bayes Classifier: Naive Bayes is a classifier technique which is presented by Thomas Bayes. This technique which is learned from data and predict class which each class have a probability [10]. Bayes hypothesis is appears in Eq. (4.3)

$$P(A/B) = \frac{P(B/A) * P(A)}{P(B)} \dots\dots\dots(4.3)$$

where $P(A)$ and $P(B)$ are probabilities of observing A and B . $P(B/A)$ is the probability of observing event B given that A is true. Naïve Bayes equation is represented by $P(A|B)$, A is a input vector that has feature and B is a class label. Based on information from training data, for each combination A and B, the final probability $P(B|A)$ of model should be trained. With that model, testing data of A'

can be declared by look for B' value by maximazing P(A'|B') value. Then for classification, Naïve Bayes formula can be declared as Eq. (4.4)

$$P(B / A) = \frac{P(Y) \prod_{i=1}^q P(A_i | B)}{P(B)} \dots(4.4)$$

where P(B| A) is probability data for A vector in Y class. P(Y) is initial probability of Y class.

4.3.4 K-nearest neighbor

In KNN the classifier has taken the investigation information into the class that has larger part of votes among k neighbors [21]. The basis by which closest neighbors are acquired amid a test is Euclidian distance. The classifier put away all the training data in the occurrence of KNN and when testing information is given to the classifier it look for its k closest information indicates and the information is marked as the new set that comprises widely held of its k neighbors. The estimations of k fluctuate upon the client.

CHAPTER 5

RESULTS AND DISCUSSION

This chapter describes the hardware tools, which is the requirement of the proposed work followed by the conclusions. This also comprises the result of the proposed work.

5.1 Following are the Hardware and Software Requirements

1) Software Used

Matlab R2017 a

Window7

2) Hardware Used

2GB RAM

Dual core processor

500 GB Hard Disk

5.2 Results and discussions

For carrying out opinion mining for reviews, A Bollywood movie ‘Padmaavat’ has been selected. Total 701 reviews have been collected from varies websites (www.imdb.com, bookmyshow.com, and other google users reviews) in which reviews were divided into three categories. Few people are in the favour of the movie and wanted to release it. And few are against of the release. Some had mix review about the movie story, caste and the release issue. All reviews were judged manually and marked into three classes as positive, negative and neutral reviews. After extraction of features, classification has been carried out by four different machine learning classifiers. The results for all classifiers using a different combination of features have been displayed in this chapter in the form of tables and graphs.

	A	B
459	Negative	action in this film. Slow motion, Overglorifying foolish acts and senseless expressions. I am not sure what did they do other than kissing and some scenes look like commercials. Its like Ads coming back to back with jewelry. Well only thing I appreciate is camera work without which it would have been unbearable to watch. I expected a lot after hearing great box office collections but this really sucks. I think this may be a fine movie to many girls as it overglorified feminism. I don't think suicide is such a great act especially when you can something else. Poor story. The climax looked like TV Serial. So slow and stupid expressions. Half of the movie you would see actors staring at each other for no reason. Remaining half is
460	Negative	Bhansali's been doing same things again and again that this was really boring for the most part. Anything good in it is parts that look carbon copies of his better films. The production, costume and cinematography are great although CGI didn't look that good. The glorification of Rajput idiocracy and patriarchy is what the film runs on and creating wooden characters where the actors can't do much more isn't helping either. The only character that is given any focus is the villain Alauddin Khilji who has been transformed into the personification of all kinds of evil that it isn't remotely convincing. Although a fictionalized tale, it is based on true historical figures and incidents which have been gone through so much contrivance and still there is nothing above archetypes. The controversy created have only helped the film and those who started it switched sides finding the glorification satisfying and at the end the common people suffer the collateral damages of the riots. Many who supported the film before the release stand as fools.
461	Negative	What a disastrous show of acting by Ranveer Singh and others. Singh made it look as if someone from today time travelled to the old times and lived there. He just have in him. Poor dialogues, cinematography make it a big flop. Need I say less about the controversy created only for political gains & in this country of vote bank politics, illiteracy can really triumph at any level.
		Not a bad movie if you want to take it as a typical Indian masala movie filled with melo-drama and crinrv. unrealistic

Figure 5.1: A screenshot from the collected reviews for movie ‘Padmaavat’

Figure 5.1 represents the snapshots of collected reviews through different sites like www.imdb.com, bookmyshow.com and many google user reviews. All these reviews are extracted manually from these sites by the author.

All the reviews are divided in three classes, namely positive, negative and neutral based from the keywords.

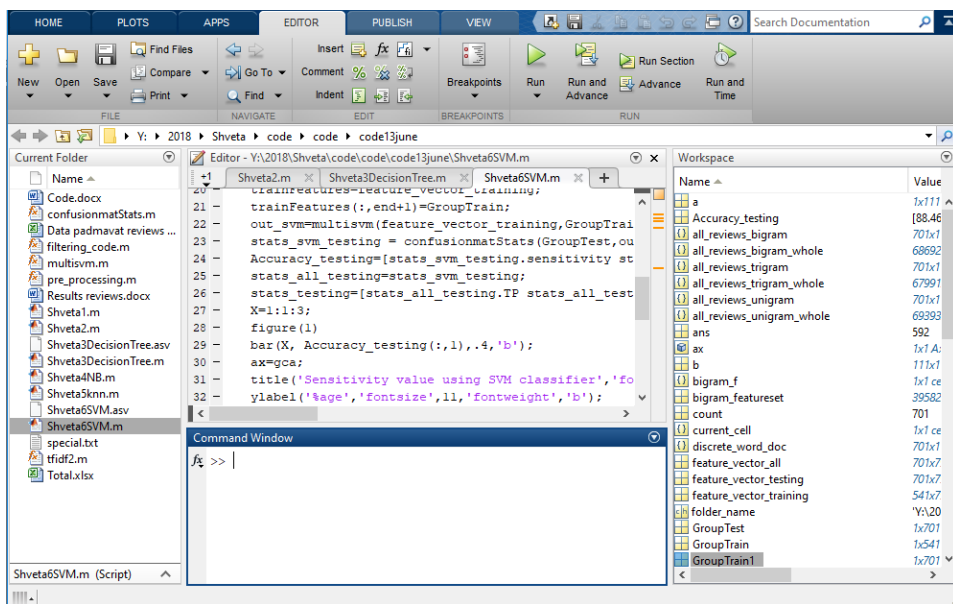


Figure 5.2: MATLAB coding of the proposed work

Figure 5.2 represents the snapshot of implementation in MATLAB.

Table 5.1: Samples of each Review category in a collected dataset of total 701 reviews

Category of reviews	Number of reviews
Positive	442
Negative	210
Neutral	49

Table 5.1 reveals the information about total number of reviews, along with the number of Positive, Negative and neutral reviews.

5.2.1 Classification Matrix

The execution of classifier is examined using confusion matrix [26]. It shows the quantity of right and wrong estimates generate by the model contrasted and genuine groupings in the test information. This is an n-x-n cluster indicating connections amongst genuine and anticipated modules, where n is quantity of classes. A confusion matrix is a representation instrument normally utilized as a part of regulated learning. All division of the network represents the instances in an anticipated class, while each column demonstrates the cases in a genuine class.

The advantage of this matrix is anything but difficult to check whether the framework is confounding two classes (i.e. regularly mislabeling one as another).

A confusion matrix comprises data of real and anticipated classifications is performed by using a classification system. A confusion matrix is made by utilizing estimations of, False Positives, True Negatives, True Positives and False Negatives.

Implementation will represent the four different possible outcomes of a single prediction. And specially when a two situation has to be described for a two-class case with classes “1” (“yes”) and “0” (“no”).

A false positive is when the outcome is incorrectly classified as “yes” (or “positive”), when it is in fact “no” (or “negative”). A false and a negative conclusion is when the outcome is incorrectly classified as negative when it is in fact positive. True positives and true negatives are obviously correct classifications.

Table 5.2: Confusion matrix for two class problem

Actual Class	Predicted Class	
	Yes	No
Yes	TP	FN
No	FP	TN

Table 5.2 represents the general confusion matrix for two class (yes and no) problem.

5.2.1.1 Accuracy

The classification accuracy is summarized in the form of a confusion matrix to the test data. It represent the extent by which classifiers correctly classify the extent. It is defined by the ratio of the number of correctly classified patterns (TP and TN) to the total number of patterns (species) classified [26].

$$Accuracy = \frac{TP + TN}{TP + TN + FP + FN} \dots(5.1)$$

5.2.1.2 Sensitivity

The sensitivity of a classifier is defined by the ratio of samples correctly classified to specific species class. It is defined by equation below [26]

$$Sensitivity = \frac{TP}{TP + FN} \dots(5.2)$$

5.2.1.3 Specificity

The specificity is the fraction of normal species correctly classified as normal class.

$$Specificity = \frac{TN}{TN + FP} \dots(5.3)$$

5.3 Classification results of Movie Reviews using TF-IDF features of unigram, bigram, trigram and SVD of 3-Gram TF-IDFs

Table 5.3: Confusion matrix for opinion classification on movie reviews using different machine learning classifiers by concatenation of proposed feature set

Decision tree classifier				
Total		Predicted Class		
		Positive	Negative	Neutral
Actual Class	Positive	391	51	0
	Negative	0	184	26
	Neutral	0	0	49
Naïve-Bayesian				
Total		Predicted Class		
		Positive	Negative	Neutral
Actual Class	Positive	442	0	0
	Negative	60	150	0
	Neutral	0	0	49
K-nearest neighbor				
Total		Predicted Class		
		Positive	Negative	Neutral
Actual Class	Positive	345	97	0
	Negative	0	203	7
	Neutral	0	0	49
Support Vector machine				
Total		Predicted Class		
		Positive	Negative	Neutral
Actual Class	Positive	442	0	0
	Negative	55	151	4
	Neutral	0	0	49

Table 5.3 represents the Confusion matrix for opinion classification on movie reviews using different machine learning classifiers (Decision tree, Naive-Bayes, knn, SVM) using TF-IDF of 3-gram features. Table 5.3 shows classification of opinions on movie reviews which is based on three classes of opinion namely positive opinion, negative opinion and neutral opinion.

Table 5.4: Sensitivity, specificity and accuracy values using different machine learning classifiers by concatenation of proposed feature set.

Decision tree classifier							
Type of Reviews	True +ive	True -ive	False +ive	False -ive	sensitivity	specificity	accuracy
Positive	391	259	0	51	0.88465	1	0.9272
Negative	184	440	51	26	0.8761	0.8961	0.8901
Neutral	49	626	26	0	1	0.9601	0.9629
Average accuracy							0.94383
Naïve-Bayesian							
Type of Reviews	True +ive	True -ive	False +ive	False -ive	sensitivity	specificity	accuracy
Positive	442	199	60	0	1	0.7683	0.9144
Negative	150	491	0	60	0.7142	1	0.9144
Neutral	49	652	0	0	1	1	1
Average accuracy							0.94293
K-nearest neighbor							
Type of Reviews	True +ive	True -ive	False +ive	False -ive	sensitivity	specificity	accuracy
Positive	345	259	0	97	0.7805	1	0.8616
Negative	203	394	97	7	0.9666	0.8024	0.8516
Neutral	49	645	7	0	1	0.9892	0.9900
Average accuracy							0.90106
Support Vector machine							
Type of Reviews	True +ive	True -ive	False +ive	False -ve	sensitivity	specificity	accuracy
Positive	442	204	55	0	1	0.7876	0.9215
Negative	151	491	0	59	0.7190	1	0.9158
Neutral	49	648	4	0	1	0.9938	0.9942
Average accuracy							0.92673

Table 5.4 represents Sensitivity, Specificity and Accuracy values using different machine learning classifiers with 3-gram TF-IDF and SVD of 3-gram TF-IDF features. Table 5.4

shows the true positive, true negative, false positive, false negative values on different opinions of user reviews by using classifiers (Decision tree, Naive-Bayes, knn, SVM).

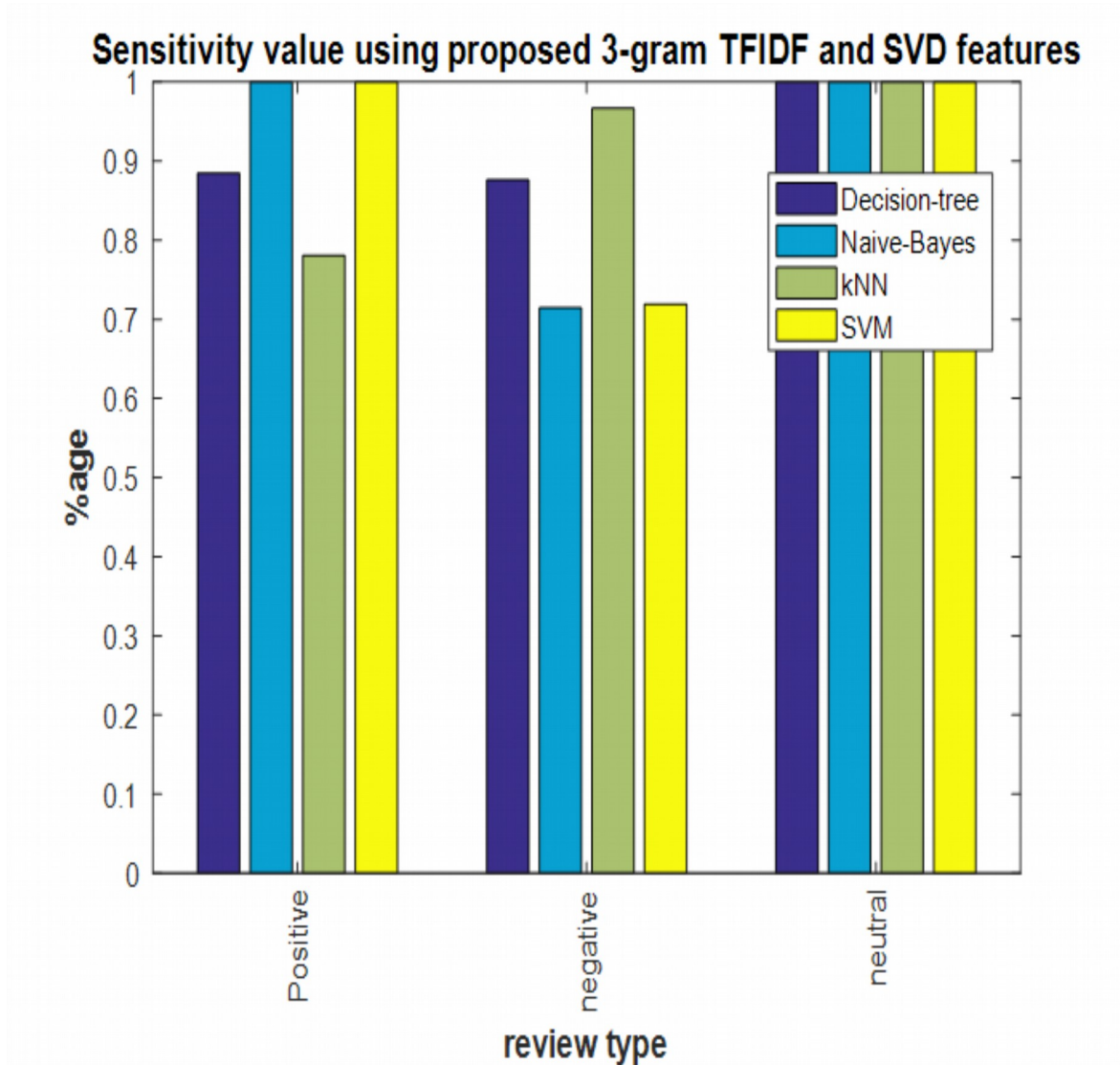


Figure 5.3: Sensitivity value using proposed 3-gram TFIDF and SVD features

Figure 5.3 represent Sensitivity value from different classifiers (Decision tree, Naive-Bayes, knn, SVM) by using 3-gram TF-IDF and Sensitivity value by using SVD of 3-gram TF-IDF features.

In figure 5.3 x-axis represent the type of opinion on movie reviews and y-axis shows the percentage of Sensitivity.

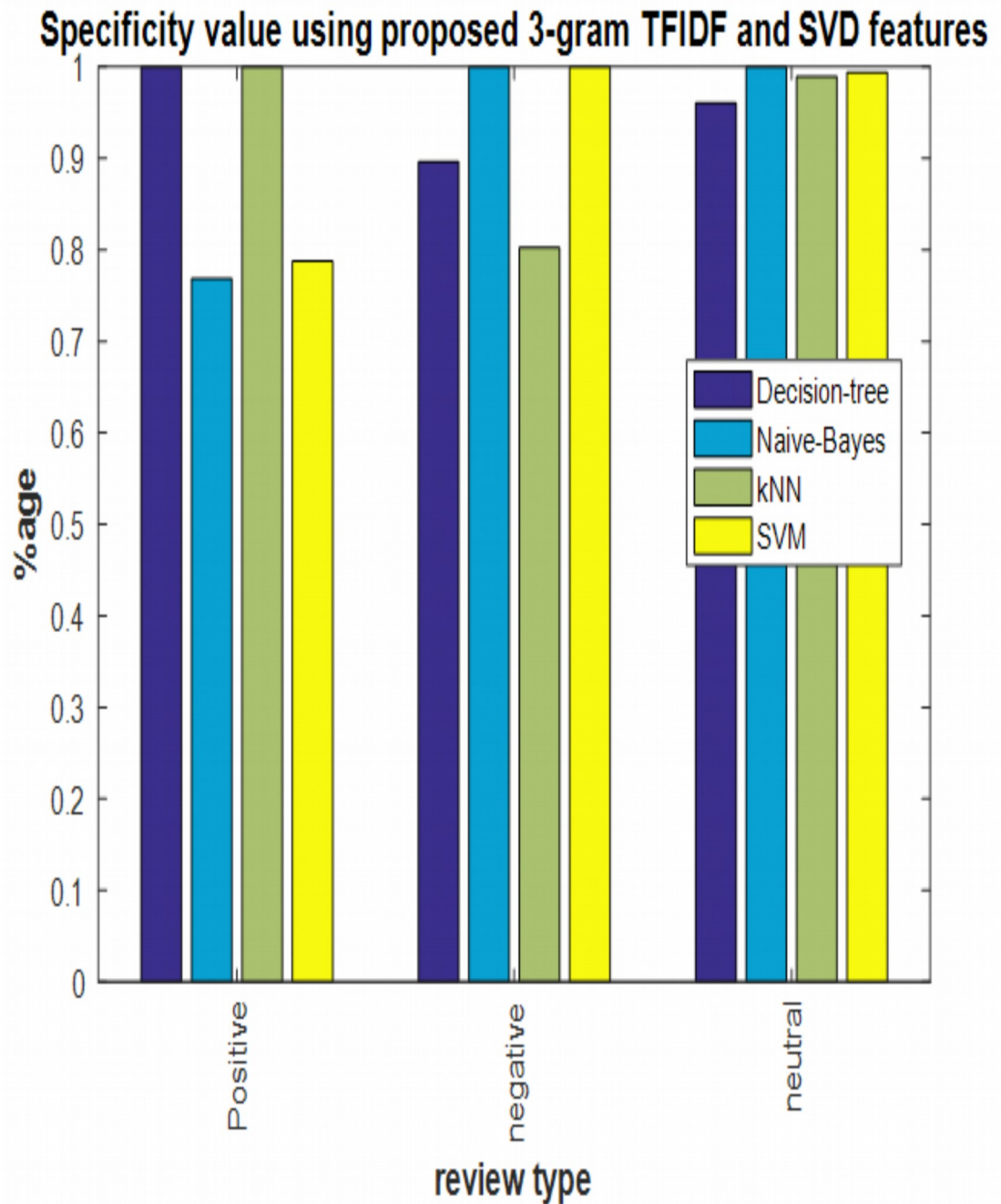


Figure 5.4: Specificity value using proposed 3-gram TFIDF and SVD features

Figure 5.4 represent Specificity value from different classifiers (Decision tree, Naive-Bayes, knn, SVM) by using 3-gram TF-IDF and Specificity value by using SVD of 3-gram TF-IDF features. In figure 5.4 x-axis represent the type of opinion on movie reviews and y-axis shows the percentage of Specificity.

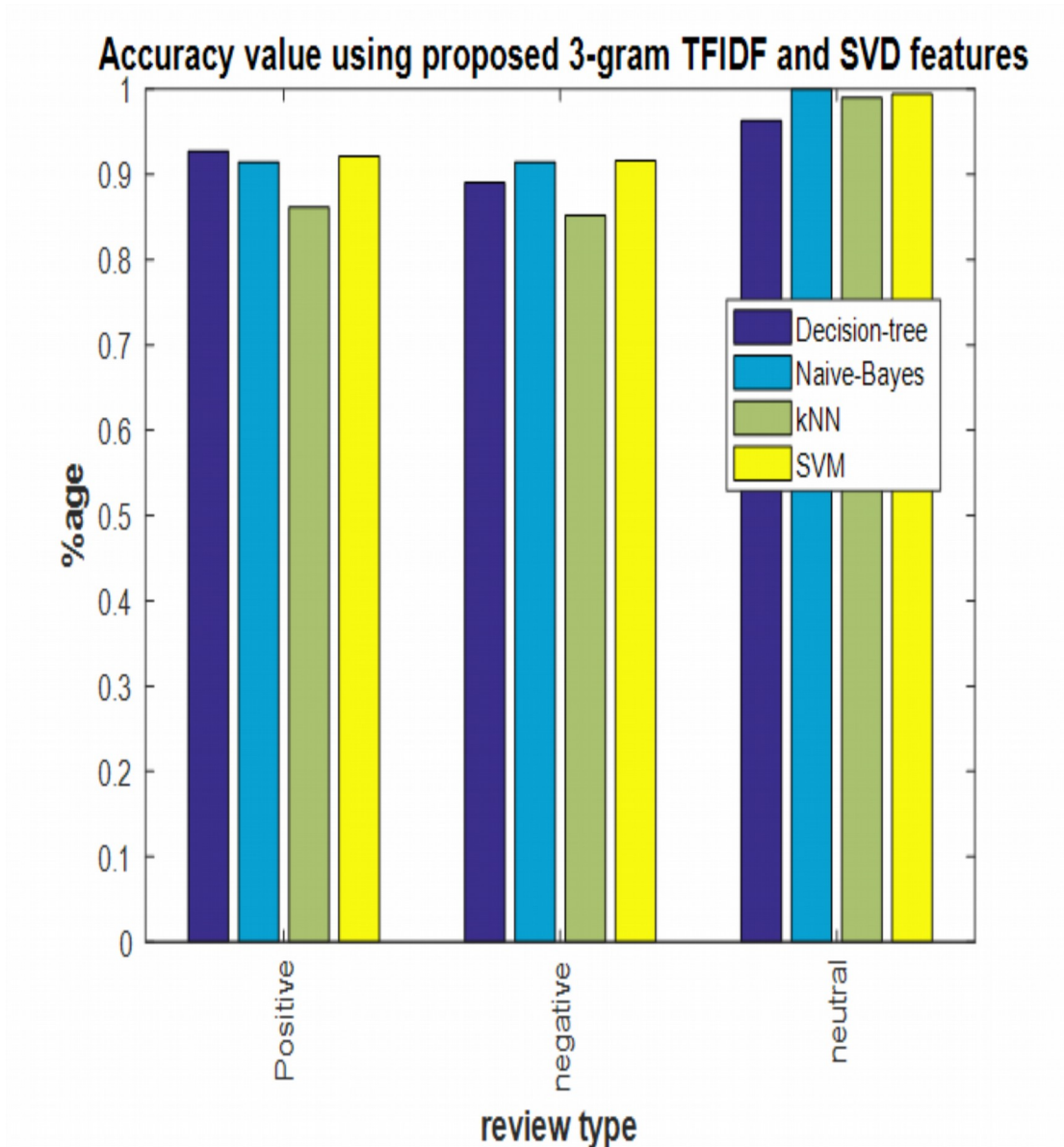


Figure 5.5: Accuracy value using proposed 3-gram TFIDF and SVD features

Figure 5.5 represent Accuracy value from different classifiers (Decision tree, Naive-Bayes, knn, SVM) by using 3-gram TF-IDF and accuracy value by using SVD of 3-gram TF-IDF features. In figure 5.5 x-axis represent the type of opinion on movie review and y-axis shows the percentage of Accuracy.

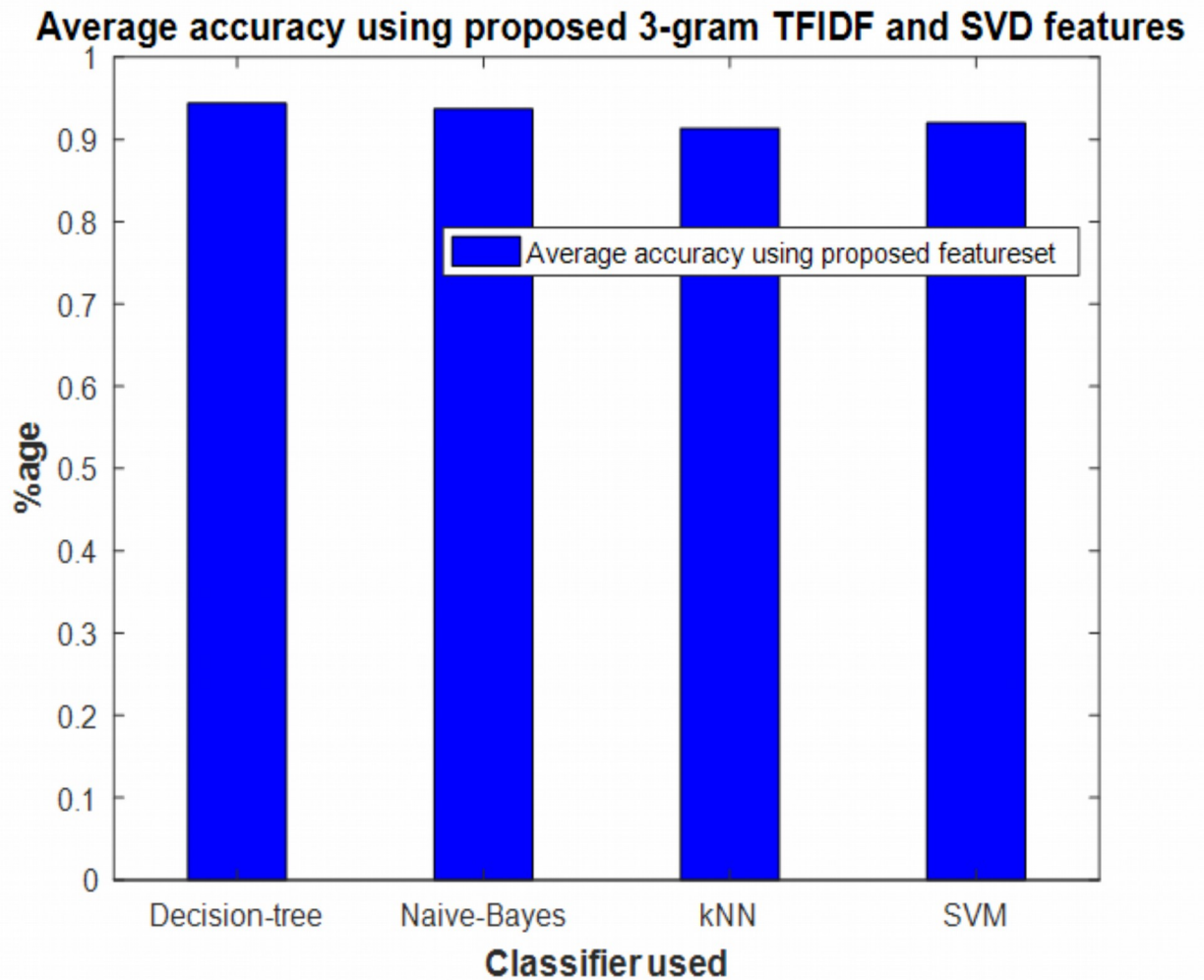


Figure 5.6: Average accuracy after opinion classification on movie reviews from different machine learning classifiers using 3-gram TF-IDF and SVD of 3-gram TF-IDF features.

Figure 5.6 represents average accuracy after opinion classification on movie reviews from different classifiers using 3-gram TF-IDF and average accuracy from different classifiers using SVD of 3-gram TF-IDF features. In figure 5.6 x-axis shows name of the classifiers (Decision tree, Naive Bayes, kNN, SVM) used in classification and y-axis shows percentage of average accuracy.

5.4 Classification results using TF-IDF features of unigram and SVD of TF-IDF unigram

Table 5.5: Confusion matrix for reviews classification using different machine learning classifiers using TF-IDF of unigram features

Decision tree classifier				
Total		Predicted Class		
		Positive	Negative	Neutral
Actual Class	Positive	391	51	0
	Negative	0	181	29
	Neutral	0	0	49
Naïve-Bayesian				
Total		Predicted Class		
		Positive	Negative	Neutral
Actual Class	Positive	442	0	0
	Negative	63	147	0
	Neutral	0	0	49
K-nearest neighbor				
Total		Predicted Class		
		Positive	Negative	Neutral
Actual Class	Positive	344	98	0
	Negative	0	200	10
	Neutral	0	0	49
Support Vector machine				
Total		Predicted Class		
		Positive	Negative	Neutral
Actual Class	Positive	439	1	2
	Negative	41	150	19
	Neutral	15	2	32

Table 5.5 represents the Confusion matrix for opinion classification on movie reviews using different machine learning classifiers (Decision tree, Naive-Bayes, kNN, SVM) using TF-IDF of unigram features. Table 5.5 shows classification of opinions on movie reviews which is based on three classes of opinion namely positive opinion, negative opinion and neutral opinion.

Table 5.6: Sensitivity, specificity and accuracy values using different machine learning classifiers using TF-IDF of unigram and SVD of unigram TF-IDF features

Decision tree classifier							
Type of Reviews	True +ive	True -ive	False +ive	False -ive	sensitivity	specificity	accuracy
Positive	391	259	0	51	0.8846	1	0.9272
Negative	181	440	51	29	0.8619	0.8961	0.8858
Neutral	49	623	29	0	1	0.9555	0.9586
Average accuracy							0.92386
Naïve-Bayesian							
Type of Reviews	True +ive	True -ive	False +ive	False -ive	sensitivity	specificity	accuracy
Positive	442	196	63	0	1	0.756757	0.9101
Negative	147	491	0	63	0.7000	1	0.9101
Neutral	49	652	0	0	1	1	1
Average accuracy							0.94006
K-nearest neighbor							
Type of opinion	True +ive	True -ive	False +ive	False -ive	sensitivity	specificity	accuracy
Positive	344	259	0	98	0.7782	1	0.8601
Negative	200	393	98	10	0.9523	0.8004	0.8459
Neutral	49	642	10	0	1	0.9846	0.9857
Average accuracy							0.89723
Support Vector machine							
Type of opinion	True +ive	True -ive	False +ive	False -ive	sensitivity	specificity	accuracy
Positive	439	203	56	3	0.9932	0.7837	0.9158
Negative	150	488	3	60	0.7142	0.9938	0.9101
Neutral	32	631	21	17	0.6530	0.9677	0.9457
Average accuracy							0.92386

Table 5.6 represents Sensitivity, Specificity and Accuracy values using different machine learning classifiers with unigram TF-IDF and SVD of unigram TF-IDF features. Table 5.6 shows the true positive, true negative, false positive, false negative

values on different opinions of user reviews by using classifiers (Decision tree, Naive-Bayes, knn, SVM).

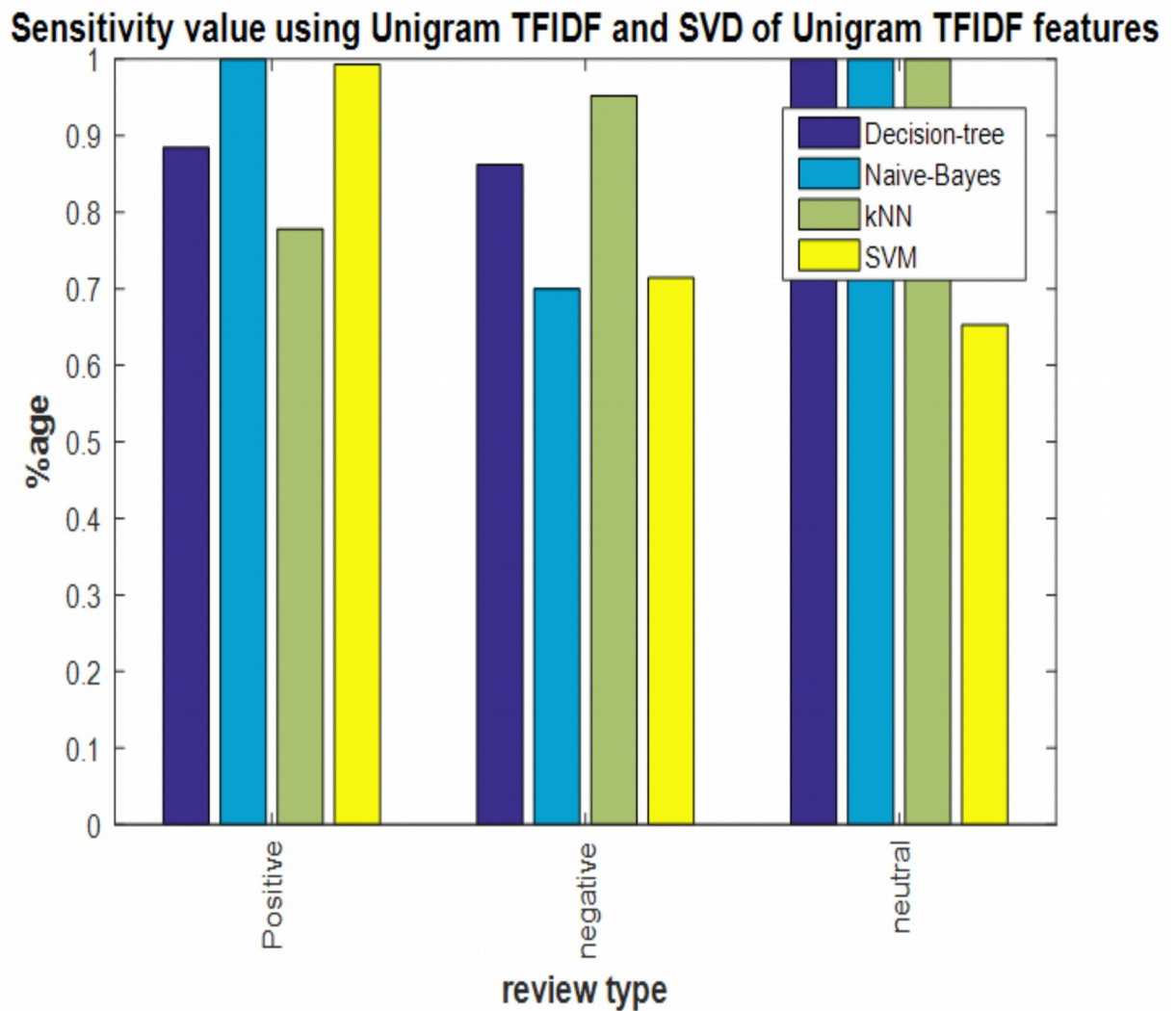


Figure 5.7: Sensitivity value using Unigram TFIDF and SVD of Unigram TFIDF features

Figure 5.7 represent Sensitivity value from different classifiers (Decision tree, Naive-Bayes, knn, SVM) by using unigram TF-IDF and accuracy value by using SVD of unigram TF-IDF features. In figure 5.7 x-axis represent the type of opinion on movie reviews and y-axis shows the percentage of Sensitivity.

Specificity value using Unigram TFIDF and SVD of Unigram TFIDF features

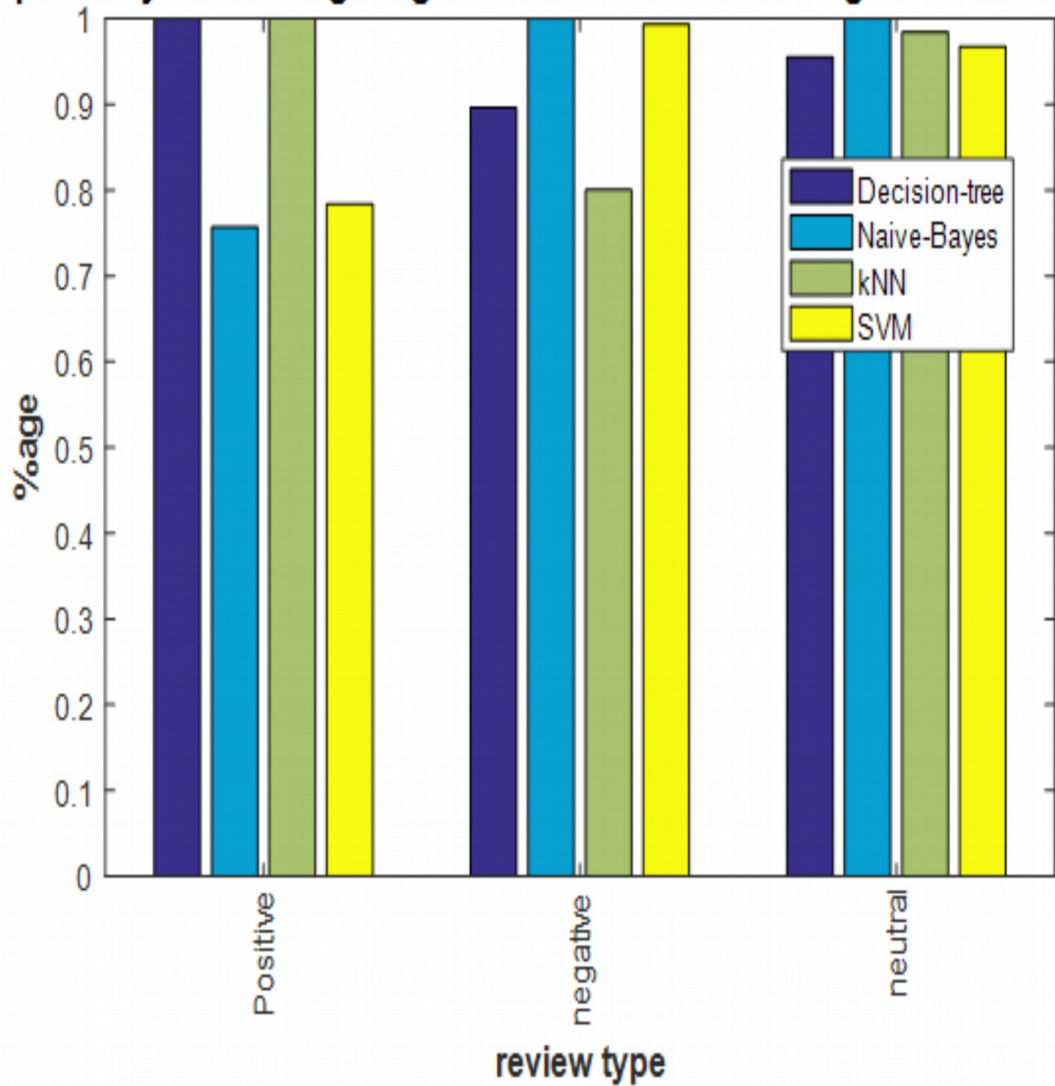


Figure 5.8: Specificity value using Unigram TF-IDF and SVD of Unigram TFIDF features

Figure 5.8 represent Specificity value from different classifiers (Decision tree, Naive-Bayes, kNN, SVM) by using unigram TF-IDF and accuracy value by using SVD of unigram TF-IDF features.

In figure 5.8 x-axis represent the type of opinion on movie reviews and y-axis shows the percentage of Specificity.

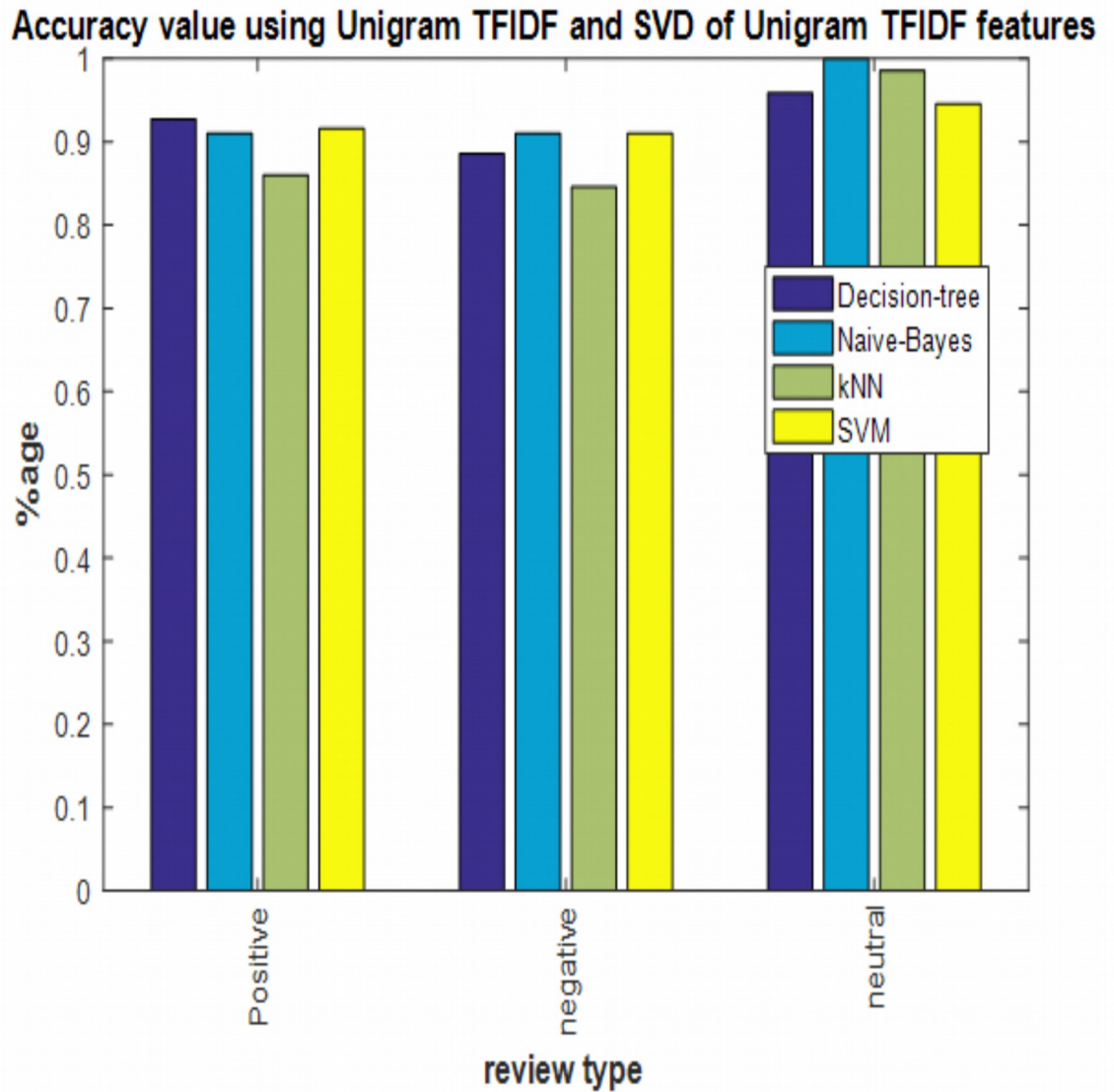


Figure 5.9: Accuracy value using Unigram TFIDF and SVD of Unigram TFIDF features

Figure 5.9 represent Accuracy value from different classifiers (Decision tree, Naive-Bayes, kNN, SVM) by using unigram TF-IDF and accuracy value by using SVD of unigram TF-IDF features.

In figure 5.9 x-axis represent the type of opinion on movie reviews and y-axis shows the percentage of Accuracy.

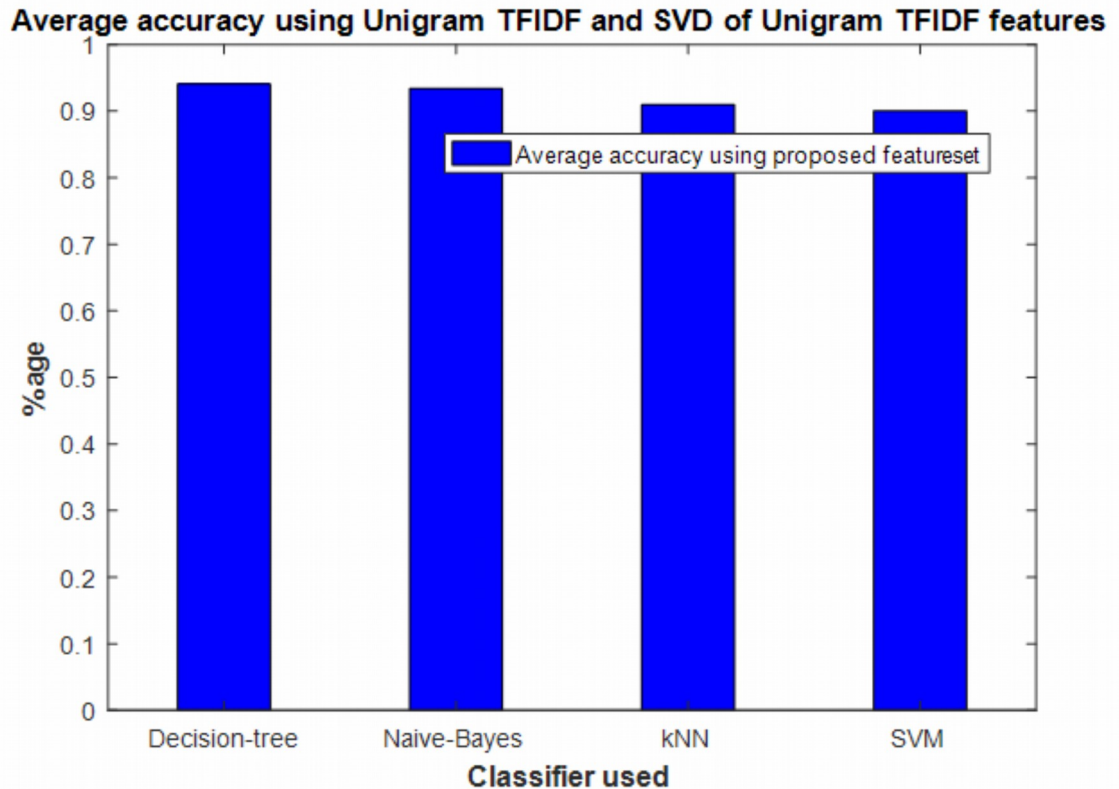


Figure 5.10: Average accuracy after opinion classification on movie reviews from different classifiers using Unigram TF-IDF and SVD of Unigram TF-IDF features

Figure 5.10 represents average accuracy after opinion classification on movie reviews from different classifiers using unigram, tf-idf, and average accuracy from different classifiers using SVD of unigram tf-idf features.

In figure 5.10 x-axis shows name of the classifiers used in classification and y-axis shows percentage of average accuracy.

Among all three tested classifiers, decision tree gives the best classification rate in true classification of category of reviews. There are three categories of 701 reviews samples out of which seventy percent has been used for training. Performance evaluation has been carried out only for tested data-sets by four different machine learning classifiers. Sensitivity, specificity and accuracy has been given in tabular as well as graphical form which shows about 94% accuracy in true classification of review documents when decision tree classifier has been used.

CHAPTER 6

CONCLUSION AND FUTURE SCOPE

6.1 Conclusion

Now a day's Sentiment Analysis have become a big issue. Explosion of social media has provided many opportunities for citizens to publicly spread their opinions , but it has created serious problems when it comes to make sense of these opinions and because of viral nature of social media importance to get an understanding of citizens, number of opinion has grown (when attention is very unevenly and frequently distributed) Some issues become important through word-of-mouth. These research works are focused on movie reviews. There are large amount of user-generated movies review which are available on the internet like IMDB, BookMyShow, Google user reviews etc.

There are many challenges like one or more bad feature of the film does not make it overall bad similarly as one or more good feature does not make it overall a good movie. Therefore, opinion mining of film review is considered more challenging than opinion mining of other type of reviews. In our work, Padmaavat movie or film reviews are collected in order to explore opinion mining. At first collected reviews are filtered in which stopwords, special characters, Numeric's etc. are removed. Then four different features named as TF-IDF of unigram, bigram and trigram word lexicons and SVD of all the concatenated TF-IDF feature vectors are extracted. Then these features are concatenated and fed to four different machine learning classifiers.

Among all three tested classifiers, decision tree gives the best classification rate in true classification of category of reviews. There are three categories of 701 reviews samples out of which seventy percent has been used for training. Performance evaluation has been carried out only for tested data-sets by four different machine learning classifiers. Sensitivity, specificity and accuracy have been given in tabular as well as graphical form which shows about 94% accuracy in true classification of review documents when decision tree classifier has been used.

6.2 Future scope

Proposed method has been tested only on collected dataset for a conflicted movie 'Padmaavat' but does not tested on standard datasets available on internet. Also feature reduction or selection is not explored which can decrease complexity and computation time of the machine learning classifiers.

REFERENCES

- [1] A. A. Hakim, A. Erwin, K. Eng, M. Galinium, and W. Muliady, Automated document classification for news article in bahasa Indonesia based on term frequency inverse document frequency (TFIDF) approach, in 6th International Conference on Information Technology and Electrical Engineering (ICITEE), 1-4, 2014.
- [2] A. M. Mahmood, N. Satuluri, and M. R. Kuppaa, An Overview of Recent and Traditional Decision Tree Classifiers in Machine Learning, International Journal of Research and Reviews in Ad Hoc Networks, Vol. 1, No.1, 2011
- [3] A. S. Manek, P. D. Shenoy, M.C. Mohan, K. R. Venugopal, Aspect term extraction for sentiment analysis in large movie reviews using Gini Index feature selection method and SVM classifier, Published in: World Wide Web, 20, 135–154, 2017
- [4] A. Kukkara, R. Mohana, A Supervised Bug Report Classification with Incorporate and Textual field Knowledge, Procedia Computer Science, 132, 352-361, 2018.
- [5] B. Bansal, S. Srivastava, Sentiment classification of online consumer reviews using word vector representations, Procedia Computer Science, 132, 1147–1153, 2018.
- [6] B. Liu, Sentiment Analysis and Opinion Mining, Human Language Technologies, 2012.
- [7] B. Trstenjak, S. Mikac, D. Donko, KNN with TF-IDF based Framework for Text Categorization, Procedia Engineering, 69, 1356-1364, 2014.
- [8] C. Akimushkin, D. R. Amancio and O. N. Oliveira, On the role of words in the network structure of texts: Application to authorship attribution, Physica A: Statistical Mechanics and its Applications, 495, 49-58, 2018.
- [9] D. Kim, D. Kim, E. Hwang and H.G. Choi, A user opinion and metadata mining scheme for predicting box office performance of movies in the social network environment, New Review of Hypermedia and Multimedia, 19, 3–4, 259–272, 2013.

- [10] E. B. Putranto, P. A. Situmorang and A. S. Girsang, Face recognition using eigenface with naive Bayes, 11th International Conference on Knowledge, Information and Creativity Support Systems (KICSS), Yogyakarta, 1-4, 2016.
- [11] E. M. Taylor, J. D. Velásquez, F. B. Marquez, A novel deterministic approach for aspect-based opinion mining in tourism products reviews, *Expert Systems with Applications*, 41, 7764–7775, 2014.
- [12] G. R. BRINDHA, P. SWAMINATHAN, B. SANTHI, Performance analysis of new word weighting procedures for opinion mining, *Front Inform Technol Electron Eng*, 17(11),1186-1198, 2016.
- [13] G. Domeniconi, G. Moro, R. Pasolini, C. Sartori, A Comparison of Term Weighting Schemes for Text Classification and Sentiment Analysis with a Supervised Variant of tf.idf, *Data Management Technologies and Applications*, 39-58, 2015.
- [14] H. Han, G. Karypis, V. Kumar, Text Categorization Using Weight Adjusted k-Nearest Neighbor Classification, *PAKDD*, 53-65, 2001.
- [15] H. Uguz, A two-stage feature selection method for text categorization by using information gain, principal component analysis and genetic algorithm, *Knowledge-Based Systems*, 24,1024–1032, 2011.
- [16] H. Niemann, M. G. Moehrle, J. Frischkorn, Use of a new patent text-mining and visualization method for identifying patenting patterns over time: Concept, method and test application, *Technological Forecasting and Social Change*,115, 210-220, 2017.
- [17] H. L. Yang, Q. F. Lin, Opinion mining for multiple types of emotion-embedded products/services through evolutionary strategy, *Expert Systems With Applications* , 99, 44–55, 2018.
- [18] J. T. Kwok, Automatic Text Categorization Using Support Vector Machine, *Proceedings of International Conference on Neural Information Processing*, 347-351, 1998.

- [19] J. Han, M. Kamber and J. Pei, Data Mining Concepts and Techniques (3rd Edition), MK MORGAN KAUFMANN , 2012
- [20] L.Wang , X. Zhao, Improved KNN classification algorithms research intext categorization, IEEE, 2012.
- [21] M. Hariri, F. Hakimi and F. Gharehbaghi, Image-Splicing Forgery Detection Based On Improved LBP and K-Nearest Neighbors Algorithm, Electronics Information and Planning, 1074-1077, 2015.
- [22] M. Kang, J. Ahn, K. Lee, Opinion mining using ensemble text hidden Markov models for text classification, Expert Systems With Applications, 94, 218–227, 2018.
- [23] M. G. Armentano, S. Schiaffino, I. Christensen and F. Boato, Movies Recommendation Based on Opinion Mining in Twitter, Mexican International Conference on Artificial Intelligence Advances in Artificial Intelligence and Its Applications, 80-91, 2015.
- [24] M. Y.Su, Using clustering to improve the KNN-based classifiers for online anomaly network traffic identification, Journal of Network and Computer Applications 34 ,722–730, 2011.
- [25] N. Jindal and B. Liu, Identifying comparative sentences in text documents, in Proceedings of the 29th Annual International ACM SIGIR Conference on Research and Development in Information Retrieval, ACM, 244–251, 2006.
- [26] N. K Bose, P. Liang, neural network fundamentals with graphs, algorithms and applications, 1995.
- [27] N. M. Shelke, S. Deshpande and V. Thakre, Survey of Techniques for Opinion Mining, International Journal of Computer Applications, 57(13), 0975 – 8887, 2012.
- [28] O. A. S. Ibrahim, D. L. Silva, Term frequency with average term occurrences for textual information retrieval, Soft Computing , 20, 8, 3045–3061, 2016.
- [29] O. S. Kwon, J. Kim, K. H. Choi, Y. Ryu and J. E. Park, Trends in deqi research: a text mining and network analysis, Integrative Medicine Research, 2018.

- [30] V. Podgorelec, P. Kokol, B. Stiglic, I. Rozman, Decision trees: An overview and their use in medicine, *Journal of Medical Systems*, 26, 445–463, 2002.
- [31] Q. I. Mahmud, A. Mohaimen, M. S. Islam and M. E. Jannat, A support vector machine mixed with statistical reasoning approach to predict movie success by analyzing public sentiments, 20th International Conference of Computer and Information Technology (ICCIT), Dhaka, 1-6, 2017.
- [32] R. Dong, M. P. O’Mahony, M. Schaal, K. McCarthy and B. Smyth, Combining similarity and sentiment in opinion mining for product recommendation, *Journal of Intelligent Information System*, 46, 285–312, 2016.
- [33] S. Liu, K. Huang and J. Chai, Research of news text with word frequency statistics and user information, 3rd IEEE International Conference on Computer and Communications (ICCC), Chengdu, 2633-2637, 2017.
- [34] S. K. Srivastava and S. K. Singh, Multiparameter Based Performance Evaluation Of Classification Algorithms, *International Journal of Computer Science and Information Technology (IJCSIT)*, 7, 3, 2015.
- [35] T. Wilson, P. Hoffmann, S. Somasundaran, J. Kessler, J. Wiebe, Y. Choi, C. Cardie, E. Riloff and S. Patwardhan, Opinionfinder: a system for subjectivity analysis, *HLT/EMNLP on Interactive Demonstrations*, 34–35, 2005.
- [36] A. K. Uysal, S. Gunal, A novel probabilistic feature selection method for text classification, *Knowledge-Based Systems*, 36, 226–235, 2012.
- [37] V. Korde and C. N. Mahender, Text classification and classifier a survey, *International Journal of Artificial Intelligence and Applications (IJAA)*, 3, 2, 2012.
- [38] W. Zhang, T. Yoshida, and X. Tang, A comparative study of (TF-IDF), LSI and multi-words for text classification, *Expert Systems with Applications*, 38, 2758-2765, 2011.
- [39] Y. H. Hua, Y. L. Chen, H. L. Choub, Opinion mining from online hotel reviews – A text summarization approach, *Information Processing and Management*, 53, 436-449, 2017.

[40] Y. Du, J. Liu, W. Ke, X. Gong, Hierarchy construction and text classification based on the relaxation strategy and least information model, *Expert Systems with Applications*, 100, 157-164, 2018.

[41] Z. Yan, EXPRS: An extended pagerank method for product feature extraction from online consumer reviews, *Information and Management*, 52, 7, 850-858, 2015.

[42] Z. W. Zhai, B. Liu, H. Xu and P. Jia, Clustering Product Features for Opinion Mining, *Proceedings of the fourth ACM international conference on Web search and data mining*, 347-354, 2011.

ORIGINALITY REPORT

13%

SIMILARITY INDEX

7%

INTERNET SOURCES

9%

PUBLICATIONS

4%

STUDENT PAPERS

PRIMARY SOURCES

- 1** Seyyed Mohammad Hossein Dadgar, Mohammad Shirzad Araghi, Morteza Mastery Farahani. "A novel text mining approach based on TF-IDF and Support Vector Machine for news classification", 2016 IEEE International Conference on Engineering and Technology (ICETECH), 2016
Publication 1%
- 2** ijaers.com
Internet Source 1%
- 3** Submitted to Thapar University, Patiala
Student Paper 1%
- 4** Ega Bima Putranto, Poldo Andreas Situmorang, Abba Suganda Girsang. "Face recognition using eigenface with naive Bayes", 2016 11th International Conference on Knowledge, Information and Creativity Support Systems (KICSS), 2016
Publication 1%
- 5** Saurabh Kr. Srivasatava, Roshan Kumari,

Sandeep Kr. Singh. "An ensemble based NLP feature assessment in binary classification", 2017 International Conference on Computing, Communication and Automation (ICCCA), 2017

Publication

1%

6

Amarouche, Kamal, Houda Benbrahim, and Ismail Kassou. "Product Opinion Mining for Competitive Intelligence", Procedia Computer Science, 2015.

Publication

1%

7

www.ijrat.org

Internet Source

1%

8

Submitted to Higher Education Commission Pakistan

Student Paper

1%

9

Rini Wongso, Ferdinand Ariandy Luwinda, Brandon Christian Trisnajaya, Olivia Rusli, Rudy. "News Article Text Classification in Indonesian Language", Procedia Computer Science, 2017

Publication

<1%

10

eprints.nottingham.ac.uk

Internet Source

<1%

11

Submitted to Universiti Sains Malaysia

Student Paper

<1%

12

www.ijer.in

Internet Source

<1%

13

airccse.org

Internet Source

<1%

14

Submitted to University of Hertfordshire

Student Paper

<1%

15

repositorio-aberto.up.pt

Internet Source

<1%

16

Submitted to University of Warwick

Student Paper

<1%

17

dspace.thapar.edu:8080

Internet Source

<1%

18

K. M. Anil Kumar, N. Rajasimha, Manovikas Reddy, A. Rajanarayana, Kewal Nadgir.
"Analysis of users' Sentiments from Kannada Web Documents", Procedia Computer Science, 2015

Publication

<1%

19

"Emerging Research in Computing, Information, Communication and Applications", Springer Nature, 2016

Publication

<1%

20

documents.mx

Internet Source

<1%

research.ijcaonline.org

21

Internet Source

<1%

22

Franca, Andre L., Ricardo Jasinski, Paulo Cemin, Volnei A. Pedroni, and Altair Olivo Santin. "The energy cost of network security: A hardware vs. software comparison", 2015 IEEE International Symposium on Circuits and Systems (ISCAS), 2015.

Publication

<1%

23

Submitted to University of Edinburgh

Student Paper

<1%

24

Submitted to City University of Hong Kong

Student Paper

<1%

25

Submitted to Guru Nanak Dev Engineering College

Student Paper

<1%

26

diva-portal.org

Internet Source

<1%

27

Submitted to Sim University

Student Paper

<1%

28

G. Desjardins, R. Proulx, R. Godin. "An Auto-Associative Neural Network for Information Retrieval", The 2006 IEEE International Joint Conference on Neural Network Proceedings, 2006

Publication

<1%

29

[Submitted to UM, Baltimore County](#)

Student Paper

<1%

30

link.springer.com

Internet Source

<1%

31

S. Geetha, Siva S. Sivatha Sindhu, S. Gobi, A. Kanna. "Evolving GA Classifier for Audio Steganalysis based on Audio Quality Metrics", 2006 Fourth International Conference on Intelligent Sensing and Information Processing, 2006

Publication

<1%

32

Xiangji Huang, Yan Huang, Miao Wen, Aijun An, Yang Liu, Josiah Poon. "Applying Data Mining to Pseudo-Relevance Feedback for High Performance Text Retrieval", Sixth International Conference on Data Mining (ICDM'06), 2006

Publication

<1%

33

ijcsit.com

Internet Source

<1%

34

www.communispace.com

Internet Source

<1%

35

chandigarhcity.citybase.in

Internet Source

<1%

36

scholarworks.sjsu.edu

Internet Source

<1%

37

www.coursehero.com

Internet Source

<1%

38

m.eurekaselect.com

Internet Source

<1%

39

Wen-Feng Xuan, Bing-Quan Liu, Cheng-Jie Sun, De-Yuan Zhang, Xiao-Long Wang.

"Finding main topics in blogosphere using document clustering based on topic model",
2011 International Conference on Machine Learning and Cybernetics, 2011

Publication

<1%

40

ijasret.com

Internet Source

<1%

41

digital.library.unt.edu

Internet Source

<1%

Exclude quotes On

Exclude matches < 8 words

Exclude bibliography On