# A STUDY ON PERSRONALISED RECOMMENDER SYSTEM

# USING SOCIAL MEDIA

BY

AISHNIVYA A/P BALAMURUGAN

A REPORT

SUBMITTED TO

Universiti Tunku Abdul Rahman

in partial fulfillment of the requirements

for the degree of

BACHELOR OF INFORMATION SYSTEMS (HONS)

INFORMATION SYSTEMS ENGINEERING

Faculty of Information and Communication Technology

(Kampar Campus)

MAY 2020

**UNIVERSITI TUNKU ABDUL RAHMAN**

# REPORT STATUS DECLARATION FORM

**Title**:  A STUDY ON PERSONALISED RECOMMENDER SYSTEM
 USING SOCIAL MEDIA

**Academic Session**: MAY 2020

I  AISHNIVYA A/P BALAMURUGAN

**(CAPITAL LETTER)**

declare that I allow this Final Year Project Report to be kept in

Universiti Tunku Abdul Rahman Library subject to the regulations as follows:

1. The dissertation is a property of the Library.
2. The Library is allowed to make copies of this dissertation for academic purposes.

Verified by,

_____  _____

(Author's signature)  (Supervisor'ssignature)

**Address**:

NO12B, JLN BUNGARAYA

 TMN CHANGKAT JAYA  DR.RAMESH KUMAR AYYASAMY

 Supervisor's name

**Date**: 9 SEPTEMBER 2020  **Date**: 9 SEPTEMBER 2020

**A STUDY ON PERSRONALISED RECOMMENDER SYSTEM**

**USING SOCIAL MEDIA**

BY

AISHNIVYA A/P BALAMURUGAN

A REPORT

SUBMITTED TO

Universiti Tunku Abdul Rahman

in partial fulfillment of the requirements

for the degree of

BACHELOR OF INFORMATION SYSTEMS (HONS)

INFORMATION SYSTEMS ENGINEERING

Faculty of Information and Communication Technology

(Kampar Campus)

MAY 2020

**DECLARATION OF ORIGINALITY**

I declare that this report entitled "**A STUDY OF PERSONALISED RECOMMENDER SYSTEM USING SOCIAL MEDIA**" is my own work except as cited in the references. The report has not been accepted for any degree and is not being submitted concurrently in candidature for any degree or other award.

Signature        :        _____

Name             :        <u>AISHNIVYA A/P BALAMURUGAN</u>

Date             :        <u>10 SEPTEMBER 2020</u>

BIS (Hons) Information Systems Engineering
Faculty of Information and Communication Technology (Kampar Campus), UTAR

# ACKNOWLEDGEMENTS

From the bottom of my heart, I would be thankful to those who had stood by my side in accomplishing this project. My supervisor DR. Ramesh Kumar Ayyasamy had done a great role in guiding on the project since the beginning regardless of his hectic day and restless schedules. He had given an opportunity to be in such a good title in current situations of the economy and the world.

Nevertheless, family. Family that had been gifted to me had tolerated me with my workloads . Always had and will be a great guidance and boost the inner strength to proceed with my tasks. They were always been who am I in the project and the backbone of my efforts that I pour in to produce an outcome that would be vital.

BIS (Hons) Information Systems Engineering
Faculty of Information and Communication Technology (Kampar Campus), UTAR

# ABSTRACT

Recommender system means to give user direction identified with the helpful services dependent on their personalised service suggestion, behaviour or neighbor's inclinations. With the popularity of social network, numerous clients like to share their perspectives via social networking media, for example, rating, sites, tweets and so on, which prescribes clients interest item. Personalised recommender framework gives a suitable recommendation like shopping, scheduling, hotels, tags, motion pictures and so on, which delivers enormous information on the web. This outcomes in the issue of data over-burden. To get over this issue, Personalized Recommendation System have been prosperously utilized. This paper studies personalised recommender system using social media. The problem that is discussed the research study is information overload which is a hot topic of current world. Objectives to condense the data , analyze the users' behaviour before recommending items is much important.

In the research study Naive Bayes Theorem classifier , k-Nearest Neighbor Classifier and Support Vector Machine classifier is used. These machine learning algorithm processes the data set obtained. The evaluation on these algorithm is done to evaluate accuracy of the algorithm. Therefore is also well stated the recommendation has a great play is almost everything. This assistance from our natural components gives us a basic technique to find the best choice without having a great deal of effort to isolating through the different choices available in the market. At this moment advancement, the Recommendation system is an application that isolated altered information and gives the best way to deal with understand a user's taste and to propose reasonable things to them by considering the models among their inclinations and research studies of various things. The research project has well stated the accuracy and determined performance evaluation metrics of the Twitter data set after going through the process of preprocessing and sentiment analyzing.

**Table of Contents**

BIS (Hons) Information Systems Engineering
Faculty of Information and Communication Technology (Kampar Campus), UTAR

BIS (Hons) Information Systems Engineering
Faculty of Information and Communication Technology (Kampar Campus), UTAR

# LIST OF FIGURES

BIS (Hons) Information Systems Engineering
Faculty of Information and Communication Technology (Kampar Campus), UTAR

BIS (Hons) Information Systems Engineering
Faculty of Information and Communication Technology (Kampar Campus), UTAR

BIS (Hons) Information Systems Engineering
Faculty of Information and Communication Technology (Kampar Campus), UTAR

## LIST OF TABLES

| Table Number | Description | Page |
|---|---|---|
| Table 2.2.1 | Analysis on Literature review 1 | 58 |
| Table 2.2.2 | Analysis on Literature review 2 | 59 |
| Table 2.2.3 | Analysis on Literature review 3 | 60 |
| Table 2.2.4 | Analysis on Literature review 4 | 60 |
| Table 2.2.5 | Analysis on Literature review 5 | 61 |
| Table 2.2.6 | Analysis on Literature review 6 | 61 |
| Table 2.2.7 | Analysis on Literature review 7 | 62 |
| Table 2.2.8 | Analysis on Literature review 8 | 62 |
| Table 2.2.9 | Analysis on Literature review 9 | 63 |
| Table 2.2.10 | Analysis on Literature review 10 | 64 |
| Table 5.1 | Dataframe of dataset 1 | 112 |
| Table 5.2 | Dataframe of dataset 2 | 113 |
| Table 5.3 | Steps Verification | 116 |
| Table 6.1 | Evaluation score of Naive Bayes Classifier | 119 |
| Table 6.2 | Evaluation score of KNN | 119 |
| Table 6.3 | Evaluation scores of Support Vector Machine | 120 |
| Table 6.4 | Classification report of Naive Bayes | 121 |
| Table 6.5 | Classification report of KNN | 121 |
| Table 6.6 | Classification report of Support Vector Machine | 122 |
| Table 6.7 | Performance Comparison | 122 |
| Table 6.8 | Performance of Naive Bayes Classifier | 124 |
| Table 6.9 | Performance of KNN Classifier | 125 |
| Table 6.10 | Performance of Support Vector Machine classifier | 125 |
| Table 6.11 | Classification report of Naive Bayes classifier | 126 |
| Table 6.12 | Classification report of KNN classifier | 126 |
| Table 6.13 | Classification report of Support Vector Machine classifier | 126 |
| Table 6.14 | Performance comparison data set 2 | 127 |

BIS (Hons) Information Systems Engineering
Faculty of Information and Communication Technology (Kampar Campus), UTAR

# LIST OF ABBREVIATIONS

| RS | Recommender System |
|---|---|
| ICT | Information Communication Technology |
| TRS | Travel Recommender System |
| Base MF | Base Matrix Factorization |
| Social MF | Social Matrix Factorization |
| Context MF | Context Matrix Factorization |
| CF | Collaborative Filtering |
| EFT | Explicit Feedback Technique |
| IFT | Implicit Feedback Technique |
| MCRS | Multi Criteria Recommender System |
| K-NN | K-Nearest Neighbor |
| SVM | Support Vector Machine |
| GPS | Global Positioning System |
| SNA | Social Network Analysis |
| URL | Uniform Resource Locator |
| RMSE | Root Mean Square Error |
| NLP | Natural language processing |
| SA | Sentiment Analysis |
| CTR | Click through rate |
| ML | Machine Learning |
| IICF | Item-item collaborative filtering |
| CBF | Crystallographic binary format |
| RFM-QA | Random forest based quality assessment |

BIS (Hons) Information Systems Engineering
Faculty of Information and Communication Technology (Kampar Campus), UTAR

| | |
|---|---|
| AI | Artificial Intelligence |
| IDF | Inverse Document Frequency |
| NLTK | Natural language toolkit |
| sklearn | Scikit Learn library |
| Libsupport | Library support |
| ASR | Automatic Speech Recognition |
| VAD | Voice Activity Detection |
| MLP | Multi Layer Perceptron |
| RF | Random Forest Classifier |
| MaxEnt | Maximum Entropy |
| SVC | Support Vector Classifier |
| CBR | Case based reasoning |

## CHAPTER 1 : INTRODUCTION

The research study concentrates on Personalized Recommender System utilizing web-based social networking. Recommender System (RS) has been viably used to deal with issue overwhelming. Relational associations, for instance, Facebook, Twitter are dealing with immense size of information by proposing user captivated things and things. RS has wide extent of uses, for instance, investigate articles, new social marks, films, music etc. The creating pervasiveness of casual networks assembles the availability of user estimations, which has become a basic impact factor on buying decisions, brand    and famous suppositions. In addition, prescribing proper reports, files, and users to seek after, has for a long while been a most cherished space for recommender structures research. A couple of new methodologies outfit constant small scale blogging exercises from administrations, for instance, Twitter , as the explanation behind distinctive user tendencies and isolating pertinent substance to explicit people. According to the user input and different quality things can be prescribed, which is immovably related to user interest. Review shows that more than 25 percent of offers made through proposition. Over 90% social orders acknowledge that things recommended by friend are useful and half people buy the endorsed things or things of their favorable position. Google+ introduced "Companions Circle" to channel the contacts as indicated by various exercises and techniques, which encourages users to be nearer to their companions. In a colossal web space, proposal finds things of user interest. Shared sifting and substance based separating are generally utilized strategies for suggestion.

Customized RS builds up factors, for instance, social interest, person's preferred position and social effect. Anil Rathod and Indiramma.M (2015) says customized RS is helpful to endorse the things on casual associations with the point that proposed things should reliant on their strength lead and social relationship of relational associations. Proposal in regular system centers around pair of (buyer, thing) while social recommendation bases on triplet (seller, buyer, thing) which improves the additionally fitting things of user intrigue. Research directed by Davide Feltoni Gurini, Fabio Gasparetti, Alessandro Micarelli and Giuseppe Sansonetti (2013) says

CHAPTER 1 : INTRODUCTION

Sentiment investigation or supposition mining is officially characterized as the computational investigation of conclusions and assessments about a substance communicated in a book. The substance is grouped into five classes: item, individual, brand, occasion, idea. Especially, right now accept the idea as the slant research study target element. Estimation research study is a troublesome errand, thus - before the arrangement of the calculation - a few suspicions are required.

Research study is conducted to study the performance of Personalised Recommender System algorithms. The machine learning algorithm covers Naive Bayes Classifier, K Nearest Neighbor Classifier and Support Vector Machine Classifier that involves artificial intelligence techniques. VADER Sentiment Analyser has also been implemented to analyse the sentiment polarity of Tweets in the data set. Thus by combining Lexicon Based approach and machine learning algorithm it forms a hybrid approach towards the classification of accuracy evaluations. Support Vector Machine Classifier depicts highest accuracy and Naive Bayes Classifier outperforming K Nearest Neighbor Classifier.

## 1.1 PROBLEM STATEMENT

Since social has an immense effect and critical as of now it is ideal to sift it through to get best and precise information. Many type of internet based life has been impacting our every day life. Actually one doesn't open news channels any longer to watch news. People hugely getting a wide range of information through internet based life in a matter of seconds.

Right now over-burden is all around watched. information Overload. In the exploration of Pietro Zanarini (2019) expressed the expression "information over-burden" . As per Gross, information over-burden is characterized as "information Overload happens when the measure of contribution to a framework surpasses its preparing limit. Chiefs have genuinely restricted intellectual handling limit. Thus, when information over-burden happens, all things considered, a decrease in choice quality will happen."

As the post and suggestion originates from a wide range of classes, the issue confronted is that trouble in recognizing content dependent on intrigue. As expressed there are huge amounts of facts accessible on the web however are they pertinent and dependent on the enthusiasm of the users. From the investigates done expressed that online life information can be valuable for recognizing new trends in the correspondence or issues which could include wild terrible exposure. Online networking is additionally utilized as a channel to speak with users. For supporting dynamic procedures, organizations utilize online networking reports, made ex post and dependent on predefined key execution pointers, or they utilize a scramble load up for jumping on-going investigations dependent on ongoing internet based life information.

CHAPTER 1 : INTRODUCTION

More often than not we acquire the correct information later or even not in any way. Rather than getting into the point the proposal recommends things that we don't require. Along these lines trouble in prescribing right thing at ideal time, setting or subject is confronted. Gigantic development of internet based life use has prompted an expanding collection of information, which has been named Social Media Big information. Social media platforms offer numerous prospects of information designs, including printed information, pictures, recordings, sounds, and locations. By and large, this information can be partitioned into unstructured information and organized information. In interpersonal organizations, the literary substance is a case of unstructured information, while the companion/adherent relationship is a case of organized information. internet based life information can be dissected to pick up bits of knowledge into patterns. Researches analysed Twitter information to concentrate how individuals' state of mind changes with time of day, weekday and season. collect a huge informational collection during a particular time period on a particular subject and examine it quantitatively. Along these lines distinguishing and suggesting right thing at perfect time, setting or subject has gotten problematic.

## 1.2 Project Scope

Right now Sentiment mining is one of the significant viewpoint toward recommender framework. information is mining significant where significant information can be mined dependent on the positive or negative points of view. Procedures will be taken care of to get most appropriate recommender by accessing the presentation if every strategy. Notion investigation techniques are utilized in the research study. informational index is gotten from Twitter as micro blogging has become the primary wellspring of news and proposal among the present ages. Notion investigation here states the utilization of normal language preparing, content research study and to acquire abstract information from Twitter.

The noteworthy viewpoint saw in the research study study is that accuracy of suggestion. Measurable strategies are utilized for assessment of the exhibition.

Factual strategies are numerical terms, models, and frameworks that are used measurable research study of unrefined research information. The utilization of factual investigation gets information from inquire about information and gives different ways to deal with study the intensity of research yields. Right now systems will be utilized to figure the exactness, accuracy ( of positive and negative). It will be the contrasted in the assessment with produce the best outcome. It is done in light of the fact that in Sentiment Analysis point of view a feeling which is dealt with positive can be negative in another announcement. The project implements qualitative analysis and quantitative analysis. Vader sentiment analyser has been implemented in analysing the polarity of sentiment of the tweets which is a lexicon based approach. In Analysing text classification algorithms Naive Bayes Classifier, K Nearest Neighbor Classfier and Support Vector Machine Classifier.

## 1.3 Project Objective

The objective aimed to be achieved conducting the research study is to condense data overload by retrieving the most suited data and services from large amount of data. The attempt of the current research study is to recommend the most suited and minimum set of exact information to the user. enormous measure of information, along these lines offering customized types of assistance. Review shows that in excess of 20 percent of deals produced through proposal. As indicated by the exploration of Himgauri D. Ambulkar , Apashabi Pathan (2013) the vast majority of the users accept that things suggested by companion or neighbors are helpful and not many purchase the prescribed items or things of their own advantage. Finding and breaking down user's past conduct, for example, appraisals of a thing it at that point predicts the things dependent on comparative sort of evaluations given by similarly invested users to the end user or focused on user.

The next purpose of the study is To find and analyse users' conveyed information using text classification algorithms that preducts higher accuracy and performance. The research is also done to find and analyse users' then predicts content based on similarity of content. Compares the content records and endeavors to prescribe things which are like that of user's favored things in the past.Users may have comparative interests yet have various assessments about them. Users tend to react same towards similar contents where their likes and dislikes depends on the content. Thus in this case, instead of recommending unwanted contents , recommending and predicting the right content according to the user preference is targeted to be achieved.

The research study is carried out to identify machine learning algorithm that recommend suitable or most similar content at a situation. "What item should be recommended in which time or context or subject?" by utilizing properties the thing arrangement, user inclinations, and proposal criteria is viewed as significant. Web-based social networking information can be valuable for identifying new trends in the correspondence or issues which could include wild awful exposure. Alan Hong (2019) says creating dynamically is challenge and when you achieved dynamic the

time you need to retrieve them is very important. Thus in this context to frequently update new data and identify user requirement to recommend suitable content is vital

## 1.4 Impact, Significance and Contribution

Research tries to date can be said to have basically centered around the investigation of the improvement of the suggestion models by using all the open information and user profiling information. In any case, barely any research studys have watched out for the issue of discovering which thing features, user interests and framework settings are the most vital exactly when precise and non-precise suggestions are delivered. In case those characteristics were distinguished, suggestion strategies might be reinforcing or disposing of their conditions with explicit speculations of user profiles and information things.

This assistance from our natural components gives us a basic technique to find the best choice without having a great deal of effort to isolating through the different choices available in the market. At this moment advancement, the Recommendation system is an application that isolated altered information and gives the best way to deal with understand a user's taste and to propose reasonable things to them by considering the models among their inclinations and research studys of various things.

This proposal frameworks recognizes the appropriate substance and things to utilize that may intrigue them as opposed to tossing everything into the substance of users. As states feeling investigation of tweets isn't easy because of multilingual and casual terms. Therefor this research study handles the issues by analysing the information got to deliver expected exact information.

## 1.5 Project Background Information

Recommender framework or a proposal framework is a stage to channel huge assortment of information that fall under the classifications of "rating" or "inclination" a user feels about a thing or a substance. It has a few sorts of approaches in particular community oriented separating, content based sifting, multi criteria recommender framework, chance mindful recommender framework, versatile recommender framework a half and half recommender framework. Every methodologies has own few different ways of conduction the approaches and capacities.

Ms.A.M.Abirami and Ms.V.Gayathri (2016) claims that Sentiment research study (or) assessment mining assumes a noteworthy job in our day by day dynamic procedure. These choices may Sentiment order can be acted in 3 phases, for example, archive level, sentence level and highlight level. Archive and sentence level it utilizes just single article and concentrates just a solitary sentiment. In any case, such suppositions are not fitting for certain circumstances. Despite the way that the fundamental establishments of recommender systems can be followed back to the broad work in the subjective science, guess hypothesis, data recovery, speculations, and furthermore have connections to the board science, and furthermore to the shopper decision displaying in promoting, recommender frameworks developed as a free research territory in the mid 1990's when analysts began concentrating on suggestion issues that expressly depend on the appraisals structure. The procedure that employments rating structure is named as "Shared Filtering" and it was presented by Goldberg et al (1992) with regards to first business recommender framework Tapestry which was intended to prescribe archives attracted from research studys to a assortment of users. One significant part of recommender framework is personalization that makes it conceivable to give customized suggestion to user and it is generally upheld by content data. In content based sifting, things or a portion of the primary substance based proposal frameworks and noticed that they utilized innovation identified with data recovery.

CHAPTER 1 : INTRODUCTION

Humans would have always wondered how the "people you may know" features on social media has been appearing every time we log in. Such feature suggest a list of people user possibly know, who are similar to your details based on your friends or followers, current location, groups, liked page or posts. These recommendation might be different to you and your friends based on the individuals' interests. You may see quotes or poems that is totally different from your friends where it is all songs recommended on the feed. It is due to the search and liked posts of the users. In some terms it is also knows data analytics. data Analytics is a craft of handling crude data to extricate some sensible data. data Analytics is broadly utilized in numerous ventures and association to settle on a superior Business choice. By applying investigation to the organized and unstructured data the ventures gets an extraordinary change their method for arranging and dynamic. data research study is the way toward checking, cleaning, and changing so as to recover valuable data from the data. This data will be more

Supportive in recommending business ends and decisions making. data Analysis has an assortment of points and strategies that consolidates numerous methods so as to give better exactness. One of the most mainstream strategies for data investigation procedure is data mining that fundamentally focuses on demonstrating and disclosure of data for expectation process as opposed to clear purposes. prescient investigation is basically utilized for foreseeing estimating/characterization where as content research study utilize measurable, phonetic and auxiliary strategies so as to recover data from content sources. This content sources are generally as unstructured data. In short recommender system in explained as shown and further explained in coming chapters in detail.

Figure 1.1 Personalised Recommender System

## 1.6   Highlight of what have been achieved

The project has achieved the proposed objective to overcome the current issue of data overload.   The steps in the algorithm involved preprocessing steps which was to handle the data for a smooth sentiment analysis and text classsification. The preprocessing steps involves remove pattern, punctuation and characters. The tweets will be then tokenized, stemming of the text is conducted followed by removing the stop words. The tweets are the rejoined to go through the process sentiment analysis and text classification. VADER sentiment analyser which is considered a suitable analyser for social media text has been imported to analyse sentiment of the tweets. Naive Bayes Classifier , K Nearest Neighbor Classifier and Support Vector Machine classifier is the implemented to analyse the accuracy of classification methods.

BIS (Hons) Information Systems Engineering
Faculty of Information and Communication Technology (Kampar Campus), UTAR

## 1.7 Report Organization

The Next Section which is Chapter 2 will be explained on literature reviews that has been gathered. It has comparison of the research study with previous analysis or studies that has been conducted. Chapter 3 explains the system design that has been used by the author. Why those methodologies has been chosen and how they function to measure accuracy. Flowchart of the algorithms has been introduced in chapter 3. Chapter 4 conducts the system implementation that has been done and analysis to measure polarity of the text in the data set that has been imported. The machine learning algorithms has also been evaluated and implemented. Chapter 5 covers system testing, another data set is implied in the algorithm to observe similar output. Chapter 6 states the results obtained and discussions are conducted. Chapter 7 concludes the project.

## CHAPTER 2 : LITERATURE REVIEW

## 2.1 Litereature Reviews

The literature review conducted on 45 research studies. The research study litrature reviews were done on Personalised Recommender System, Machine Learning Algorithm, Artificial Intelligence Techniques, Sentiment Analysis of Twitter, Naive Bayes Classifier, K Nearest Algorithm, Support Vector Machine classifier, Content based Recommender, Collaborative filtering techniques and hybrid methodologies.

## 2.1.1 Literature Review Recommender system and sentiment Analysis

**Fabio,et al (2015)** The research study explored how the explosion in online life requires the utilization of social media analytic. Explains the fundamental phases of the social media process and depict the most widely recognized social media systematic strategies being used; and we talk about the manners by which social media analytics make business esteem. The reason for the research study directed was to deal with the circumstance worldwide, multilingual online life checking and investigate organization, to analyze the in excess of 5,000 user feelings that are posted about different brands every month on travel sites.

The issue statement that must be recognized strain to build consumer loyalty and quality of service in the midst of monetary downturn. The main challenge was to having the option to rapidly recognize user disappointment and afterward right the issue at their source. The target expressed for the exploration was to collect, monitor analyze, summarize and imagine social media information and to encourage discussions and cooperations to separate helpful examples and insight. In this way to expedite a triumph the target a device was made extraordinarily intended to follow the online notoriety of 12,000 lodgings, including its competitors. Social media research study includes a three-arrange process: catch, comprehend, and present.

Subsequently it immediately uncovered various issues that users were encountering. In determination a set up of rewards and preparing program that urged individual lodgings to interface with users through online discussions and simpler to in a flash and extensively share encounters all consolidate to make a web based life scene that is quickly developing and turning out to be perpetually part of the texture of organizations.

**Fan, W. and Gordon, MD., (2014)** The research study research study quantifies new proposal that consider the literary fragment of user audits. The exploration is about a played out an all around request of a genuine true eatery survey informational index and report the systems and revelations. To grasp the condition by applying content assessment to a proposition circumstance and show that the more low down printed information can improve rating desire quality.

The issue distinguished was sites giving customer audits are incredibly mechanically poor, customers routinely should pick the alternative to peruse through immense proportions of substance to find a particular piece of fascinating information. Getting to and glancing through substance studies is baffling when customers simply have a dark idea of the thing or its features and they need a proposal or closest organize. Catchphrase look normally don't give incredible results, as comparative watchwords routinely appear in extraordinary and in horrendous reviews. In this way, understanding audits is that a commentator's general rating may be to a great extent intelligent of item includes in which the hunt user isn't intrigued. By and by, the issue expressed is recognizing organized data from freestyle content is a difficult undertaking as users routinely enter casual content with poor spelling and punctuation.

The target must be accomplished was to develop procedures to group and break down content and structure-based web surveys, and utilize the subsequent research study to improve customized proposals for web user just as to recognize prevalent order subjects and conclusions. Kappa coefficient (K) used to gauge pairwise understanding among a lot of annotators making class decisions, remedying for anticipated

possibility understanding. A Kappa estimation of 1 infers immaculate understanding, the lower the worth, the lower the understanding performed 7-overlap cross approval and utilized exactness, accuracy and review to assess the nature of our characterization

To use the conclusion data into a solitary score for each survey, made an interpretation of commented on content audits into a book rating score that can without much of a stretch be contrasted with the meta data star rating. Content rating recipe of [ P/P + N] ∗ 4 + (1) utilized where P is the quantity of Positive sentences in the survey, and N is the quantity of Negative sentences. Another strategy was by regression loads to figure content rating scores can bring about scores that lay outside of the [1:5] territory.

**Kywe, S.M., Lim, E.P. and Zhu, F., (2012)** The explanation depicted was to develop a clever consistent constant user explicit travel recommender structure (IRTUSTRS) through solidifying customers' relational association profile and current territory by misusing worldwide situating framework (GPS) data for development proposal age. The issue perceived was the development proposition issue ceaselessly circumstance.

The direct of research was to communicate the possibility of areas based frameworks includes three sorts of charts, for instance, area diagrams, user area charts, user charts. Three social affairs of past region based long range casual correspondence organizations are Geo-labeled media-based organizations, Point-area based organizations and Trajectory-based administrations. Recommender structure generally gives things proposition as an overview, according to the interests of the customer. The investigation needs to predicts the customer's response for each proposed things as recommender systems help customers in the decision method to pick the particular thing and make the strategy easier.

The objective or augmentation intended to be practiced is to develop a wise constant user explicit travel recommender framework (IRTUSTRS) for customized suggestions. To misused the customers' interpersonal organization profile and GPS data for the time of development proposals. To probably evaluate made IRTUSTRS through two

15

constant enormous scope data sets of Yelp and Trip Advisor. Thus the framework incorporates proposition part of the extraction of appropriate information from the relational association information and relating update of the way of thinking. It is named as theory masses, and incorporates the extraction and request of the various associations, thoughts and their events as indicated by the definitions made in the mysticism.

The results saw was the procedure for anticipating customer express travel proposal is broadly apportioned into two phases. In the underlying advance RS comprehends the customer direct through Social Network Analysis, SNA data and research study it with the customer tendencies. The future work moreover joins the customized travel suggestion through abusing multi-operator structures to assemble customer data from various sources.

**Berka, T. and Plößnig, M., (2004)** The point of the research study was to give a knowledge into the urgent choice about prescribing approach, calculation and engineering, we will give a concise outline of the four essential methodologies for recommender frameworks. These four essential approaches vary in the strategies utilized, yet additionally in "this present reality understanding" of the fundamental calculation. The issue proclamation has been expressed as pursue, items and administrations in the field of the travel industry (like lodgings, bundles, and so on.) are basically not physical and normally exists generally as data. Therefore, they are predestinate for electronic deal. ICT permits effectively to give the travel industry contributions more extravagant portrayals to empower voyagers to settle on increasingly educated decisions. As outcome, the multifaceted nature of item depictions is developing.

The venture was controlled by recommender frameworks a famous research territory and progressively utilized for eCommerce. There the object of research was to copy disconnected travel specialists by giving users knowledgeable travel proposals to encourage their basic leadership procedures and utilizing one separating method can miss the mark when attempting to make suggestions for complex items. To be certain that solitary a mix of various separating procedures can offer noteworthy enhancements for the basic leadership process.For accomplishing degree, the strategy

BIS (Hons) Information Systems Engineering
Faculty of Information and Communication Technology (Kampar Campus), UTAR

conveyed was by database predetermination gives an underlying choice of things, in light of basic database cooperations (like guideline based sifting utilizing boolean rationale). This yields an extremely proficient decrease of the quantity of things, which must be processed, at a beginning period of the work flow.

Knowledge based filtering strategy was handled with. The utilization of information based sifting enables engineers to utilize express area information.Collaborative filtering, content-based filtering strategies, or for all intents and purposes some other type of thing rating method to acquire at least one numerical evaluations for each thing. This permits to join algorithmic or certain area information. As to picture execution a sketch of work stream for the plan procedure of a travel industry recommender framework was directed by first deciding data, apply monetary criteria. Following determine separating systems and adjust point by point AI structure and execution.

After beginning investigation, look into had inferred that can utilize the accompanying AI methods in a TRS: Database preselection - Quick decrease of a lot of information. Content-based scoring that grants to use typical vacationer media information to incorporate scores based substance and its semantics to the general deduction process. Consolidating irrefutable region data that grants to organize deliberately amassed space data and expertise.

**Rathod, A. and Indiramma, M., (2015)** The inspiration driving the investigation is that building one thing is simply not adequate any more. Associations ought to have the alternative to, at any rate, develop various things that address the different issues of various customers. The improvement toward E-business has empowered associations to give customers more options. In any case, in stretching out to this new level of customization, business increase the proportion of information that customer must process before they have the alternative to pick which thing tends to their issue.

BIS (Hons) Information Systems Engineering
Faculty of Information and Communication Technology (Kampar Campus), UTAR

The objective was communicated as to give a ton of recommender structure models that length the extent of different employments of recommender systems in E-business. Second, to research the way all of the models uses the recommender structure to overhaul pay on the site. Third, to depict a mapping from employments of recommender systems to a logical classification of strategies for executing the applications. Fourth, to take a gander at the effort required from customers to find recommendations. Fifth, to delineate a great deal of proposals for new recommender system applications reliant on parts of logical order that have not been explored by the present applications.

The method used was non-customized recommender structures since they require little customer effort to make the proposition, and are Ephemeral, in light of the fact that the system doesn't see the customer beginning with one meeting then onto the following since the recommendations are not established on the customer.Attribute based recommender structures endorse things to customers reliant on syntactic properties of the things.Individuals to individuals relationship recommender systems prescribe things to a customer subject to the association between that customer and various customers who have gained things from the E-business site. To finish up the research study it was resolved that recommender frameworks are a key method to robotize mass customization for E-business destinations.

**Interaction Design Foundation (2019)** The outline is driven since Twitter customers make tweets each day, these customers are also overwhelmed by the huge proportion of information open and the epic number of people they associate with. To overcome the above information over-trouble issue, recommender structure can be familiar with help customers make the best possible assurance. There is no complete diagram for the area of suggestion in Twitter to arrange the existing fills in similarly as to recognize regions that need to also inspected. The research study right now to full the opening by introducing a logical order of recommendation assignments in Twitter, and to use the logical characterization to portray the appropriate works of late.

BIS (Hons) Information Systems Engineering
Faculty of Information and Communication Technology (Kampar Campus), UTAR

The customized recommender systems use characteristics of things, profiles of customers and the connection or trades among customers and things to predict the customers' future thing determinations. Synergistic sifting and substance based approachs are every now and again used technique in modified proposition. The fundamental supposition of the user to user based network arranged isolating methodology is that if an individual X has a comparative end as an individual Y on an issue A, X will undoubtedly get Y 's appraisal on an other issue B than an arbitrarily picked person. The recommender system finds people with tantamount tastes or tendencies, as showed by their past assessments or certain correspondences. By then, the system predicts the tendency of a customer on an unrated thing using the tendencies of similar customers. Another modified proposal approach is thing to-thing network situated isolating which is used by Amazon.com's recommender structure. Things An and B are significantly similar if a tolerably colossal piece of the customers who purchase thing A furthermore buy thing B. By then, the tendency of a customer over an unrated thing B is foreseen subject to the customer's assessed thing

There are four systems used specifically topology based methodology, basic procedure, weighted substance based strategy and twittomender. With respect to the ebb and flow explore weighted substance based strategy was helpful as it utilizes five highlights, in particular notoriety, action, area, companions in like manner and substance of the tweets, are anticipated to be applicable for proposal, just fame and action have been assessed. The impulse of the research study is that if a target customer has various notable and dynamic followees, other predominant and dynamic followees should be endorsed to the customer. If the target customer has quite recently notable followees, simply celebrated followees should be recommended.

Utilizing the scientific classification, the research study have over viewed a few suggestion techniques exceptionally produced for Twitter. To the best of information, this is the first run through a scientific categorization is utilized to arrange suggestion errands in Twitter. As future proposal the present hash label suggestion frameworks

just consider the substance of tweets yet not user inclinations or viability of hash labels in spreading data. There are likewise not many takes a shot at notice or retweet suggestion. At the point when answers for these proposal errands are created and assessed with high correctness, one can conceive a progressively thorough scope of suggestions customizing the utilization of Twitter.

**Prerana Khurana , Shabnam Parveen (2016)** The made research study intends to convey the noteworthiness of Recommendation Sytems, different approaches and social segments, which effect Personalized Recommendation System. The maker states community oriented sifting and substance based separating are commonly used procedures for proposition. Notwithstanding the way that there are various calculations to work away at Data Mining, cold beginning has made people to step back in breaking down the convenience of those computations lead to little decrease in creativity and enhancements in data mining estimations. In any case, chilly beginning can be depicted as detachment of data for showing computations. Three sorts of cold starting issues are new customer issue, new thing issue and new structure issue.

Sparsity issue is one of the troublesome issues experienced by recommender structure is information sparsity has exceptional effect on the possibility of proposal. The fundamental job for information sparsity is that most users don't rate the greater part of the things and the open research studies are normally insufficient. As to fathom one can utilize customer profile data while enlisting user similarity thing with others. Flexibility is the property of framework demonstrates its ability to oversee making extent of data in an in the present style manner.Approximation thinking is endorsed to pound versatility. Overspecialization. happens if framework underwrites just things which having high surveying against the user profile. Resemblance measures, inspecting and estimation decrease in synergistic filtering vanquishes the overspecialization issue.

BIS (Hons) Information Systems Engineering
Faculty of Information and Communication Technology (Kampar Campus), UTAR

CHAPTER 2 : LITERATURE REVIEW

The idea of altered recommender framework is that it is useful to propose the things on easygoing relationship with the point that embraced things should dependent on their chronicled lead and social relationship of social affiliations.

Recommendation in customary structure rotates around pair of (purchaser, thing) anyway social suggestion bases on triplet (dealer, purchaser, thing) which upgrades the more suitable things of user interest. The method managed for General Recommendation most remarkable frameworks are content based filtering and thing based confining. Both of these frameworks are defenseless against cold start and sparcity issue. To defeat these issues adjusted proposal structure utilizes the social interest, social profile, and so forth to endorse customer fascinated things. Essential Matrix Factorization the task of RS is to lessen the screw up of foreseen a motivation to the veritable rating regard. Subsequently, the BaseMF model is set up on the watched rating data by constraining the objective work. This relies upon the probabilistic framework factorization, which uses the low position arrange.

The CircleCon model has been found to beat BaseMF and SocialMF concerning precision of the RS. In the first place, trust circle determination, this can be identified with the guide of different characterization (classes) with certain breaking point regard. Second, trust regard task, this should be conceivable with proportional trust, try based trust and trust separating. CMF Model to endorse good things from sender to beneficiary. Neighbors can be recognized by Bayesian inference which helps with perceive customer

**Abdel-Hafez, A. and Xu, Y., (2013)** The diagram intends to inspect the available customer showing methodologies for web based life locales, and to highlight the inadequacy and nature of these strategies and to give a fantasy to future work in customer showing in internet based life destinations. Critical issues in personalisation is building customers' profiles, which depend upon various segments, for instance, the pre-owned information, the application space they plan to serve, the depiction procedure and the development system. Online networking is groups into social networking, collaborative work, content sharing, blog, rating and audit, social

bookmarking, virtual world and virtual game world. Issues has been tended to was the dynamicity issue. To accomplish this objective viably, particularly with the assorted variety of nature of the accessible internet based life sites; one research recommended to regard them as an environment of related components when you build up a web based life methodology as opposed to regarding them as independent frameworks.

The initial method is the information assortment, which assembles users' information from social media sites remembering filled for structures information, log document information, and associations with others in the framework. The subsequent stage is the profile development, where the users' advantages will be removed and spoke to utilizing various strategies; loads additionally will be implanted with each enthusiasm demonstrating the level of intrigue. .

Targets to add progressively related watchwords to the profile so as to improve the last forecast outcomes; numerous sources can be utilized so as to extricate the additional catchphrases, for example, WordNet synsets, Web sites and similarly invested users or companions users.
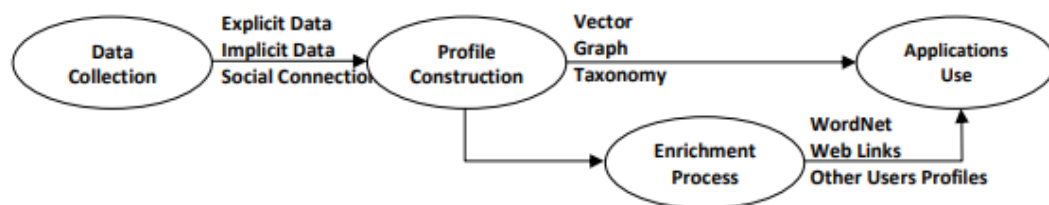


Figure 2.1 Cleaning stages

Future of this methodology is concentrating on customary profiling procedures without considering the assorted variety of components gave by various online networking site.

CHAPTER 2 : LITERATURE REVIEW

**Seth, A. and Zhang, J., (2008)** The issue distinguished is to sort out and deal with the a huge number of approaches in recommending music. Music Information Retrieval Technique (MIR) been created to tackle issues, for example, kind characterization, artist classicification and insturement acknowledgment.

The advantage position is that music recommender is to assist users with separating and find melodies as indicated by their preferences. A decent music recommender framework ought to have the option to naturally recognize playlist and produce play lists appropriately. Then, the improvement of recommender frameworks gives an extraordinary chance to industry to total the users who are keen on music. All the more critically, it raises difficulties for us to more readily comprehend and show users' inclinations in music.

The shortcoming saw in MIR is it is quick and precise, the downsides are obvious.First of all, the user needs to think about the article data for a specific music thing. Besides, it is additionally tedious to keep up the expanding metadata. Besides, the proposition results is commonly poor, since it can simply recommend music reliant on distribution metadata and none of the customers' information has been thought of. Perhaps the best methodologies in proposition systems, it anticipate that if customer X and Y rate n things moreover or have similar direct, they will rate or follow up on various things nearly.

Content-based approach makes figure by separating the tune tracks. It is built up in information recuperation and data sifting that suggests a song which resembles those the customer has checked out in the past rather than what the customer have evaluated 'like'. Request by Humming (QBSH) Humming and singing are the regular strategy to communicate the tunes steps: improvement of the tunes database, translation of the customers' melodic information question and model planning computations which are used to get the closest outcomes from arrangements. Feeling based Model Music feeling has propelled lots of research and it has become the essential example for music revelation and proposition. A business web organization called 'Musicovery' uses the basic inclination model (2D valence-fervor) found by clinicians. It empowers customers to locate their typical saw feeling in a 2D space: valence (how positive or negative) and fervor (how stimulating or calming).

BIS (Hons) Information Systems Engineering
Faculty of Information and Communication Technology (Kampar Campus), UTAR

Customer Listening Experience Modeling Depending on the level of music capacity, their wants in music are varied suitably. Jennings inspected the different sorts of crowd individuals whose age reach out from 16-45 and characterized the crowd individuals into four social events:  fans, casuals, indifferent. Thing Profiling – Editorial metadata: Metadata got by a singular ace or assembling of experts. Social metadata: Metadata procured from the research study of corpora of artistic information, generally from the Internet or other open sources. Acoustic meta data: Metadata got from an assessment of the sound sign. This should be with no reference to a scholarly or supported information. For instance Beat, musicality, pitch, instrument, outlook, etc.

Anyway this exploration is in a beginning period. As should be obvious from the advancement of music recommenders over the previous years, the given outcomes will in general be increasingly customized and emotional. Just thinking about the music itself and human appraisals are not, at this point adequate. A lot of work as of late have been done in music discernment, brain science, neuroscience and game which study the connection among music and the effect of human behaviour.Main benefits for tuning in to the music which incorporate work yield augmentation, execution upgrade, and separation from upsetting sentiments and so forth. Users' inclination in music is connected to their personality.Future music recommender ought to have the option to lead the users sensibly pick music. As far as possible, we are trusting that through this investigation we can construct the scaffold among detached research in the various controls.

 **Mrs.M.Sridevi, Dr.R.Rajeswara Rao. Et al (2016)** Recommender systems are significant decision to glance through counts as it help customers with discovering things they most likely won't have discovered free from any other individual. Abnormally, recommender structures are routinely realized using web search device requesting nontraditional data. According to this research study recommender system

is apportioned into three groupings synergistic filtering(CF), content based sifting and half and half models (CF + content based separating).

Synergistic separating assembles a model from the user's past conduct (things as of late got or picked or conceivably numerical rating given to those things) and practically identical decisions made by various customers. This model is then used to anticipate things (or assessments for things) that customer may have an interest in.This approach experience the evil impacts of three issues: cold beginning, versatility, and sparsity. Content-based sifting approach utilize a movement of discrete characteristics of a thing in order to prescribe additional things to customers with similar properties. A segment recommender system give proposition subject to the measurement profile of customer. Suggested things can be made for different measurement claims to fame, by consolidating appraisals of customers in those specialties.

Information based: knowledge set up recommender proposes things based regarding inductions about customer's needs and likes. This data will to a great extent contain thing incorporates that address customer issues. Half breed recommender: Hybrid recommender structure is the one that consolidates different proposal procedures to make the yield. There is no inspiration driving why some different strategies of a comparative kind couldn't be hybridized, for example, two assorted substance based recommender can coordinate, NewsDude, which uses both credulous Bayes and kNN classifier in its news recommendations is just one instance of it. In future the study will be continued to chip away at half breed recommender utilizing grouping and similitude for better execution.

**Ganu, G. (2009)** The reason for explore is to recognize applicable individuals to pursue among a huge number of users that connect. Spotlights on measurements for estimating the theme closeness among Twitter users, the last endeavors the chart of connections among users to surmise relationships. The fundamental thought behind this work is that users may have comparable interests however have various sentiments about them. Along these lines, the exploration expands the substance based proposal by methods for the opinions and sentiments extracted from the user smaller scale presents all together on improve the exactness of the recommendations.

25

This prompts characterize a novel weighting capacity so as to improve content-based user profiles.

The issue discovered is that there are no endeavors towards supposition user suggestion in informal organizations. As needs be the objective perceived estimation assessment structure is to procure a yield esteem that addresses how a great deal of positive, negative or unprejudiced is the supposition conveyed in a tweet. In this manner, ask about completed a Supervised Machine Learning count subject to a Naıve Bayes classifier. With the ultimate objective of setting up our computation, required an enlightening assortment with stamped tweets.

Subject-activity word object (SVO) weighting limit thought behind this work is that considering customer mindsets towards his own one of a kind focal points can yield benefits in recommending allies to seek after. Specifically, Authors consider  which is the presumption imparted by the customer for a given thought, the sum he is enthusiastic about that thought, and (iii) the sum he conveys target comments on it.

This technique empowered makers to collect more finish customer profiles than traditional substance based systems. Preliminary results show the upsides of proposed model differentiated and some top tier procedures.

**Schafer, J.B., Konstan, J. (2012)** A recommender framework to be great in the event that it can deliver helpful suggestions, yet in addition clarify qualities of the proposals created by it as far as components watched and examined by media scholars. Authors attempt to do as such in this research study, to address addresses, for example, given a lot of messages about a recent development that have just been browsed by a user, what least arrangement of extra messages ought to be prescribed to the user, with the goal that she gets basic yet assorted data the occasion?  intend to suggest messages about the story that the user will see as generally valuable as far as how much rearrangements and conclusion decent variety the messages will give to the user. Explicitly center around recommender frameworks for short notes.

The goal of the exploration done was to demonstrate the highlights that give simplification and assorted variety in a novel way, in view of the hidden interpersonal organization of the message writers, members, and browsers. Data about the informal

community of users was not accessible as of not long ago with the pattern towards interpersonal interaction sites and the extraction of informal organizations from email graphs. Authors utilized the informal community based highlights that give simplification and decent variety to build up a customized user model spoke to as a Bayesian system, which demonstrates the preferences of its user towards these highlights.media content.

$$sim(u_i, u_j) = sim(p(u_i), p(u_j)) =$$
$$= \frac{\sum_{c \in C_{u_i} \cup C_{u_j}} \omega(u_i, c) \cdot \omega(u_j, c)}{\sqrt{\sum_{c \in C_{u_i}} \omega(u_i, c)^2} \cdot \sqrt{\sum_{c \in C_{u_j}} \omega(u_j, c)^2}}$$

Figure 2.2 : Formula of algorithm

proposed and assessed a way to deal with customized suggestion of participatory media content utilizing interpersonal organizations and a Bayesian user model.

**Ambulkar, H. and Pathan, A., (2015)** Review coordinated because in making recommender systems is to thick information over-trouble by recuperating the most fit information and administrations from a ton of data, in this way giving tweaked organizations.

Methods subject to how customer profile information is used. In content-based, the recommender will recoup things whose substance is match to those of the profile. Network situated filtering relies upon appraisals of various customers. User-based balance content records with customer profiles. Model based (Item-based) use data mining strategies to develop a model of customer evaluations. Customized RS is useful to suggest the things on interpersonal organizations with the point that prescribed things should reliant on their past direct and social relationship of informal organizations. Propose a better than average approach to manage improve the precision of recommendation structure by introducing "found companion systems". Relational enthusiasm of customer in casual association improves the idea of the proposition.

Cold beginning issues have three sorts: new user issue, new thing issue and new structure problem.Interpersonal relations, for instance, "companion circle" and

individual inclinations are critical factor in casual association , which deals with the issue of cold beginning. Sparsity issue is seen with a huge course of action of things. It is difficult to find the model if there ought to be an event of new system, due to availability of less information about customer and as of late entered thing or things. Versatility issue in recommendation is caused in light of the extending number of customers, pack customers into bunches reliant on their inclinations toward relative data articles, for instance, films, jokes, etc. To handle flexibility issues rely upon estimation part.

Cooperative separating, gathering and pre-getting ready are used to deal with versatility issue. Loud information is cleaned and changed to the recommendation system in preprocessing. Sparsity issue is one of the difficult issues experienced by recommender system is data sparsity has remarkable impact on the idea of proposition. The essential clarification for data sparsity is that most customers don't rate most of the things and the available appraisals are ordinarily pitiful. One can use customer profile information while discovering customer likeness thing with others.

The systems dealt with to deal with the issues was Basic Matrix factorization is gotten together with interpersonal organization data in a proposition structure where, the probabilistic network factorization doesn't consider any social factor. Circle Con model has been found to beat BaseMF and SocialMF concerning precision of the RS. ContextMF model, since it more straightforward for the recommended things of the model to be changed into purchase rate than the got things in Facebook style informal communities. Setting Matrix Factorization recognizes near eagerness by means of planning objective work. Direct neighbors can be recognized by Bayesian finding which perceives customer singular interest clearly related to appraised things.

**Debashis Das (2017)** This investigation research study contemplates and subtleties the diverse kind of recommender structure and acclaimed suggestion calculations and its livelihoods. Suggestion framework is arranged into customized and non-customized. Customized also requested into content based, synergistic separating, segment, information based and half breed. Collective separating has portrayal known as customer based and thing based.

Content based filtering worked with respect to the user likes and portrayal about the thing. The advantage of this procedure is that other user's information not required.No information sparsity just as cool start. The inconvenience is that substance research study is fundamental to characterize the item features. The greatness of the item can't be assessed. The likeness calculation is deficient to the item depiction. Collaborative filtering fabricate the framework by thinking about the user's past conduct. The greatness of the thing can be assessed through user evaluations anyway chilly start issue for various users and new items. Stability versus versatility issue. The procedure is area free since Item include isn't required. Collection of demographic data offer ascent to protection issue. Versatility versus strength issue. Information based calculation thinks about the information about the thing and its component, user inclination (asked expressly), proposal criteria before giving the recommendation.Hybrid worked by joining diverse recommender frameworks to construct an increasingly hearty system.

There are a few of suggestion framework input methods expressed. Implicit Feedback technique(IFT) is acquired without user's awareness yet got dependent on the user's activity during the procedure. IFT can be gathered at a lot of lower cost. IFT is unproblematic, it doesn't put a heap on the user of the recuperation system.It is less exact as contrast with EFT however enormous data can be gathered at a lower cost. IFT is powerless against clamor. IFT is less exact contrasted with EFT. Hard to decipher. Explicit Feedback Technique includes the users for allocating either numeric or score rating for assessing the item.

Ratings given on these measures license these decisions to be taken care of measurably to give circulations, normal, etc. EFT is easy to utilize. The exactness of EFT is by all accounts higher than IFT. EFT is supreme. Nosiness is one of the difficulties influencing EFT. Utilizing numerical scale can be confounding as the user probably won't be steady in giving their rating. User's appraising probably won't show the genuine assessment of users. EFT is defenseless to clamor. It is touchy to user context.Hybrid Feedback Technique is mix of both IFT and EFT. HFT recuperates the

estimate rating precision. HFT is the blend of both Implicit and express techniques. HFT isn't modest. It is computationally concentrated.

In future, different qualities and strategies can be created and assessed for proficient usage of proposal frameworks.

**Logesh, R et al (2018)** Recommender system are predominantly classified into two kinds Content-Based and Collaborative Filtering. Content-Based frameworks; which recommends quantity of things like past understanding of the user,who enjoyed previously. Collaborative filtering recommender frameworks depend on the idea of closeness of users.

There are a couple of philosophies used in recommender framework. Communitarian Filtering Technique major supposition that can't avoid being that the people who have inclinations in the past may moreover prefer to have same kind of tendencies later on. Content-based Filtering has limited substance assessment. Overspecialization limits customers to things like the ones portrayed in their specific profiles and right now things and various decisions are not found. Versatile Recommender Systems are used. Hazard Aware Recommender Systems existing recommender structures focused on suggesting things/administrations to customers by not considering into a record of upsetting the customer in a specific time, day and situation. Regardless, in various progressing applications it isn't only critical to prescribe customized content yet moreover consider to consider the threat of pushing proposals in explicit conditions and upsetting the usage. Thus, the degree of hazard is genuinely influence on the introduction of the structure during the time spent proposition. Multi-Criteria Recommender Systems (MCRS) solidifies tendencies reliant on various worldview.

The general execution of recommendations is improved by joining distinctive establishment into the Recommender structures. A couple of benefactors MCRS considered as a Multi-criteria Decision Making (MCDM) issue, and apply MCDM strategies to execute MCRS systems. Setting Aware Recommender Systems are setting pre-filtering, which channel the customer likes data before applying suggestion calculation as demonstrated by the objective suggestion objective, setting post

30

isolating which apply figuring on customer tendency data and a short time later change evaluations according to proposition setting. Social occasion Recommender Systems a get-together of individuals having with tantamount taste or direct. There are various issues related to the arrangement of Group Recommender systems which recall sparsity and vulnerability of data for customer profile, exhibiting bunch level likes and evaluating bunch proposals. Master Based Recommender System - addresses an issue of recognizing pros in the important field of information the executives is the most reassuring undertaking. Master is one who has more information in their specific region than an ordinary customer.

Non–Independent and Identically Distributed Recommended Systems. Chart Based Recommended Systems builds business applications that have monstrous measure of information with a higher request level of connection, another arrangement.. Data Based Recommender Systems. Cross breed based Recommender framework incorporates   strategies including weighted, course, meta level.

## 2.1.2 Literature Review on Hybrid classification

**Khan, F.H, et al (2014)** The research study study directed on account of nature of miniaturized scale sites scale composes on which people post constant messages about their opinions on a variety of focuses, talk about current issues, gripe, and express constructive estimation for things they use in consistently life. The goal of this research study study can be condensed as follows: Introduces and completes a cross variety approach for choosing the slant of eachtweet. Shows the estimation of pre-having data using acknowledgment and assessment of slangs/constrictions, lemmatization, amendment and stop words evacuation. Tests the exactness of slant identifification Resolves the data sparsity issue using space free methods. Correlation with various techniques to exhibit the ampleness of the proposed half breed strategy

Prehandling steps include: evacuation of URLs, hash-labels, username and remarkable characters; spelling amendment using a word reference; truncations and slangs with expansions, lemmatization and stop words expulsion. The proposed arrangement figuring wires a creamer plan using an improved kind of emoji research study, SentiWordNet investigation and an improved extremity classifier using summary of positive/negative words. Author have moreover discussed troubles that are looked and proposed the computation that settle these issues and additions the grouping exactness effectively decreasing the amount of ordered neutrals.

## 2.1.3 Literature Review on The Text Classification Approaches handled

**Abirami, A.M. and Gayathri, V., (2017)** The research study states the purpose as sentiment Analysis utilizes 3 terms so as to bring the sentiment .That is item and highlight, conclusion holder, accuracy what's more, direction. Feeling Analysis manages a few specialized difficulties, for example, object recognizable proof, conclusion direction arrangement, and highlight extraction. Ordinarily sentiment investigation can be performed utilizing directed and unaided learning. Among procedures SVM is viewed as increasingly reasonable for    Analysis.

Extremity Shift is a most huge issue to be tended to in Sentiment Analysis. extremity Shift suggests that Polarity (finish) of the sentence is resolved in different route From the extremity truly imparted in the Sentence paired. Research study proposed technique to improve the exactness of precision assessment in the twitter Post . Proposed method outfit and author for oversee Forced character limit problem.Most of the web based life For instance, twitter are confined to certain character limit which Prompts abnormality in the outflow of sentence. The proposed framework was made to manufacture the Precision of feeling mining by brushing customary machine

SVM - Support Vector Machine alongside the Fuzzy area cosmology. Conventional AI computation sets up obstruction arranged by sentiment.That is overview features will be requested unmistakably to Restricted class, for instance, Positive/Negative. The proposed system brings a sensible understanding that The feathery based cosmology is more space unequivocal as differentiated and conventional crisp philosophy.

**Dey, et al (2016)**   Difficulties saw in the research study is Sentiment Analysis are an estimation word which is treated as positive side may be considered as negative in another situation. Also the degree of motivation or cynicism in like manner incredibly affects the notions. For example "incredible" and "by and large amazing" can't be managed same. The system expressed was using two Supervised Machine Learning computations : Naive Bayes' and K-Nearest Neighbor to calculate the correctness, precision (of positive and negative corpus) and survey estimations (of positive and negative corpus).

K-NN is a kind of case based learning, or lazy acknowledging where the limit is simply approximated locally and all estimation is yielded until game plan. It is non parametric procedure used for game plan or backslide. On the off chance that there ought to emerge an event of collection the yield is class enlistment (the most inescapable gathering may be returned) , the article is described by a larger part vote of its neighbors, with the item being designated to the class commonly typical among its k nearest neighbor.

Precision, Precision and audit are system used for evaluating the presentation of slant mining..To finish up results expresses that the classifiers shows consequences of the informational index with the Naive Bayes' Theorem.

**Recupero,. et al (2015)** The purpose of the research study is consequently mining sentiments from natural language content using Sentilo anaysis is increasingly pulling

in the interest of the scholarly community and industry. This is the objective of Sentiment Analysis (SA) , which has a significant cover with sentiment mining, a somewhat as of late developing research field whose point is to distinguish and separate emotional data, for example, sentiments or feelings,   materials.

Issue addressed with these methodologies is that the y predominantly depend on parts of Content in which sentiments are unequivocally communicated through positive and negative terms. In numerous cases, sentiments are expressed certainly through setting and area subordinate concepts,making The presentation of NLP-based methodologies limited. This has been the fundamental Inspiration driving the possibility of sentic registering.

The objective of the study was to give an enhanced conventional portrayal of feeling sentences,  a SA algorithm to process theme level just as sentence-level accuracy scores.Author have assessed Sentilo execution at registering by and large sentence sentiment polarity on a corpus of user based lodging audit selected from TripAdvisor. On going work focuses on removing feeling highlights from all diagram Examples that are created by Sentilo, and on structuring a calculation to calculate perspective based feeling scores. Computational knowledge methods ,including fuzzy thinking, combinatorial advancement, learning mechanisms, sentic pattern and discovery and analogical thinking Are under scrutiny as potential augmentations of Sentilo's approach.

## 2.1.4 Literature Review on Collaborative Technique

**Chen,. et al (2012)** The point is recommend useful tweets to a user is a noteworthy test and the point of convergence of this exploration study. Intuitively, a tweet is useful to a user, if the user is enthusiastic about or ready to scrutinize the tweet. Whether or not a user is enthusiastic about a tweet is directed by various segments, for instance, the nature of the tweet and the Authority of the distributer . Individual is furthermore a noteworthy factor to pick whether a tweet is before long important.

The issue is conventional methodologies dismember the substance of users' introduced or retweet statuses on discover the subjects of energy for Twitter users. Regardless, profiling users' individual advantages right now extremely troublesome. Network community oriented separating and Collaborative situating are promising progressions for recommender systems. Collective situating is a situating variation of helpful separating and gives thing positioning outcomes as demonstrated by relative inclination rather than user rating estimation. It works by finding the relationship among users and things reliant on watched user inclinations with the objective that clandestinely user inclinations can be translated from the watched ones.

The target of the research study are as per the following : Author use point level latent factors of tweets to get users' typical advantages over tweet content, which makes us handle the issue of information sparsity in users' retweet exercises. This grants us to change the synergistic separating framework to deal with the proposal issue. Author acquaint inactive segments with show users' social relations, which colossally influence users' decisions in a casual association. The model solidifies express features, for instance, Authority of the distribute and nature of the tweet, which can help further improve the proposition results.

Here Author present updating the situating worldview for per tweet proposition. Helpful situating is an increase of the idle factor model with situating streamlining measure . In the inactive factor model, each user u moreover, thing I have a low dimensional depiction in the idle component space. The rating score is foreseen by assessing the prejudice among user and thing:

$$\hat{y}_{u,i} = \mu + b_u + b_i + p_u^T q_i$$

Where y is the foreseen inclination of user u for thing I. μ is the general ordinary rating, and bu and bi are user inclination and thing inclination on rating score.

To conform to the circumstance of tweet situating, Author change the model for the network community oriented situating setting as showed by situating improvement

worldview. Give a user u and two tweets k and h the pairwise situating model for tweet inclination is characterized as follows:

$$P\left(r(k) > r(h)|u\right) = \frac{1}{1 + e^{-(\hat{y}_{u,k} - \hat{y}_{u,h})}}$$

Where ˆyu,k is the foreseen inclination of user u for tweet k, r is short documentation for rank solicitation. Condition 2 models the probability of thing sets' rank solicitations for a given user. Author can get inclination pair of things for a given user by tolerating a user lean toward the Tweets the retweet to the rest of tweets. Formally, Author define rank inclination set D as follows: D = {< u, k, h > |k ∈ Re(u),h/∈ Re(u)}

$$\mathcal{D} = \{< u, k, h > |k \in Re(u), h \notin Re(u)\}$$

Where Re(u) is the course of action of tweets user u retweeted. Be cause the amount of potential choices of negative model h is gigantic, Author use testing methods to get negative models in the planning framework.

$$\min \sum_{<u,k,h> \in \mathcal{D}} \ln\left(1 + e^{-(\hat{y}_{u,k} - \hat{y}_{u,h})}\right) + \text{regularization}$$

Author have exhibited that introduced CTR strategy staggeringly improves the recommendation execution. Since the CTR model involves three significant fragments — tweet point torpid segments, user social connection factors and unequivocal features — Author should know the adequacy of each part.

**Wang, K. and Tan, Y. (2011)** The motivation behind the research study depends on the improvement of accuracy. Innocent Bayesian system is a famous order calculation and it could similarly be used in the recommendation field. Right when components influencing the order results are unexpected self-sufficient, gullible Bayesian procedure is wind up being the course of action with the best execution.

The goal expressed was organized another synergistic separating count reliant on unsophisticated Bayesian system. The new figuring has a near multifaceted nature to blameless . Bayesian system. Regardless, it has a change of the independence which makes it possible to be applied to the model where prohibitive opportunity supposition that isn't obeyed cautiously. The new figuring outfits us with another fundamental answer for the nonappearance of opportunity other than Bayesian frameworks. The extraordinary presentation of the estimation will outfit customers with progressively accurate proposition. F-measure is the common mean of precision and review. Precision is the amount of right proposals isolated by the amount of all returned proposition and audit is the amount of right recommendations confined by the amount of all the acknowledged interests expected to be found .

## 2.1.5 Literature Review on Content based Technique

**Van Meteren., et al (2000)** In this research study recommender framework PRES (acronym for Personalized Recommender System) is presented. PRES makes dynamic hyperlinks for a site that contains an assortment of exhorts about do it without anyone's help home improvement. The reason for these dynamic hyperlinks is to make it simpler for a user to discover fascinating things and in this manner improving the connection between the framework and the user. A site can be organized by partitioning its Author site pages into content pages and route pages. The content pages furnish the user with the intrigue things while the route pages help the user to look for the intrigue things. This is certainly not an exacting characterization be that as it may. Pages can likewise be cross breed as in the two of them give content just as route offices.

Content-based isolating is utilized as Recommender structures are a one of a kind sort of information filtering systems. Information isolating deals with the movement of things looked over a colossal arrangement that the user is most likely going to find fascinating or important and can be viewed as a request task. PRES A user profile is

discovered from analysis given by the user. The recommender system differentiates the user profile and the reports in the collection. The reports are then situated dependent on specific criteria, for instance, closeness, interest, proximity and noteworthiness and the best situated reports appear as hyperlinks on the current site page.



Figure 2.3 The PRES architecture

A content based separating framework chooses things dependent on the connection between the things and the  intrigues rather than a community collaborative separating framework that picks things dependent on the connection betweenwith comparative preferences. PRES is a content based filtering  framework. It makes proposals by contrasting a user profile and the content of each report in the assortment. The content of a report can be spoken to with a lot of terms.

To evaluate PRES three imaginary users author made which all had various interests. Two themes about home improvement author relegated to the principal user, three to the second user and the to the third user. For every subject all the pertinent reports in the assortment author chosen and delegated important in the event that it contained data that was related with the point. The level of records that the users chose per theme additionally contrasted.

## 2.1.6 Literature review on Machine Learning Technique

**Bellogín, et al (2008)** The research study study shows a recommender system takes a gander at the user's profile or use data to a couple of reference properties. These characteristics may have a spot with the information things (in the substance based system), or to the user's social condition (in the aggregate filtering approach). Blends of the two philosophies have been in addition investigated in the supposed creamer recommender systems

Author propose the use of Machine Learning (ML) techniques to analyze what's more, acknowledge which user and structure features/settings in a modified recommender system are progressively adequate for right recommendations.

Logically difficult to fathom the reasons and conditions under which the given proposition arrange the user profile. The central target of the Authors' investigation are the creation of models that identify with positive (important) and negative (non-significant) proposition cases, and the research study of these plans with ML systems in order to make sense of which features (inclinations) have all the earmarks of being commonly basic to give either positive and negative recommendations. With this preliminary Author in like manner expected to look into the effect of each segment in the consolidated structure, assessing the precision of the proposals when a blend of the two models is used.

For every proposal surveyed (assessed) by a user Author make a model. The properties of the model identify with the characteristics Author need to inspect, and their qualities are gotten from log information databases. System adaptable to the current status of the separated inclinations, and survey whether the proposals improve with the changes.

39

**Naveed, N. (2011)** The research study is stated that right now on the content of a tweet and train an expectation model to gauge for a given tweet its probability of being retweeted dependent on its content. The problem statement question of what makes a message be retweeted has frequently been tended to, yet for the most part in a situation of retweet prephrasing for a given user and with an attention on the structure of the informal organization.

The objective of the Research study is that the author consider the issue of realizing which tweets are re-tweeted, in light of a wide scope of content highlights and independent of setting Information, for example, the user's position Int he social net work and the time stamp of a tweet. That it is conceivable to anticipate which tweets are retweeted. By breaking down the parameters learned in the expectation model, author recognize the highlights that contribute most emphatically to the Likelihood of a tweet being retweeted.This takes into consideration a More profound in sight in to what is of enthusiasm for the Twitter community.

Regression Analysis. Author utilize calculated relapse to figure the likelihood of another tweet being retweeted. Logistic regression is a summed up direct Relapse technique for taking in a mapping from any number of numeric factors to a paired or probabilistic variable. In the Twitter setting,author take in a mapping from the highlights of a tweet to the paired worth demonstrating retweets.

Author acquire include loads that show their impact on the likelihood of a Message being retweeted. By taking a gander at these authorights,author can comprehend what impacts the retweet conduct in Twitter and in End can reason accuracys on what the users considerate laying on a worldwide scale. By computing the highlights for another message and applying the Capacity characterized in Equation. Author obtains probability for this New message to be retweeted.

The result shows measure the accuracy of retweet prediction.Therefore,author split the arrangement of tweets Into a training and data set dependent on the time stamps of the tweets. The preparation set comprises of all tweets with the most reduced time stamp Qualities and contains 75% of the accessible information set.The remaining

40

25% of the information are held for the test set on which author assess The forecast quality. Conclusion, author utilized a learning approach dependent on unadulterated Highlights to foresee the likelihood of a message to be retweeted

**Chettri., et al (2015)**   K Nearest Neighbors is an essential estimation that stocks each current case and portrays new cases/data subject to a relative measure (e.g., division limits). It has generally used for a significant long time as a convincing gathering mode. It acknowledge that the data is in a part space.   naive Bayes request point is to build up a standard which license assigning future things to a class, given only the vectors of components portraying the future articles.

Among three figuring that have been analyzed case base reasoning is better as demonstrated by concentrate as case base reasoning reductions the computational time and see progressively tangled cases and besides work with the gigantic data similarly as on immaterial course of action of appreciated cases equipping the case base.

**KumarR. and VermaR., (2012)** In the current research study, author have focused on the strategies important. Specifically, this work is worried about grouping issues in which the yield of occasions concedes just discrete, unordered qualities.

C4.5 algorithm is explained as take as information a variety of cases, each having a spot with one of not many classes and depicted by its characteristics for a fixed plan of attributes, and yield a classifier that can unequivocally predict the class to which another case has a spot. It has limitations such as empty branches, inisgnificant branches and over fitting.

KNN is assumed that an article is tested with a lot of various qualities, yet the gathering to which the item has a place is obscure. Accepting its gathering can be resolved from its qualities; various calculations can be utilized to computerize the

BIS (Hons) Information Systems Engineering
Faculty of Information and Communication Technology (Kampar Campus), UTAR

order procedure. Naive Bayes one of which has a spot with a known class, and all of which has a known vector of variables, the point construct a standard which will allow us to give out future things to a class, given only the vectors of elements portraying the future articles. SVM ( Support Vector Machine)  extremely successful strategy for relapse, order and general example acknowledgment. It is viewed as a decent classifier as a result of its high speculation execution without the need to include from the earlier information, in any event, when the element of the info space is extremely high.

**Jadhav, et al (2016)** This research study bases on research study of various request frameworks, their inclinations and weaknesses. The essential goal of data mining is to remove the significant information from gigantic unrefined data and changing over it to a sensible structure for its effective and capable use. The goal is Characterization count gives out each case to a particular class with the ultimate objective that request bungle will be least.

K-Nearest Neighbor Classification all around lies in close area. K-Nearest Neighbor is event based learning strategy. Model based classifiers are also called lazy understudies as they store the total of the readiness tests and don't build a classifier until another, unlabeled model ought to be requested. A Naive Bayes classifier ponders that the proximity (or nonappearance) of a particular feature(attribute) of a class is immaterial to the closeness (or nonattendance) of whatever other component when the class variable is given.

Choice Tree Induction. Decision tree learning uses a decision tree as a judicious model which maps recognitions about a thing to choices about the thing's goal worth. Decision tree count is a data mining acknowledgment frameworks that recursively fragments a data set of records using significance first voracious procedure or extensiveness first technique until all the data things have a spot with a particular class. The information in choice tree is spoken to as [IF-THEN] rules which is simpler

for people to comprehend. None of the calculation can fulfill all obliges and criteria. Contingent upon application and necessities, explicit calculation can be picked.

**Pernkopf (2005)**   The motivation behind the exploration study is that the k-NN classifier performs well for the circumstance where the amount of tests for learning the parameters of the Bayesian framework is close to nothing. Bayesian framework classifiers defeat explicit k-NN methodologies to the extent memory necessities and computational solicitations. This research study shows the nature of Bayesian frameworks for classification..

Each vertex (center point) of the chart addresses a discretionary   while the edges get the quick conditions between the variables. Naive Bayes classifier •Tree extended straightforward Bayes classifier. Selective unhindered Bayesian framework classifier Search-and-score structure learning.

The filter approach reviews the significance of the features from the instructive record and the assurance is basically established on true measures. The effects of the picked remembers for the introduction of a particular classifier are excused. On the other hand, the wrapper approach uses the classification execution of the classifier itself as a significant part of the mission for surveying the segment subsets.

The particular k-NN classifier on incessant qualities imperceptibly beats the CFS-SUN classifier. In any case, the cultivated classification execution on the hold-out enlightening record is the identical for both. As referenced over, the Bayesian framework classifiers use a discrete incorporate space.

This research study examines Bayesian framework classifiers to the specific k-NN classifier. The particular k-NN classifier uses a subset of features which is set up by strategies for sequential incorporate assurance methods. In order to get acquainted with the structure of the Bayesian frameworks, the incline climbing search and the back to back forward floating figuring are used.

**Pak,. et al (2010)** Author perform semantic assessment of the accumulated corpus and explain discovered miracles. The goal is to a procedure to accumulate a corpus with positive and negative decisions, and corpus of compositions. the technique licenses to assemble negative and positive inclinations to such a degree, that no human effort is required for orchestrating the records and perform true etymological research study of the gathered corpus.

Directed trial assessments on a lot of genuine microblogging presents on demonstrate that the introduced technique is proficient and performs superior to recently proposed strategies.The gathered data set is utilized to remove includes that is destined to be used to prepare the supposition classifier. Explored different avenues regarding unigrams, bigrams, and trigrams. The review talks about the most ideal approach to thus assemble a corpus for end assessment and feeling mining purposes.

Author construct a slant classifier utilizing the multinomial Naive Bayes classifier. As to Increase precision, expand the exactness of the arrangement, author ought to dispose of normal n-grams, for example n-grams that don't unequivocally demonstrate any conclusion nor show objectivity of a sentence. Such n-grams show up equitably in all datasets.

**Çano,, et al (2017)** The author addressed the most pertinent issues considered and present the related information mining and suggestion systems used to conquer them. author moreover investigate the hybridization classes every half and half recommender has a place with, the application areas, the assessment procedure and proposed future research bearings. In view of the findings, the vast majority of the investigations join cooperative filtering with another system frequently in a manner. Hybrid recommenders speak to a decent premise with which to react in like manner by investigating more current openings, for example, contextualizing suggestions, including equal mixture calculations, preparing bigger datasets, and so on.

CF-X

Here author report examines that consolidate CF with one other strategy which isn't CBF (those are considered CF-CBF). A case of this blend is where the Authors go half and half to improve the exhibition of a multi-criteria recommender. They base

BIS (Hons) Information Systems Engineering
Faculty of Information and Communication Technology (Kampar Campus), UTAR

their authorr on the suspicion that normally just a couple of choice criteria are the ones which sway user inclinations about things and their relating appraisals.

CF-CBF

This is an exceptionally mainstream hybrid RS using the two best suggestion strategies. By and large the proposals of the two frameworks are authorighted to create the last rundown of expectations. In different cases the mixture RS changes from CF to CBF or is comprised of an increasingly mind boggling kind of blend

CF-CBF-X

Those are cases in which CF and CBF are joined together with a third methodology. These sort of suggestions are especially valuable in on-line informal communities (e.g., for promoting). The objective of the Authors is to give acceptable recommendations in information sparsity circumstances.

IICF-UUCF

Thing Item CF and User-User CF are two types of CF recommender, varying in transit the areas are shaped. A few research study consolidate them two to improve generally CF performance. CBF-X

### 2.1.7 Literature review comparison of Machine Learning Algorithms

**Thanh Noi, P. and Kappas, M., (2018)** It is reasonable for an examination to analyze and assess the presentation of RF, support vector machine and kNN for land use. The problem that has been discussed is Distant detecting satellite pictures are considered as one of the most significant information hotspots for
land use planning because of their broad geological inclusion at a productive expense while giving indispensable data on the world's surface. Be that as it may, the precision and preparing season of land use/spread guides utilizing distant detecting pictures is as yet a test to the far off detecting network.

The objective is to assess the presentation of the three most expanding classifiers, RF, kNN, and support vector machine. At the point when applied to a Sentinel-2 picture and to survey the impacts of the preparation tests size, techniques, and type (adjusted/imbalanced) on the precision of the grouping aftereffects of the three previously mentioned classifiers.

Methodology that has been use in the study region was chosen dependent on the land spread attributes and the accessibility of far off detecting symbolism information. The distant detecting picture was preprocessed, climatically adjusted, and cut to the investigation region.The preparation information (preparing and testing tests) was gathered dependent on the manual understanding of the first Sentinel-2 information and high-goal symbolism accessible from Google Earth.

For a precise appraisal of the order results, 650 focuses for each land spread class was gathered. Notwithstanding, to guarantee that the preparation and testing datasets Authorre autonomous, Author supported 15 m for all point tests (testing dataset) and eliminated focuses which had supported focuses crossing with (or having a place with) polygon tests. As results, Author got the quantity of testing focuses (pixels).

**Wang, W., Xi, J., Chong, A. and Li, L., (2017)** A solid requirement for an efficient technique to demonstrate, group, and comprehend different driving styles. So as to upgrade the presentation of driving style classification and diminish naming endeavors, a semi supervised technique, to be specific, a semi supervised uphold vector machine (S3VM), was created by joining the benefits of administered and undirected techniques.

In this research study, Author predominantly center around drivers' longitudinal driving practices with respect to forceful and typical driving styles when driving on stunning streets. The S3VM strategy considers the extra data of unlabeled driving information to catch more fundamental attributes of the driving information. Second, the S3VM approach produces an ideal choice limit by completely using the information on the unlabeled information and named information.At last, tests are

46

directed to analyze the S3VM and support vector machine approaches. Feature Selection: To research drivers' driving styles, speed and quickening/deceleration Authorre utilized to catch driver qualities.

**Zendehboudi, A., Baseer, M.A. and Saidur, R., (2018)** Giving a thorough and basic investigation of the best in class audit on the use of help vector mama chine models for sun based and wind vitality anticipating.

Be that as it may, there are still issues, including huge range vari capable measure of intensity age, because of variety in wind speed also, heading. For sun poAuthorred vitality, various boundaries, for example, sunlight based rise point, fog impact and overcast spread, will cause fluctuations in yield. The irregular and variable yield could prompt Authorighty negative effects on framework, poAuthorr transmission and circulation hardware, which forestall across the board utilization of environmentally friendly poAuthorr vitality age.

To evaluate the exhibition of support vector machine in the fifield of inexhaustible vitality, the creators have done a far reaching survey on the support vector machine models in wind and sun based vitality assets. By receiving a precise writing audit procedure, fifirstly, the distributions in the ideal zones are removed and chosen; and, also, the examination of the distributions was completed by gathering them into various classes dependent on the support vector machine technique application.

VR, uphold vector relapse, is the support vector machine usage for work guess and relapse. Diverse fundamental part works are utilized in support vector machine models. The capacities can be classified as polynomial (Poly), exponential outspread premise work (ERBF), outspread premise work (RBF), sigmoid and straight. Building up an all inclusive steady model for wind and sun oriented vitality forecast isn't doable. Both sunlight based and wind re sources display variable examples dependent on an assortment of elements.

**Manavalan, B. and Lee, J., (2017)** Author applied support vector machineQA on different benchmarking datasets and the outcomes show that support vector

machineQA performed altogether better than other single model techniques in positioning protein 3D models just as in choosing the best model from the pool. Besides, support vector machineQA was aimlessly tried in the CASP12 test. current test was to prepare a support vector machine to accur ately map input highlights extricated from a 3D model to its TM score/GDT_TS score; this is vieAuthord as a relapse issue. The most urgent aspect of this undertaking is to remove a lot of significant highlights.

Normal trial methods such as X-beam crystallography, NMR and electron microscopy are expen sive and frequently tedious methods of deciding the 3D structures of uncharacterized protein arrangements. Therefore, an immense hole exists between the quantity of realized protein groupings and their e perimentally illuminated 3D structures. Assessment. Model precision was assessed utilizing three integral meas ures: (I) Pearson's relationship coefficient (CCrank), Spearman's position relationship (qrank) and Kendall's tau connection (srank) between the real positioning and anticipated positioning; Average TM-score or GDT_TS misfortune; and Z-score. Execution of support vector machineQA on I-TASSER imitations. Assessed the exhibition of support vector machineQA on the I-TASSER set which comprises of 56 non-homologous focuses with each target containing around 300–500 models support vector machineQA is better than the vitality based techniques.

Execution of support vector machineQA on 3DRobot baits. Joined execution of support vector machineQA on I-TASSER and 3DRobot sets show that the support vector machineQA's blend of different factual potential vitality based terms and consistency-based terms between anticipated furthermore, determined estimations of 3D models improves the exhibition of QA.

Execution of support vector machineQA on in-house CASP11 models. . Execution of support vector machineQA on CASP11 targets. Author assessed the exhibition of support vector machineQA on CASP11 targets. For this reason, Author utilized 88 focuses for both Stage1 and Stage2, as utilized in the authority CASP11 evaluation. Performance of different strategies on Stage2 CASP11 targets. Author note that support vector machineQA beat our past technique RFMQA both in positioning basic

models and in model determination. Author note that support vector machineQA beat our past technique  RFMQA both in positioning basic models and in model determination. result demonstrates that in spite of the fact that support vector machineQA was prepared on single area targets, it per shaped sensibly Authorll for multi-space models. Execution of support vector machineQA in CASP12 . These outcomes demonstrate that support vector machineQA can add to the fruitful 3D demonstrating of troublesome objective proteins as far as model choice.

**Jayakumar Sadhasivam., Ramesh Babu Kalivaradhan., and Senthil Jayavel(2019)** The research study is conducted to assess execution of the calculation. It is trying to accomplish high precision in Twitter slant investigation.The element extraction comprises of two stages. Making Dataset. As the language utilized in the tweets is unique in relation to standard language. The spelling and the utilization of the words may contrast. As it a managed learning strategy, it needs a preparation set to hinder mine the extremity of the tAuthoret.

Information Preprocessing. The tweets are first preprocessed, and afterward it is gone to the classifier. Estimation Classification. In the wake of playing out the above advances, grouping is finished utilizing naive Bayes, Support Vector Machine, Maximum Entropy and Group Classifiers.

In a real world situation, the exhibition of Naïve Bayes classifier is better com pared to different calculations. The Naïve Bayes classifier makes string suspicion about the freedom of highlight which is the first hindrance of the Naïve Bayes classifier with high precision and produces strong outcomes regardless of whether there are a few bugs in the preparation test. Support Vector Machine gives a remarkable arrangement greatest entropy classifiers particularly like Naïve Bayes classifier, aside from it doesn't expect the freedom of the highlights.

**Talbot, R., Acheampong, C. and Wicentowski, R., (2015)** For a mes sage passing on both a positive and negative sentiment, whichever is the more grounded opinion ought to be picked' a framework which uses a Credulous Bayes classifier to decide the supposition of tweets. The preparation information Author utilized incorporates 8142 tweets, each marked as certain, negative or unbiased. Vocabulary of approximately 6800 tokens named as either certain negative

Preprocessing Steps highlights utilized in our classifier are unigrams, invalidated unigrams, and two uncommon labels showing the nearness or nonattendance of words in the tAuthoret be ing found in the estimation vocabulary utilized a Naive Bayes classifier to arrange the tweets. naive Bayes depends on the supposition of contingent autonomy among the highlights, some falsehood here

Human-created rundown of positive and negative words and images, whose presence naturally abrogated the classifier's choice, ought to be additionally investigated. It is almost certain that more words and images exist whose nearness is profoundly demonstrative of a negative or positive tAuthoret sentiment. Programmed production of these rundowns would likely improve execution and be more experiment count justified.

**Acosta, J., Lamaute, N., Luo, M., Finkelstein, E. and Andreea, C., (2017)** This examination explores different avenues regarding various calculations to discover which one yields the most elevated prediction exactness. Examination breaks down the similitude between the posts of a client so as to help check the cause of the post. Notwithstanding recognizing similitude between their notions. Goal of this investigation is to decide if utilizing the word embedding created by the word2vec calculation could be utilized to order supposition. huge thought is that by utilizing word embedding, there will be no compelling reason to physically make highlights , ased off stylometry so as to group assumption precisely

BIS (Hons) Information Systems Engineering
Faculty of Information and Communication Technology (Kampar Campus), UTAR

Philosophy 14,640 tweets and 15 properties including the first tAuthoret text, Twitter client related information, and the class slant name.

scikit-learn's train_test_split work was utilized to part the Twitter posts and their individual estimation names. 70% of the examples Authorre designated for preparing and 30% for testing purposes, one model utilizing constant sack of-words (CBOW) and 2) another model - utilizing skip-gram (SG). Models Authorre tried utilizing both various leveled softmax and negative testing. Various leveled softmax gave higher precision scores and in this way it was picked. SG preparing model beat CBOW model in each classifier prepared. Gaussian Naïve Bayes had a precision rate of 64% for CBOW, anyway when utilized with the SG model roduced a precision pace of 69%. The most elevated precision rate yielded by a classifier was 72% utilizing Support Vector Classifier and SG as the word2vec preparing model.

**Joshi, S. and Deshpande, D., (2018)** Author have endeavored to direct notion investigation on "tweets" utilizing extraordinary AI calculations. Author endeavor to group the Polarity of the tAuthoret where it is either certain or negative. On the off chance that the tAuthoret has both positive and negative components, the more predominant notion ought to be picked as the last name.

Endeavored to arrange human opinion into two classes to be specific positive and negative System machine learning and characteristic language handling (NLP) calculations, unigrams and bigrams which is a type of portrayal of the "tAuthoret". Author utilize different AI calculations dependent on NLP (Natural Language Processing) to lead notion examination utilizing the extricated highlights. test dataset have 800000 and 200000 tweets individually extricate two kinds of highlights from our dataset, in particular unigrams and bigrams. naive Bayes, Maximum Entropy or support vector machine which figures out how to make expectations. To assess our framework, Author use Baseline Classification which is our assessment metric in which test information is taken care of to the educated calculation which consequently

creates suggested forecast appraisals of words. With help of pre-ordered brilliant set and assessment metric Author check the precision of our model.

**Kharde, V. and Sonawane, P., (2016)** measure of substance produced by clients is too vasT for a typical client to break down. So there is a need to computerize this, different assumption examination strategies are broadly utilized. The twitter dataset utilized in this overview work is as of now named into two classes viz. negative and positive extremity and hence the supposition investigation of the information turns out to be anything but difficult to watch the impact of different highlights Highlight Extraction

AI Naive Bayes: It is a probabilistic classifier and can gain proficiency with the example of looking at a lot of records that has been order To prepare and arrange utilizing Naïve Bayes Machine Learning method ,Author can utilize the Python NLTK library Most extreme Entropy Classifier, no suppositions are taken as to relationship in the middle of the highlights extricated from dataset. This classifier consistently attempts to expand the entropy of the framework by assessing the restrictive dispersion of the class mark.

Backing vector machine examines the information, characterize the choice Boundaries and utilizes the bits for calculation which are Performed in input space. The info information are two arrangements of vectors of size m each. At that point each information which spoke to as a vector is grouped into a class
Vocabulary Based. Vocabulary based technique utilizes assumption word reference with supposition words and match them with the information to decide polarity. They relegates slant scores to the assessment words depicting how Positive, Negative and Objective the words contained in the word reference are. corpus-based methodology have target of giving word references identified with a particular area. These word references are produced from a lot of seed feeling terms that becomes through the hunt of related words by methods for the utilization of either factual or semantic procedures.

**Maipradit, R., Hata, H. and Matsumoto, K., (2019)** Author used an AI based methodology utilizing n-gram highlights and a computerized AI device for slant classifification. Despite the fact that n-gram phrases are vieAuthord as useful and helpful contrasted with single words, utilizing all n-gram phrases is anything but a smart thought on the grounds that of the huge volume of information and numerous pointless highlights address this issue, Author use n-gram IDF, a hypothetical augmentation of Inverse Document Frequency (IDF). IDF quantifies how much data. Used auto-sklearn, which contains 15 classifification calculations (arbitrary woodland, piece support vector machine, and so forth.),

Highlight extraction utilizing N-gram IDF. Mechanized AI. To group sentences into positive, unbiased, and negative, Author use auto-sklearn, an au tomated AI instrument. SentiStrength NLTK is a characteristic language toolbox and can do sentimen investigation dependent on dictionary and rule-based VADER (Valence Mindful Dictionary and sEntiment Reasoner), which is specififi- cally tuned for notions communicated in online media. Stanford CoreNLP .SentiStrength-SE .Stanford CoreNLP SO. Performances in auto-sklearn for each dataset. Stack Overflow: Linear Discriminant Analysis, Libsupport vector machine Backing Vector Classification, and Liblinear Support Vector Classification. App surveys: Random woods, Libsupport vector machine Support Vector Order, and Naive Bayes classifier for multinomial models.

**Ahuja, R., Solanki, A. and Nayyar, A., (2019)** The proposed recommender framework predicts the client's inclination of a film based on various boundaries. The recommender framework takes a shot at the idea that individuals are having basic inclination or decision. These client will inflfluence on one another's feelings. This cycle streamline the cycle and having loAuthorr RMSE.
K Nearest Neighbor Inside Clustered Sum of the Squared strategy is used to fifind the privilege no bunches with the goal that K-implies grouping can be applied to the film At that point the likeness between the clients is determined utilizing the Pearson

Correlation Matrix. At that point, utilizing the KNN prediction for the film evaluations for top N clients is finished.

Root Mean Square Error (RMSE) It is seen from the chart that for the current strategy the RMSE esteem is 1.23154 for bunch equivalent to 68, RMSE esteem utilizing proposed procedure is 1.233 to 19 groups furthermore, RMSE esteem utilizing proposed method is 1.081648 to 2 bunches. Seen that subsequent to actualizing the framework in the python programming language the RMSE esteem of the proposed strategy is superior to the current method.

**Xydas, E.S., Marmaras, C.E., Cipcigan,n.d (2013)**

The point of a force framework is to oblige the clients' requirement for electric vitality at the least expense and guaranteeing the nature of gracefully to downstream clients. Anticipating the poAuthorr charging request dependent on numerical models and notable burden information contrasts from predicting the charging request profile dependent on suppositions for example, travel and driving examples.

to introduce and look at a man-made brainpoAuthorr EV load estimating strategy utilizing Support Vector Machine

Short term Load Forecasting (STLF) is utilized to anticipate the heap request from one hour to multi Authorek. Medium-Term Load Forecasting (MTLF) is utilized to anticipate the poAuthorr load request from multi Authorek to one year in advance.

The preparing part of the Support Vector Machine model is the charging information for the initial 50 Authoreks of the year Support Vector Machine model is more delicate to hourly vacillations and gives better execution.

Mean absolute percentage error MAPE rule is the acknowledged business standard for estimating load anticipating precision and a worth under 5% is Worthy the family member blunders of the Support Vector Machine figure are littler than the Monte Carlo ones, and this legitimizes the enhancement for RMSE. MAPE for support vector machine 3.69 while monte carlo 8.99 RMSE for support vector machine 5013

and monte carlo 105.52    The outcomes demonstrate the viability and exactness of the support vector machine proposed model, over a more factual methodology.

**Chaithra, V.D., (2019)** Estimation examination of versatile unpacking video remark helps investigating the client's response towards the cell phone. The meta data of the video express the watcher's response towards the telephone. Meta data, for example, different preferences don't give nitty gritty opinion of the watchers, accordingly remarks are considered to dissect the watcher's conclusion. Feeling VADER which is a standard based methodology is applied utilizing the python bundle SentimentIntensityAnalyzer. This relegates the power of individual remark how much is positive, negative and unbiased. Managed AI calculation Naive Bayes is applied on the pre-prepared information to predict the assumption of the remarks. The classifier is prepared with information and tried on the inconspicuous information or the testing information. The data set is separated in the proportion of 7:3 with 70% information for preparing and 30% for testing. From the information all the unbiased remarks are taken out making it a parallel arrangement

Estimation VADER applied to name the remarks with the range estimation of 0.2 and negative 0.2 outcomes viably marking the remarks. Table 1 shows the consequence of applying the Sentiment VADER on the data set. The classifier was then tried on the 30% inconspicuous information and an exactness of 79.78% and F1 Score of 83.72% was accomplished.

**Jia, Y. and SungChu, S., (2020)** The target and utilize true information of administration calls, which represents extra difficulties regarding the artifificial datasets that have been normally noise Vocal articulation is an essential transporter of affec tive signs in human correspondence. Discourse as signals contains a few highlights that can extract semantic, speaker-specifific data, and emo tional. Related work about sound based sentiment examination alongside multimodal combination

All the calls are introduced in Wave design with an example pace of 8000 Hertz what's more, span shifting from 4 minutes to 26 minutes.

The start of each call contains a presentation of the clients' organization name by a robot. To address this issue, the fragment before the first delay (quiet length > 1 second) is taken out from each call.

500 calls are gathered from the call place information base covering assorted subjects, for example, protection plan in development, protection id card, subordinate inclusion, and so on. The call informational index had female and male speakers arbitrarily chose with their age running approximately from 16-80.

To dispose of commotion and long delays (quiet length > 5 seconds) in calls, Voice Activity Detection (VAD) is applied, trailed by the application of Automatic Speech Recognition (ASR) and Automatic Speaker Diarization (ASD) to decipher the verbal explanations, remove the beginning and end time of every expression, and distinguish the speaker of every articulation.

Assumption explanation is a difficult errand as the name relies upon the annotators' viewpoint, and the distinctions inborn in the manner individuals express Feelings To display data for notion examination from calls, we first acquire the streams relating to every methodology. Loudness is the abstract impression of sound pressure which is identified with sound force. Sharpness is a proportion of the high-recurrence substance of a sound, the more noteworthy the extent of high frequencies the more honed the sound. Talking rate is regularly defined as the number of words expressed every moment The test data set comprises of 21 calls with 1,890 utterances, which are physically clarified for negative (848) and non negative (1,042).

**Renault, T., (2019)** Author utilize a dataset of 1,000,000 messages sent on the microblogging stage StockTwits to assess the presentation of a wide scope of machine learning calculations and preprocessing strategies for feeling investigation in finance

BIS (Hons) Information Systems Engineering
Faculty of Information and Communication Technology (Kampar Campus), UTAR

Author build two datas ets: one adjusted data set containing 500,000 positive messages and 500,000 negative messages, and one lopsided data set con training 800,000 positive messages and 200,000 negative messages. Initially, and considering a basic multinomial Naive Bayes model (ten times cross-approval), author investigate the effect of the size of the data set on the precision of the classification. find that the exactness increments firmly with the size of the dataset: the presentation of the classifier increments by almost 10 rate focuses when the size of the data set increments

from 1000 messages to 10,000 messages, and by 4.3 rate focuses when the size of the data set increments from 10,000 messages to 100,000 messages. The mar ginal exactness improvement is diminishing and the precision arrives at a level around 500,000 messages, the effect of thinking about consistent grouping of words (ngrams) rather than unigrams considering bigrams firmly improves the exactness of the classifcation. For instance, for a data set of 250,000 messages, including bigrams improves the exactness of the classifcation by 2.2 rate focuses contrasted with a preparing with unigrams as it were Trigrams and four-gram have no significant sway on the precision. Author analyze a wide scope of AI strategies utilized in the writing (most extreme entropy, uphold vector machine, irregular woods and multilayer perceptron) to our benchmark Naive Bayes model play out a network look for hyperparameter advancement and we fnd that the best exhibition is accomplished by a Maximum Entropy classifier, firmly followed by a Support Vector Machine classifier.

the precision of the classifcation unequivocally increments with the size of the dataset the precision increments from 59.6% for a dataset of 500 messages up to 73.08% for a dataset of 1,000,000 messages.

The minimal improvement is diminishing: multiplying the size of the dataset from 500,000 messages to 1 million messages just expands the precision by 0.3 percent age focuses, while multiplying the size of the dataset from 25,000 to 50,000 increments

the precision by 1.34 rate focuses. Naive Bayes calculation (NB), a Maximum entropy classifer (MaxEnt), a straight Support Vector Classifer (SVC), a Random Forest classifer (RF) and a MultiLayer Perceptron classifer (MLP).

## 2.2 In-depth Analysis of Literature Review

In depth analysis is conducted to review the details of objective, methodology and result of evaluation of latest literature reviews that has been taken from chapter 2.1

| Author | Chettri, R., Pradhan, S. and Chettri, L., 2015. |
|---|---|
| Objective | focusing on differentiating supervised learning algorithms i.e. K-NN, Naive Bayes and Cased Based Reasoning (CBR) Classifier |
| Methodology | K Nearest Neighbors<br>Naive Bayes classification aim is to construct a rule<br>Case based reasoning |
| Result | KNN- provides accuracy of around 72%, More Parocesing time<br>Naive Bayes-<br>accuracy of about 85%, less compared to knn running time<br>Case Based Reasoning - accuracy of about 92%, Less running time compared to knn and naive bayes |

Table 2.2.1 : Analysis on Literature review 1

| Author | Jadhav, S.D. and Channe, H.P., 2016. |
|---|---|
| Objective | I. study of various classification techniques, their advantages and disadvantages. Ii. extract the useful information from huge raw data and converting it to an understandable form for its effective and efficient use |
| Methodology | Naive Bayes Classification<br>Decision Tree Induction |
| Result | KNN- 0 sec for small data set, and 100% accuracy,<br>Medium-100%<br>Large 89.842%<br>Naive Bayes- 92.85 for small data set ,<br>Medium - 81.667%<br>Large- 63.71%<br>Decision Tree- small 100%<br>Medium- 99%<br>Large - 63.713% |

Table 2.2.2 : Analysis on Literature review 2

| Author | Thanh Noi, P. and Kappas, M., 2018 |
|---|---|
| Objective | examination to analyze and assess the presentation of RF, SVM and kNN for land use/spread |
| Methodology | RF, KNN, Support Vector Machine<br>Sentinel-2 image |
| Result | Support Vector Machine has better performance<br>Sentinel-2 image |

Table 2.2.3 : Analysis on Literature review 3

BIS (Hons) Information Systems Engineering
Faculty of Information and Communication Technology (Kampar Campus), UTAR

| Author | Jayakumar Sadhasivam, Ramesh Babu Kalivaradha2, Senthil Jayavel , 2019 |
|---|---|
| Objective | classification is done using Naïve Bayes, Support Vector Machine, Maximum Entropy and Ensemble Classifiers |
| Methodology | Creating Dataset Data Preprocessing Sentiment classification |
| Result | Datasets : Data Repositiories Twitter API Support Vector Machine higher accuracy than Naive Bayes |
| Conclusion | Support Vector Machine gives a remarkable arrangement greatest entropy classifiers |

Table 2.2.4 : Analysis on Literature review 4

| Author | Talbot, R., Acheampong, C. and Wicentowski, R., 2015 |
|---|---|
| Objective | Message Polarity Classification |
| Methodology | Sentiment Lexicon Naive Bayes |
| Result | 8142 Tweets, Automatic creation of these lists would likely improve performance and be more experimentally justified compared to manual identification. |

Table 2.2.5 : Analysis on Literature review 5

BIS (Hons) Information Systems Engineering
Faculty of Information and Communication Technology (Kampar Campus), UTAR

| Author | Acosta, J., Lamaute, N., Luo, M., Finkelstein, E. and Andreea, C., 2017 |
|---|---|
| Objective | experiments with . <br><br> different algorithms to find which one yields the highest <br><br> rediction accuracy. |
| Methodology | 4,640 tweets and 15 attributes including the original tweet text <br><br> Word2Vec <br><br> Scikit Learn <br><br> Continuous Bag of Words Model <br><br> Skip gram <br><br> Support Vector Classifier |
| Result | SG preparing model beat CBOW model in each classifier prepared. Gaussian Naïve Bayes had a precision rate of 64% for CBOW, anyway when utilized with the SG model gives a precision pace of 69%. <br><br> most elevated precision rate yielded by a classifier was 72% utilizing Support <br><br> Vector Classifier and SG as the word2vec preparing model. |

Table 2.2.6 : Analysis on Literature review 6

| Author | Joshi, S. and Deshpande, D., 2018 |
|---|---|
| Objective | endeavor to group the polarity of the tweet where it is either certain or negative |
| Methodology | natural language processing (NLP) algorithms, <br><br> unigrams and bigrams <br><br> test dataset have 800000 and 200000 tweets respectively |
| Result | Naive Bayes, Maximum Entropy, Support Vector Machine to make predictions. |

Table 2.2.7 : Analysis on Literature review 7

| Author | Xydas, E.S., Marmaras, C.E., Cipcigan, L.M., Hassan, A.S. and Jenkins, N., 2013 |
|---|---|
| Objective | to oblige the clients' requirement for electric vitality at the least expense and guaranteeing the nature of gracefully to downstream clients. |
| Methodology | Short Term Load Forecasting (STLF) is utilized to anticipate the heap request from one hour to multi week. Medium-Term Load Forecasting (MTLF) is utilized to anticipate the power load request from multi week to one year in advance. Support Vector Machine |
| Result | Mean absolute percentage error (MAPE) MAPE for Support Vector Machine 3.69 while monte carlo 8.99 RMSE for Support Vector Machine 5013 and monte carlo 105.52 |

Table 2.2.8 : Analysis on Literature review 8

| Author | Ahuja, R., Solanki, A. and Nayyar, A., 2019, January |
|---|---|
| Objective | recommender framework predicts the client's inclination of a film based on various boundaries |
| Methodology | preprocessing of data, separate data frames training set and testing set<br><br>K-means clustering is used to build a separate data frame |
| Result | Root Mean Square Error (RMSE)<br><br>RMSE value is 1.23154 for cluster    68,<br><br>RMSE value using proposed technique is 1.233 to 19 clusters<br><br>and RMSE value using proposed technique is 1.081648 to 2 clusters |

Table 2.2.9 : Analysis on Literature review 9

| Author | Kharde, V. and Sonawane, P., 2016 |
|---|---|
| Objective | measure of substance produced by clients is too vast for a typical client to break down. So there is a need to computerize this, different assumption examination strategies are broadly utilized. |
| Methodology | **Machine Learning**<br><br>Naive Bayes:<br><br>Maximum Entropy<br><br>Support Vector Machine<br><br>**Lexcion Based** |
| Result | Train dataset<br><br>45000<br><br>Test Dataset<br><br>44832<br><br>**Unigram** |

BIS (Hons) Information Systems Engineering
Faculty of Information and Communication Technology (Kampar Campus), UTAR

| | **Naive Bayes** = 74.56 |
| | **Maximum Entropy** = 74.93 |
| | **Support Vector Machine** = 76.68 |
| | **Baseline Algorithm** = 73.65 |

Table 2.2.10 : Analysis on Literature review 10

BIS (Hons) Information Systems Engineering
Faculty of Information and Communication Technology (Kampar Campus), UTAR

## 2.3   Data Collection

Data collection is procedure of get together and assessing information on elements of enthusiasm, in a developed effective structure that enables one to react to communicated examine questions, test hypotheses, and assess results. Information assortment causes one to respond to questions, assess yield and even make expectations. Exact data collection is essential to keeping up the uprightness of research, choosing taught business decisions and ensuring quality insistence. For example, in retail bargains, data might be assembled from flexible applications, webpage visits, dedication programs and online diagrams to get acquainted with customers. The data has been collected using literature reviews and tsudies that has been conducted. The Twitter data sets were obtained from Data Words.

There are two sorts information assortment techniques named essential information and optional information. Essential information will be data gathered by the individual for a reason. It is unique in nature and is explicit to an exploration issue under investigation. Auxiliary information is gathered by others for some other reason. Alison Wolf (2016) says auxiliary information is classified into government insights, industry affiliations, exchange productions, organization sites and statistical surveying reports. It is said the more is the better, information assortment measure of 30 is viewed as adequate to direct vivid history hallway module in booth framework. The information right now was gathered structure studies and diaries distributed in Google Scholar. The benefit of overview strategy is that it advantageous for information assortment. Nowadays, the online investigation procedure has been the most renowned technique for social event data from target individuals. Next to the solace of data gathering, researchers can assemble data from people the world over. It has great factual centrality. Dependability is gotten the exploration considers is constantly given to all members a typical upgrade by adding literature reviews on latest topics and latest reviews been collected to be updated as Information Technology been growing fast.

BIS (Hons) Information Systems Engineering
Faculty of Information and Communication Technology (Kampar Campus), UTAR

## 2.4    Critical remarks of previous work

**Schafer, J.B., Konstan, J. (2012)** The benefit of the research study was making useful recommendations, certain bits of information from media look into, in an improved way, assessment grouped assortment in proposals had outline the foundations of such recommender structures, with the goal that the conduct of the systems can be seen even more eagerly, and changed if essential. The objective of the research study study was to assess the framework dependent on a Bayesian user model, the weakness watched was there no away from of the outcomes and elaboration given.

**Prerana Khurana , Shabnam Parveen (2016)** The advantage obtained from this research study is that it states all the possible challenges to deal with personalised recommnedation system for example overspeacialoization, scalability and cold start problem and has enough defintions towards it, yet it was not proven with any methods or evaluations. The approached stated has no more elaboration and limited. Instead the reseaech study has talked about Basic Matrix Factorization with not enough examples and defintions as to compare with this research study conducted in this Research study the definiton of each term is well explained priorly so that the reader gets knowldege or general clarification what it is about.

**Dey, et al (2016)** This research study has well explained data, results and evaluations. It used two supervised machine learning algorithm Naive Bayes and K-Neared Neighbour to prove the result obtained and gives the user an acurate result. The study has even stated the diffculties faced during the phasse and the ways used to overcome them. The study also elaborated the results obtained in with better explanation rather than tables. This terms has been taken into account while conducting this research.

**Debashis Das (2017)** The advantage visualised is that each and every explanation is guided with flowchart or diagrams to make sure that the user gets better image of what is being explain.

BIS (Hons) Information Systems Engineering
Faculty of Information and Communication Technology (Kampar Campus), UTAR

Each of the recommendation given comes with its own merits and demerits to be able to predict the situation that it can be used. Most of the research studies reviewed has lesser information on the feedback techniques, however in this study it was well defined with the merits and demerits stated in points.

**Fan, W. and Gordon, MD., (2014)** The study propose new ad-hoc and regression measuers. The resultsproved that using textual information results better than star ratings. The study has well described table and data. Each and every method was proved. Every term is defined instead regression methodology is well defined in detailed explanation before proceeding to methodology. These advantages has been guiding to define details of the survey.

**Abirami, A.M. and Gayathri (2017)** The research study well explains approaches of machine learning, dictionary based and ontology based and the challenges while stating clearly the advanatge and disadvantages of the approach in a table as it makes the reader easier to be compared.

Dictionary Based, sentiWordNet explanation and examples was useful for the study. The explanations given has been used , it is one of the significant component of sentiment analysis. It has table level comparisons that compares the reviews conducted by the author in a simplified way.

**CHAPTER 3 : SYSTEM DESIGN**

Research can be characterized as an activity that includes discovering, in a more or less much precise way, items you did not know. Research Methodology part of an exploration depicts look into strategies, approaches and designs in detail highlighting those utilized all through the investigation, justifying decision through portraying advantages and disadvantages of each approach and configuration considering ability in consideration. Methodology ought to be the most fitting to accomplish goals of the research study. The methodology chosen should likewise be made possible.

The structure of research began with a broad question, the focus was narrowed to observe the data needed. The methodologies and results author analysed. The general hypothesis was obtained to make a stand. The quantitative analysis been implemented in assessing the performance of classifiers.

According to Ashley Crossman (2019)Subjective research is a sort of sociology explore that gathers and works with non-numerical information and that tries to decipher significance from these information that help comprehend public activity through the investigation of focused populaces or spots. Individuals frequently outline it in restriction to quantitative research, which utilizes numerical information to distinguish enormous scope patterns and employs statistical operations to decide causal and correlative relationships between factors. Inside human science, subjective research is ordinarily centered around the micro-level of social interaction that forms regular day to day existence, though quantitative research commonly centers around large scale level patterns and wonders.

Qualitattive methodology is observed in obtaining literature reviews and data collection as right data and documentation is needed in the hypothesis statement generation. Subjective research assembles data that isn't in numerical structure. It is huge for investigations of individual level and was profiting to discover top to

bottom the manners in which individuals minds work. There author re several steps carried to in this research to obtain data, first of all by identifying a research problem where exploration of the research and understand was held. Follow author d by reviewing literature. Then a purpose was specifies to obtain a general idea. Data was collected with the review of research research studies. The collected data was analysed and interpreted for understanding the problem statements and objectives. The quality of data collected was determined by further review on literature with verification, dependability and affirmation. As final stage the research was reported with emerging information.

The strategy led right now was subjective approach. It is utilized to watch feeling, considerations and conviction of the mass society. Subjective research was additionally classified without hesitation investigate, contextual analysis explore and grounded hypothesis look into. There are a few sources to pick up information utilizing subjective technique including perception, interviews, survey, reports and messages, the analyst's impressions and content research study.

As to legitimize Qualitative research has the two advantages and downsides. It make an inside and out comprehension of the mentalities, practices, collaborations, occasions and social stages. It enables social researchers to see how things and web based life impacts society. This procedure is adaptable and helpful to adjust to changes in the exploration condition and led with zero or negligible expenses. The disadvantage of the philosophy is that the degree is restricted to the looks into perception. The alert of this technique is to guarantee that they don't impact data in methods for individual predisposition. Thorough preparing is led on scientists to maintain a strategic distance from these issues.

For the investigation explore research studies writer re read to legitimize the issue experienced and upgrade current recommender framework strategies. The exploration directed utilizing calculation of wistful research study. The information research study plan of the goal expressed would be accomplished with specific strategies.

Quantitative examination is the way toward gathering and assessing quantifiable and obvious information, for example, incomes, piece of the pie, and wages so as to comprehend the conduct and execution of a business. Quantitative investigation helps in assessing execution, surveying monetary instruments, and making expectations. It includes three fundamental strategies of estimating information: relapse examination, direct programming, and information mining.

In Short, Qualitative analysis has been vital in the generation of hypothesis. Literature reviews that has been conducted is used in prediction of hypothesis while quanittiative analysis been used in proving the evaluation on performance of the classifiers. While QA fills in as a valuable assessment instrument, it is regularly joined with the correlative examination and assessment device subjective investigation. It is regular for an organization to utilize quantitative investigation to assess figures, for example, deals income, overall revenues, or profit for resources (ROA).

Nonetheless, to show signs of improvement image of an organization's presentation, examiners likewise assess data that isn't effectively quantifiable or diminished to numeric qualities, for example, notoriety or representative spirit. Subjective examination centers around implications, includes affectability to setting as opposed to the craving to acquire widespread speculations, and sets up rich portrayals instead of quantifiable measurements. Subjective investigation tries to answer the "why" and "how" of human conduct.

Quantitative analysis is not something contrary to subjective examination; they are simply various methods of reasoning. Utilized together, they give valuable data to educated choices that advance a superior society, improve money related positions, and upgrade business activities.

Quantitative methodology has been handled in terms of evaluating the machine learning algorithms. Though literature reviews had proven the hypothesis, in the research in terms of analysis the algorithms for Twitter Sentiment Analysis

quantitative analysis is much needed. Naive Bayes algorithm, K Nearest Neighbor algorithm been handle my quantitative or mathematical term. Libraries from Scikit Learn were been imported for analysis of these algorithms. To identify the performance accuracy sore, precision score and recall score been implemented.

## 3.1 Design specifications

The recommender system is done based on sentiment analysis. Opinion investigation is intended to decide the demeanor of a speaker or an author as for some subject or the general extremity of the relevant hypothesis.

The algorithm used is Naive Bayes Theorem based on machine learning methods.

```
┌──────────┐   ┌──────────┐   ┌──────────┐   ┌──────────┐   ┌──────────┐
│Collection│   │   Text   │   │Detection │   │Classificatio│ │Presentat │
│ of data  │──▶│ training │──▶│   of     │──▶│  n of     │─▶│ ion of   │
│          │   │          │   │ sentiment│   │ sentiment │   │ output   │
└──────────┘   └──────────┘   └──────────┘   └──────────┘   └──────────┘
```

Figure 3.1: Stages of obtaining output

## 3.1.1 Collaborative filtering

Collaborative Filtering is as of now one of  frequently utilized methodologies and normally gives better outcomes over substance based recommendations. A few instances of this are found in the recommendation frameworks of Youtube, Netflix, and Spotify.

BIS (Hons) Information Systems Engineering
Faculty of Information and Communication Technology (Kampar Campus), UTAR

These sorts of frameworks use user connections to channel for items of interest. We can envision the arrangement of connections with a grid, where every entry speaks to communication between user and item. An intriguing perspective on Filtering is to consider it a speculation of classification and regression. While in these cases we intend to anticipate a variable that legitimately relies upon different factors (highlights), in shared Filtering there is no such differentiation of feature variables and class variables.Types of collaborative filtering named memory based and model based.

**Memory Based**

There are two methodologies: the first recognizes clusters of users and uses the collaborations of one specific user to foresee the communications of other comparative users. The subsequent methodology distinguishes groups of items that have been evaluated by user A and uses them to anticipate the cooperation of user A with an alternate however comparable item B.



Figure 3.2 Collaborative Filtering, Source ( Carlos Pinela, Medium.com,2017)

BIS (Hons) Information Systems Engineering
Faculty of Information and Communication Technology (Kampar Campus), UTAR

**Model Based**

These strategies depend on AI and information mining methods. The objective is to prepare models to have the option to make forecasts. For instance, we could utilize existing user item associations to prepare a model to foresee the best items that a user may like the most. One advantage of these strategies is that they can recommend a bigger number of items to a bigger number of users, compared to other techniques like memory-based.
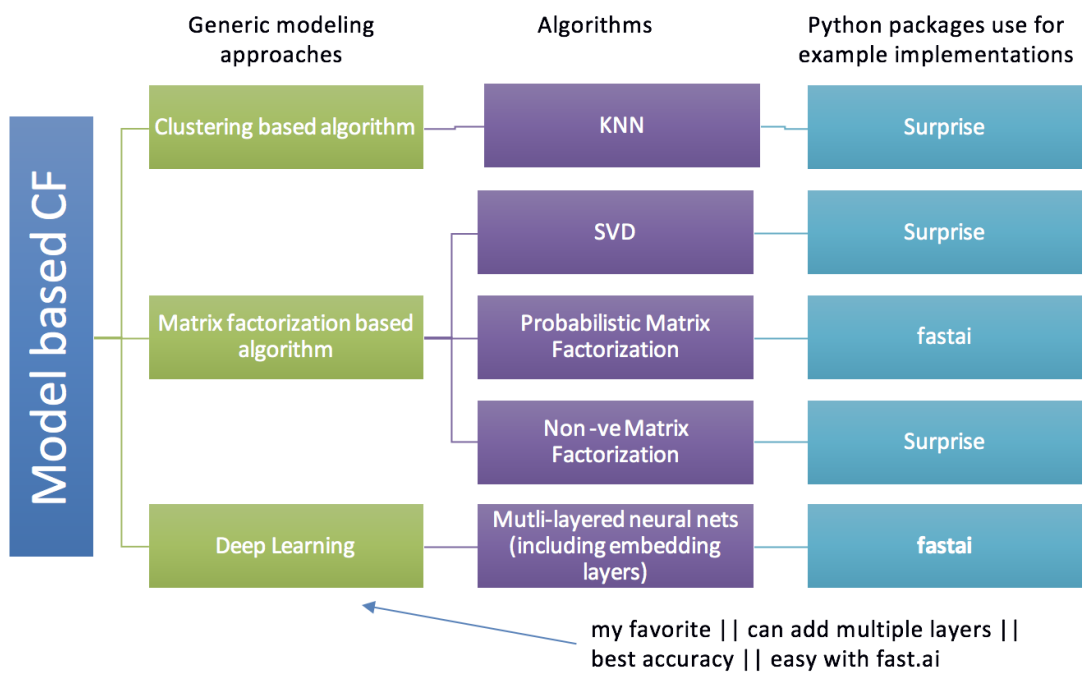


Figure 3.3 Model Based Approach, Source (Prince Gover, towardsdatascience, 2009)

BIS (Hons) Information Systems Engineering
Faculty of Information and Communication Technology (Kampar Campus), UTAR

## 3.1.2 Content Based Filtering

Content-based proposal motor works with existing profiles of users. A profile has information about a user and their taste. Taste relies upon user rating for different things. Overall, at whatever point a user makes his profile, Recommendation motor does a user survey to get starting information about the user in order to keep away from new user issue. In the recommendation strategy, the motor ponders the things that are starting at now quite assessed by the user with the things he didn't rate and looks for comparable qualities. Things like the unequivocally assessed ones will be recommended to the user. Here, considering user's taste and lead a substance based model can be worked by recommending articles noteworthy to user's taste. Such a model is powerful and redone now it needs something.



Figure 3.4 Content Based Approach, source (medium,2019)

BIS (Hons) Information Systems Engineering
Faculty of Information and Communication Technology (Kampar Campus), UTAR

### 3.1.3 Hybrid Recommender Systems

Both substance based separating and community oriented sifting have there characteristics and weaknesses. Content depiction. In specific spaces delivering an accommodating portrayal of the substance can be inconvenient. In territories where the things contain music or video for example a depiction of the substance isn't continually possible with the current advancement. Over-specialization. A substance based isolating system won't select things if the past user direct doesn't offer proof to this. Additional systems must be added to enable the structure to make proposal outside the degree of what the customer has quite recently shown excitement for. Abstract area issue. Content-based filtering methodologies experience issues in perceiving emotional information, for instance, motivations behind viewpoints and astuteness. The implementation of the system consists of machine learning algorithm as text classification that involved Naive Bayes Classifier , K Nearest Neighbor Classifier and Support Vector Machine classifier. Lexicon Based approach is used by import VADER Sentiment Analyser to analyse sentiment polarity of the tweets and identify positive, negative and neutral tweets.

### Software Tools

The tools or software in detail that will be used in the project is Anaconda 2020.02. Machine Learning Algorithm Naive Bayes Theorem, K-NN (K- Nearest Neighbor) classifier and Support Vector Machine Classifier will be used to test the accuracy among the data set produced. Lexicon Based Approach is used to analyse polarity of the tweets.

**Data set**

Data set is obtained from Twitter. AI techniques gain from models. It is critical to have great handle of information and the various terminology utilized while describing data. Here the data set used is a text dataset, where the tweets from users are obtained. The tweets are mainly reviews of Apple products from twitter. These tweet will be used to analyse polarity of the text how negative and positive it is. The another data set used Twitter reviews on self driving cars. There are many approaches that can indicate the polarity of positivity and negativity of the text. In this research research study Naive Bayes Theorem and K-NN classifier is used to measure the accuracy and performance of the approaches on the data set. The positive sentiments will be classified '1' and negative sentiments will be classified '-1'.

At the point when considering information, it is usually classified into rows and columns, similar to a database table or an Excel spreadsheet. The data set used is .csv format. This is a conventional structure for information and is what is normal in the field of AI.



Figure 3.5 Imported .csv dataset

## 3.2 System Design / Overview

### 3.2.1 Process of Text Classification

Exactly when Author talk about content grouping, Author commonly talk about the managed arrangement, which has two stages: the preparation compose and the testing stage. For the most part the preparation arrange consolidates making the named corpora informational index, pre-handling the preparation content, vectorization of the substance, and preparing of the classifier. The testing stage fuses handling of testing substance, vectorization and characterization of the testing content.

Creating Corpus. Assortment of content reliant on characterizations. Every content has a spot with one class and has been named. Now and again Author confined this data set into two different sets: the train data set and the testing data set. Pre-handling. Oust all the pointless highlights in the content, for instance, stopword, accentuation, or confused substance. This movement is huge considering the way that it will impact the preparation set of the classifier. Vectorization of Text. Change the content into vector that can be seen for PC. All content will be addressed as feature vector reliant on the highlights Author picked. Training of the Classifier. Pick one of the content characterization calculations and feed the preparation data set to the classifier to get a preparation model. Classification. Author get the preparation model, Author can deal with the testing data set into it and get the expectation of characters.
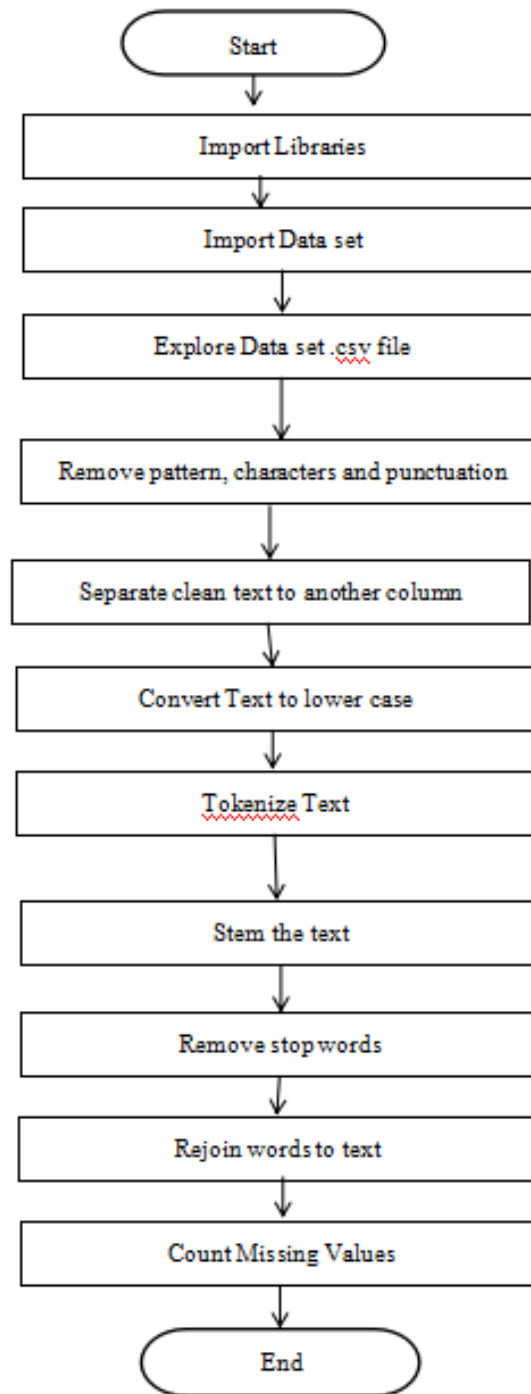
```
                          ┌─────────────┐
                          │    Start    │
                          └─────────────┘
                                 │
                                 ▼
                   ┌───────────────────────────┐
                   │      Import Libraries      │
                   └───────────────────────────┘
                                 │
                                 ▼
                   ┌───────────────────────────┐
                   │      Import Data set       │
                   └───────────────────────────┘
                                 │
                                 ▼
                   ┌───────────────────────────┐
                   │  Explore Data set .csv file│
                   └───────────────────────────┘
                                 │
                                 ▼
      ┌──────────────────────────────────────────────────┐
      │  Remove pattern, characters and punctuation       │
      └──────────────────────────────────────────────────┘
                                 │
                                 ▼
      ┌──────────────────────────────────────────────────┐
      │     Separate clean text to another column         │
      └──────────────────────────────────────────────────┘
                                 │
                                 ▼
      ┌──────────────────────────────────────────────────┐
      │         Convert Text to lower case                │
      └──────────────────────────────────────────────────┘
                                 │
                                 ▼
      ┌──────────────────────────────────────────────────┐
      │              Tokenize Text                        │
      └──────────────────────────────────────────────────┘
                                 │
                                 ▼
      ┌──────────────────────────────────────────────────┐
      │              Stem the text                        │
      └──────────────────────────────────────────────────┘
                                 │
                                 ▼
      ┌──────────────────────────────────────────────────┐
      │            Remove stop words                      │
      └──────────────────────────────────────────────────┘
                                 │
                                 ▼
      ┌──────────────────────────────────────────────────┐
      │           Rejoin words to text                    │
      └──────────────────────────────────────────────────┘
                                 │
                                 ▼
      ┌──────────────────────────────────────────────────┐
      │          Count Missing Values                     │
      └──────────────────────────────────────────────────┘
                                 │
                                 ▼
                          ┌─────────────┐
                          │     End     │
                          └─────────────┘
```

Figure 3.6 : Stages of Pre-process

### 3.2.2 Sentiment Analysis

Estimation investigation is a technique for data mining to choose the demeanor of a speaker or an essayist with respect to some subject or the general logical extremity.

Various evaluation research study strategies have been made starting late. The least demanding arrangement is to arrange content into either a positive or a negative idea order reliant on content characterization. The major philosophy is jargon based, which is to separate tweets reliant on the words that the content contains. The writings are sifted and checked if some specific opinion words are contained. It has been portrayed in a jargon that a couple of words are certain and some are negative and all of them is allocated a feeling score. The whole content will be settled subject to the score.

Notwithstanding, it is difficult to keep up a vocabulary of word reference of watchwords to evaluate conclusion score. Along these lines, some directed and unaided counts are similarly developed and used for content order, for instance, Naive Bayes, Decision Tree, K-Nearest Neighbors, Maximum Entropy, and Support Vector Machines. For these AI counts, in order to do the course of action, sufficient named data ought to be dealt with into the classifier to set up the classifier. Considering the preparation informational index, the classifier will build a likelihood model that can give a desire for next information.

In this case sentiment analysis is handled to analyse microblog data set. Now the question arise how machine learning is linked with sentiment analysis ? The essential job of Artificial Intelligence(AI) in feeling research study is to improve and mechanize the low-level content investigation works that estimation investigation. For instance, data analysts can prepare an AI model to distinguish items by taking care of it a huge volume of text document containing pre-labeled models. Utilizing

supervised and unsupervised AI systems, for example, neural systems and profound learning, the model will realize what items resemble.

When the model is prepared, similar data analysts can apply those preparation strategies towards building new models to recognize different data. The outcome is quick and reliable that helps the bigger content investigation framework recognize conclusion bearing expressions all the more successfully.

AI likewise helps data scientists tackle unique issues brought about by the development of language. Making an assessment research study rule set to represent each potential importance is impossible. Be that as it may, on the off chance that you feed an AI model with two or three thousand labeled models, it can figure out how to comprehend what the word actually means signifies with regards to video gaming, versus with regards to social media. Also, you can apply comparable training methods to comprehend other double-meaning implications too.

## 3.3 Machine Learning methods

### 3.3.1 Naive Bayes Theorem

Naive Bayes is one of the most generally perceived oversaw gathering methodologies that can be used to perform content grouping. Before that, Author need first to have an investigate what the component vector is. In order to perform characterization, Author need to pick highlights from the information first. For content arrangement, the component vector is in like manner called the term vector, which is the most huge structure during the preparation and order process. All tweets writings will be changed to term vectors to be taken care of by classifier. Generally, term vector is made reliant on a one of a one of a kind jargon, which is delivered from the preparation informational collection, and there are no duplicate words in the jargon. The size of the term vector is the size of the jargon. There are two sorts of Naive Bayes executions: Naive Bayes – Bernoulli and Naive Bayes – Multinomial. The essential difference between them is the way includes are expelled from the archives

Vector will instated with all highlights equal to zero. By then check each word in the jargon to check whether the word exists in the tweets. If it exists, by then imprint the contrasting element in the term vector with 1, if not, mark the relating highlight in the term vector to 0. As such, if the jargon is adequately huge, every single tweets can be addressed using a term vector with 0s and 1s.

Generally, in content grouping, Author will ignore the solicitation for words in the document. Or maybe, Author consider about the nearness or nonappearance of the single word, for case, whether or not a word in the jargon is incorporated for the archive or not. This model is known as a pack of words. It takes after Author hurl all words in a pack and they could be in any solicitation dealt with.

The element of term vector doesn't just speak to the nearness or nonattendance of a word, it can similarly address the recurrence of the word. Author can think each term vector as a n-measurement point in a n-estimation encourage structure, where n is the size of the jargon. For the preparation information, Author views it as a couple of classes of a great deal of n-measurement focuses. By then the content order issue becomes customary focuses characterization issue, figured the measurement might be immense.

For example, Author utilizes a report D, and Author have classes C, which contains a couple of classes. By then in order to get which class that the report D has a place, Author need to calculate the likelihood P(C|D), and picks the biggest one. P(C|D) can be handled by Bayes' Theorem:

$$P(C|D) = \frac{P(D|C)P(C)}{P(D)} \propto P(D|C)P(C)$$

Where the probability and likelihood can be handled from the marked informational collection.

Bayesian system brought under an alternate name into the content recovery network and remains a popular(baseline) technique for content arranging, the issue of settling on a choice about archives as having a spot with one arrangement or the other with word frequencies as the element. A favorable position of Naive Bayes' is that it just requires a constrained amount of preparing information to assess the parameters significant for arrangement.

Remarkably Naive Bayes' is a restrictive likelihood model. Despite its straightforwardness and solid suspicions, the Naive Bayes' classifier has been shown to work adequately in various areas. Bayesian arrangement gives reasonable learning calculations and earlier.

Figure 3.7 : Naive Bayes' Flowchart

BIS (Hons) Information Systems Engineering
Faculty of Information and Communication Technology (Kampar Campus), UTAR

**Multinomial Naive Bayes Theorem**

A multinomial dissemination is helpful to demonstrate include vectors where each worth speaks to, for instance, the quantity of events of a term or its relative recurrence. In the event that the element vectors have n components and every one of them can expect k various qualities with likelihood pk, at that point:

$$P(X_1 = x_1 \cap X_2 = x_2 \cap \ldots \cap X_k = x_k) = \frac{n!}{\prod_i x_i!} \prod_i p_i{}^{x_i}$$

The restrictive probabilities P(xi|y) are figured with a recurrence tally (which compares to applying a most extreme probability approach), however for this situation, it's essential to think about the alpha boundary (called Laplace smoothing factor) which default esteem is 1.0 and keeps the model from setting invalid probabilities when the recurrence is zero. It's conceivable to dole out all non-negative qualities, in any case, bigger qualities will allot higher probabilities to the missing highlights and this decision could change the steadiness of the model. In our model, we will consider the default estimation of 1.0.

BIS (Hons) Information Systems Engineering
Faculty of Information and Communication Technology (Kampar Campus), UTAR

### 3.3.3K- Nearest Neighbor Algorithm (K-NN)

kNN algorithm is used in both classification and regression problems. The KNN calculation accept that comparable things exist in nearness. Resulting, comparable things are close to one another. similar data centers are close to each other. The KNN estimation on this supposition that being veritable enough for the count to be significant. KNN gets the chance of likeness with some number juggling we may have learned in our childhood—discovering the partition between centers around a diagram.



Figure 3.8 Calculation of distance

Assumed the location of A and B is known. The distance is found with the formula, distance 'C' is,

$$c = \sqrt{(x_A - x_B)^2 + (y_A - y_B)^2}$$

K-NN is a sort of instance based learning, where the capacity is just approximated locally and all calculation is conceded until classification. It is non parametric strategy

BIS (Hons) Information Systems Engineering
Faculty of Information and Communication Technology (Kampar Campus), UTAR

used for order or backslide. In case of grouping the yield is class investment (the most well-known bunch may be returned).

The Nearest Neighbor rule (NN) is the least unpredictable structure of K-NN when K = 1. Given obscure model and a getting ready set, all the partitions between the obscure test and all the models in the readiness set can be figured. The separation with the littlest worth identifies with the model in the readiness set closest to the obscure model. Consequently, obscure model may be arranged reliant on the portrayal of this closest neighbor. The K-NN is a basic computation to understand and execute , and a necessary resource we have accessible to us for sentiment research study. KNN is simple of the way that it doesn't acknowledge anything about the data, other than a separation measure can be resolved dependably between two events.

Figure 3.9 K-Nearest Neighbor Classifier Flowchart

BIS (Hons) Information Systems Engineering
Faculty of Information and Communication Technology (Kampar Campus), UTAR

**Support Vector Machine**


Support Vector Machine is used to isolate the two classes of information focuses, there are numerous conceivable hyperplanes that could be picked. Our goal is to locate a plane that has the greatest edge, i.e the most extreme separation between information purposes of the two classes. Boosting the edge separation gives some support so future information focuses can be ordered with more certainty.


**Hyperplanes and Support Vectors**

Hyperplanes are choice limits that help characterize the information focuses. Information focuses falling on either side of the hyperplane can be ascribed to various classes. Additionally, the component of the hyperplane relies on the quantity of highlights. On the off chance that the quantity of info highlights is 2, at that point the hyperplane is only a line. On the off chance that the quantity of information highlights is 3, at that point the hyperplane turns into a two-dimensional plane. It gets hard to envision when the quantity of highlights surpasses 3.


Backing vectors are information focuses that are nearer to the hyperplane and impact the position and direction of the hyperplane. Utilizing these help vectors, we amplify the edge of the classifier. Erasing the help vectors will change the situation of the hyperplane. These are the focuses that assist us with building our SVM.


Margin: This is the separation between every one of the help vectors, as appeared previously. The calculation plans to boost the edge. The issue of finding the most extreme edge (and thus, the best hyperplane) is an advancement issue, and can be tackled by enhancement strategies.


Kernel: The part is a sort of capacity which is applied to the information focuses so as to plan the first non-direct information focuses into a high dimensional space where they are detachable. As a rule there won't be a direct choice limit, which implies that no single consecutive line will isolates the two classes. Pieces address this issue.
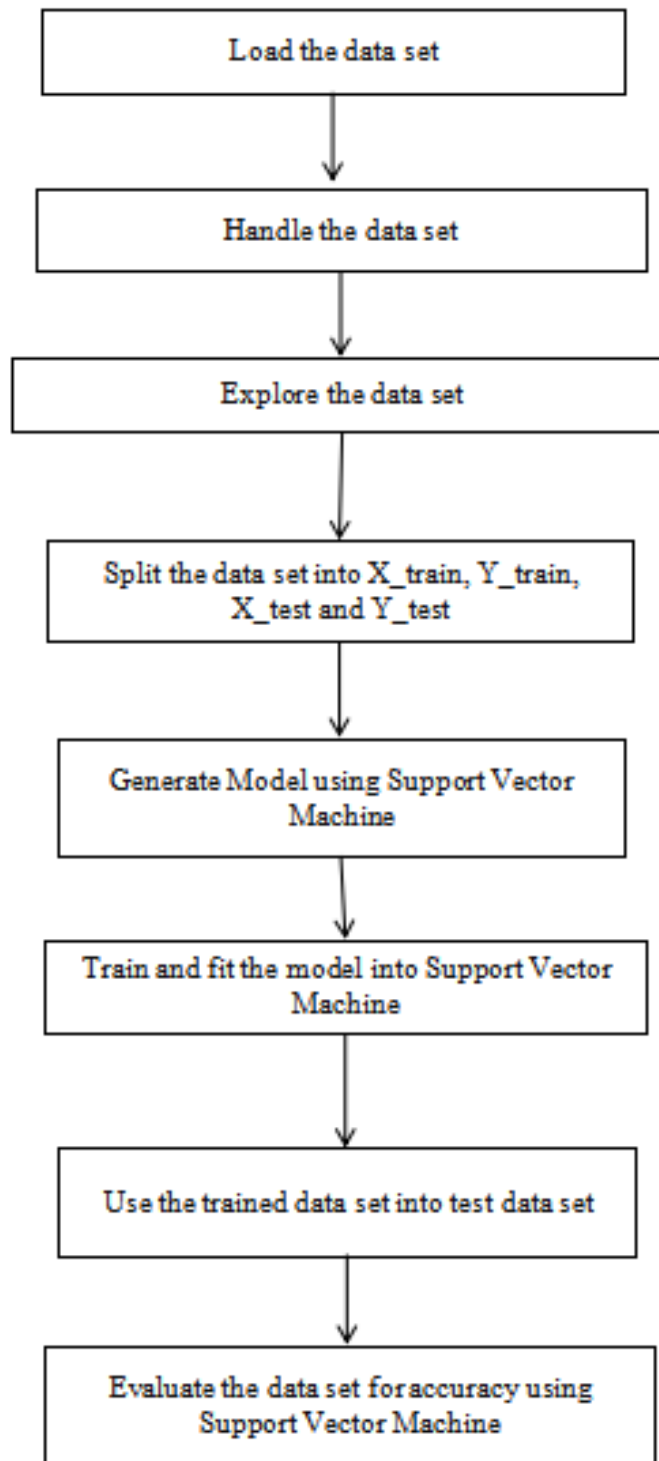
Figure 3.10 Flowchart of Support Vector Machine Classifier

BIS (Hons) Information Systems Engineering
Faculty of Information and Communication Technology (Kampar Campus), UTAR

## 3.4 VADER Sentiment Analyser

VADER (Valence Aware Dictionary and sEntiment Reasoner) is a vocabulary and rule-based notion investigation instrument that is explicitly sensitive to assessments communicated in online media. VADER utilizes a mix of A conclusion dictionary is a rundown of lexical highlights which are commonly marked by their semantic direction as either sure or negative. VADER not just tells about the Positivity and Negativity score yet additionally educates us regarding how sure or negative a slant is. Vader Sentiment analyser is lexicon based approach to analyse polarity of the tweets in the data sets. The labeled tweets is then used in the text classification. This forms a hybrid approach towards the text classification.

```
┌─────────────────────────────────┐
│     Extracted Twitter Data set  │
└─────────────────────────────────┘
                │
                ▼
┌─────────────────────────────────┐
│       Pre-Processing Steps      │
└─────────────────────────────────┘
                │
                ▼
┌─────────────────────────────────┐
│   Sentiment Analyzer using VADER│
│       Sentiment Analysis        │
└─────────────────────────────────┘
                │
                ▼
┌─────────────────────────────────┐
│    Sentiment Polarity of tweets │
└─────────────────────────────────┘
                │
                ▼
┌─────────────────────────────────┐
│  Categorize into Positive,      │
│  Negative and Neutral           │
└─────────────────────────────────┘
                │
                ▼
┌─────────────────────────────────┐
│  Label of the polarity as 1     │
│  [positive], -1[negative] and 0 │
│  for neutral                    │
└─────────────────────────────────┘
```

Figure 3.11 Vader Sentiment Analyser

## 3.5 Implementation Issues / Challenges

In a huge bit of the internet based life, language used by the users is easygoing. users make their own words and spelling alternate routes and accentuation, erroneous spellings, slang, new words, URLs and abbreviated structures. Right now kind of content solicitations to be amended. Right now inspecting the content HTML characters, slang words, emoji, stop words, accentuations, urls are ought to have been expelled. Parting of joined words are moreover be seen for purging. Users who are furthermore assessing the thing, administrations and offices Author gave by various destinations are ought to have been tended to. Various systems for looking at users conduct, sees, mentality are ought to be examined and demands to be institutionalized. Looking for right data set was also a challenger because it needs right data, sentiment polarity is observed in the research study. Finding tweets and data set with sentiments that match the polarity was difficult. As most of the data set comes with sentiment ratings with stars or even numbers that rates the scale from one to five. Obtained data set also contain irrelevant data that does not match the situation or some files could not be loaded into python as it has been polluted.

## CHAPTER 4 SYSTEM IMPLEMENTATION

### 4.1 Data Preprocessing

Data preprocessing is a data mining technique that incorporates changing unrefined data into a legitimate game plan. Genuine data is consistently inadequate, clashing, and also debilitated in explicit practices or floats, and is presumably going to contain various goofs. Data preprocessing is an exhibited procedure for settling such issues.

Content pre-processing is a stage that happens after content mining. Content information can be sourced from contrast places; content can emerge out of online books, content can be web scratched and it might likewise originate from online documentation. Content pre processing is fundamental so as to additionally control your content information. In characteristic language handling, one thing to remember is that whatever you do to the crude information may affect how your model will be prepared. For instance, stripping accentuation and evacuating upper cases may change the importance of your sentences. This is something to remember while experiencing your information and what you need to have as a final product.

In current world information commonly deficient: lacking trademark characteristics, missing the mark on explicit attributes of interest, or containing simply complete data. Uproarious: containing errors or special cases. Clashing: containing mistakes in codes or names.

**Import libraries**

```
10
11    #import libraries in |
12    import nltk
13    import numpy as np
14    import re
15    import pandas as pd
16    import pylab as pl
17    import matplotlib.pyplot as plt
18
19    from nltk.tokenize import WordPunctTokenizer
20    from bs4 import BeautifulSoup
21
22    from sklearn.feature_extraction.text import TfidfVectorizer
23    from sklearn.decomposition import PCA, TruncatedSVD
24
25    from vaderSentiment.vaderSentiment import SentimentIntensityAnalyzer
26    analyser = SentimentIntensityAnalyzer()
27
28    from sklearn import metrics
29    from mpl_toolkits.mplot3d import Axes3D
30
31    from matplotlib import pyplot
32
33
34
```

Figure 4.1.1:    Imported libraries

This is the way we import libraries in Python using import catchphrase and this is the most popular libraries which any Data Scientist used. NumPy is the focal pack for consistent figuring with Python. It contains notwithstanding different things: A unimaginable N-dimensional display object. Pandas is for data control and research study. Pandas is an open source, BSD-approved library giving tip top, easy to-use data structures and data assessment mechanical assemblies for the Python programming language. Pandas is a NumFOCUS upheld adventure. This will help ensure the accomplishment of progress of pandas as a world-class open-source undertaking, and makes it possible to provide for the endeavor.   Seaborn is a Python data portrayal library reliant on matplotlib. It gives a noteworthy level interface to drawing charming and instructive authentic structures.

Advised messages are ordinarily given in conditions where it is important to alert the customer of some condition in a program, where that condition (routinely) doesn't warrant raising an exclusion and consummation the program. For example, one should give a counsel when a program uses an old module.

**Import dataset**

Load the Data

When the libraries are imported, our following stage is to stack the information, put away in the GitHub vault connected here. You can download the information and keep it in your neighborhood envelope. From that point onward, we can utilize the read_csv strategy for Pandas to stack the information into a Pandas information outline df, as demonstrated as follows:

```
df = pd.read_csv('appletwitter.csv')
print(len(df))
```

Figure 4.1.2: Imported dataset

By utilizing Pandas author imports  informational index and the record author utilized here is .csv document.   Be that as it may, to get to and to utilize fastly author use CSV records as a result of their light loads.

```
df["_unit_id"].value_counts(ascending=False).head(2)

def remove_pattern(input_txt, pattern):
    r = re.findall(pattern, input_txt)
    for i in r:
        input_txt = re.sub(i, '', input_txt)

    return input_txt
```

Figure 4.1.3 Remove Pattern defined

Patterns that are unwanted has to be removed from the text. The query shows the text column in the dataset is chosen to remove punctuation such as !"#$%&'()*+, -./:;<=>?[\]^_`{|}~ . It wil also form another column named 'Tweet_punct'.

```
df['Clean_text'] = np.vectorize(remove_pattern)(df['text'], "@[\w]*")
```

Figure 4.1.4 Clean Text Statement

A column named clean text will be generated that has removed pattern from text column.

```
46    start = time.time()
47    df = pd.read_csv('appletwitter.csv')
48    print(len(df))
49
50    df["_unit_id"].value_counts(ascending=False).head(2)
51
52
53    def remove_pattern(input_txt, pattern):
54        r = re.findall(pattern, input_txt)
55        for i in r:
56            input_txt = re.sub(i, '', input_txt)
57
58        return input_txt
59
60
61    df['Clean_text'] = np.vectorize(remove_pattern)(df['text'], "@[\w]*")
62
63
64    df['Clean_text'] = df ['text'].str.lower()
65    example_review = df.iloc[0]
66    print(example_review['Clean_text'])
67
68
69    print(nltk.word_tokenize(example_review['Clean_text']))
70
71
```

Figure 4.1.5 Replace String

The clean text column will be analysed further to remove unwanted symbols.

BIS (Hons) Information Systems Engineering
Faculty of Information and Communication Technology (Kampar Campus), UTAR

**Sentence Tokenization**

```
76
77   ▾ def identify_token(row):
78         Clean_text=row['Clean_text']
79         tokens = nltk.word_tokenize(Clean_text)
80
81         #taken only words and not punctuation
82         token_words =[w for w in tokens if w.isalpha()]
83         return token_words
84   df['words'] = df.apply(identify_token, axis=1)
85
86
87
88   from nltk.stem import PorterStemmer
89   stemming = PorterStemmer()
90   ▾ def stem_list(row):
91         my_list = row['words']
92         stemmed_list = [stemming.stem(word) for word in my_list]
93         return(stemmed_list)
94   df['stemmed_words'] = df.apply(stem_list,axis=1)
95
96
97
98   from nltk.corpus import stopwords
99   stops = set(stopwords.words("english"))
100
101  ▾ def remove_stops(row):
102        my_list = row['stemmed_words']
103        meaningful_words = [w for w in my_list if not w in stops]
104        return (meaningful_words)
105  df['stem_meaningful'] = df.apply(remove_stops, axis=1)
106
107
```

Figure 4.1.6 Sentence Tokenization

The importance of tokenization is to hack up some current content into littler lumps. For instance, a passage can be tokenized into sentences and further into words. Author will initially begin with a straightforward string that we might want to partition into sentences. In the event that preprocessing a panda information outline, will need to circle the tokenization methodology over the entirety of your lines. This is accomplished by making a capacity. On the off chance that you are preprocessing a NumPy cluster, will need Tokenization of circle for every component in the exhibit so as to tokenize the sentences. Author will utilize word_tokenize strategy from NLTK to part the audit text into singular words (and you will see that accentuation is additionally delivered as isolated 'words').Let's gander at our model line. Apply the word_tokenize to all records, making another segment in the DataFrame.

```
97
98      from nltk.corpus import stopwords
99      stops = set(stopwords.words("english"))
100
101   ▼ def remove_stops(row):
102         my_list = row['stemmed_words']
103         meaningful_words = [w for w in my_list if not w in stops]
104         return (meaningful_words)
105     df['stem_meaningful'] = df.apply(remove_stops, axis=1)
106
107
108   ▼ def rejoin_words(row):
109         my_list = row['stem_meaningful']
110         joined_words = (" ".join(my_list))
111         return joined_words
112     df['Tweets_clean'] = df.apply(rejoin_words, axis=1)
113
114     missing_values_count=df.isnull().sum()
115     missing_values_count[0:13]
116
117
118
119     #To determine length of clean text
120     df['Clean_text_length'] = df['Tweets_clean'].apply(len)
121
122
```

Figure 4.1.7 Text Stemming, Remove stop words and rejoin words

**Stemming**

Stemming lessens related words to a typical stem. It is a discretionary cycle step, and it is valuable to test exactness with and without stemming. We should take a gander at a model.

**Removing stop words**

'Stop words' are normally utilized words that are probably not going to have any advantage in characteristic language preparing. These incorporates words, for example, 'a', 'the', 'is'. As before author characterize a capacity and apply it to our DataFrame. Author make a lot of words that we will call 'quits' (utilizing a set assists with accelerating eliminating stop words).

**Rejoin words**

The meaningful words will be concatenated to a string

BIS (Hons) Information Systems Engineering
Faculty of Information and Communication Technology (Kampar Campus), UTAR

## 4.2 Sentiment Analysis

Vader sentiment analyser is imported in this case. The Python code for the standard based notion investigation motor. Actualizes the linguistic and grammatical guidelines portrayed in the research study, consolidating studyly determined measurements for the effect of each standard on the apparent force of slant in sentence-level content. Significantly, these heuristics go past what might ordinarily be caught in a regular sack of-words model. They fuse **word-order sensitive relevants** between terms. For instance, degree modifiers (likewise called intensifiers, sponsor words, or degree qualifiers) sway slant force by either expanding or diminishing the power.

```python
124    from textblob import TextBlob
125
126  ▾ def calculate_sentiment(Tweets_clean):
127        return TextBlob(Tweets_clean).sentiment
128
129  ▾ def calculate_sentiment_analyser(Tweets_clean):
130        return analyser.polarity_scores(Tweets_clean)
131
132    df['sentiment']=df.Tweets_clean.apply(calculate_sentiment)
133    df['sentiment_analyser']=df.Tweets_clean.apply(calculate_sentiment_analyser)
134
135
136    s = pd.DataFrame(index = range(0,len(df)),columns= ['compound_score','compound_score_sentiment'])
137
138  ▾ for i in range(0,len(df)):
139      s['compound_score'][i] = df['sentiment_analyser'][i]['compound']
140
141  ▾   if (df['sentiment_analyser'][i]['compound'] <= -0.05):
142        s['compound_score_sentiment'][i] = 'Negative'
143  ▾   if (df['sentiment_analyser'][i]['compound'] >= 0.05):
144        s['compound_score_sentiment'][i] = 'Positive'
145  ▾   if ((df['sentiment_analyser'][i]['compound'] >= -0.05) & (df['sentiment_analyser'][i]['compound'] <= 0.05)):
146        s['compound_score_sentiment'][i] = 'Neutral'
147
148    df['compound_score'] = s['compound_score']
149    df['compound_score_sentiment'] = s['compound_score_sentiment']
150    df.head()
151
152
153    count=df['compound_score_sentiment'].value_counts()
154    print(count)
155
156    #drop neutral tweets
157    df = df.set_index("compound_score_sentiment")
158    df = df.drop("Neutral", axis=0)
159
160
```

Figure 4.1.8 Sentiment Analyser codes

The figure above shows the for loop used in the calculations of sentiment analyser.

BIS (Hons) Information Systems Engineering
Faculty of Information and Communication Technology (Kampar Campus), UTAR

**Data Splitting**

Train_test_split is a capacity in Sklearn model determination for parting information clusters into two subsets: for preparing information and for testing information. With this capacity, you do not have to isolate the dataset physically.

Of course, Sklearn train_test_split    make irregular allotments for the two subsets. In any case,    can likewise determine an arbitrary state for the activity

```
df['label'] = df['compound_score'].apply(lambda x: -1 if x < 0 else 1)

#View total number of rows in df
rows=len(df.axes[0])
print('Total number of positive and negative', rows)

#Count numbers of positive and negative tweets
count=df['label'].value_counts()
print(count)

#spliiting dataset into train set and test set
from sklearn.model_selection import train_test_split
X_train, X_test, y_train, y_test = train_test_split(df['Tweets_clean'], df['label'], test_size=0.2, random_st
```

Figure 4.1.9 Splitting data set

The figure shows the query to split the dataset into train set and test set.

**Parameters**

X, y. The essential parameter is the dataset you're deciding to use. train_size. This parameter sets the size of the planning dataset. There are other options: None, which is the default, Int, which requires the particular number of tests, and float, which ranges from 0.1 to 1.0.

test_size. This parameter demonstrates the size of the testing dataset. The default state tests the readiness size. It will be set to 0.25 if the planning size is set default.

random_state. The default mode plays out an unpredictable split using np.random. On the other hand, you can incorporate an entire number using a cautious number.

BIS (Hons) Information Systems Engineering
Faculty of Information and Communication Technology (Kampar Campus), UTAR

**Output of Preprocessing Steps**



Figure 4.1.10 Visualised dataframe

The dataset that has been imported is viewed as shown above where is has column on id, Query and tweets.

Figure 4.1.11 Column of tweets

The columns of text and clean text is shown in figure labelled 'G1' and 'G2'.



Figure 4.1.12 Tokenized words

The labels 'G3', 'G4' and 'G5' shows tokenized words, stemmed text and removing stop words

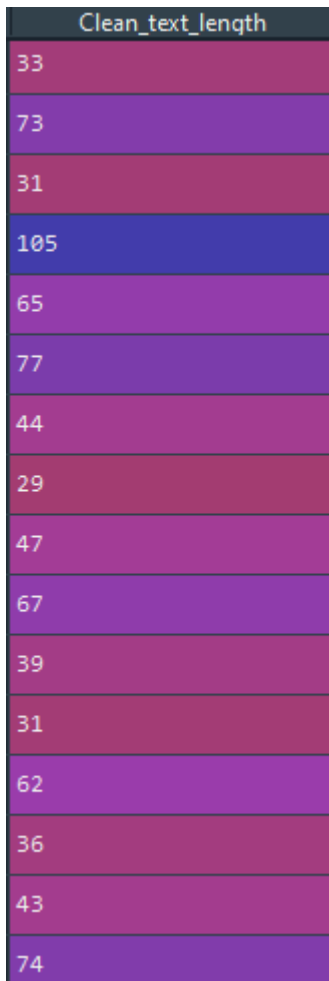Figure 4.1.13 Clean text length

The length of clean text length is visualised in a separate column as above.

BIS (Hons) Information Systems Engineering
Faculty of Information and Communication Technology (Kampar Campus), UTAR

| sentiment | sentiment_analyser | compound_score | label |
|---|---|---|---|
| Sentiment(polarity=~~~~~~~~~~, subjectivity=0.55) | {'neg': 0.181, 'neu': 0.601,… | 0.1265 | 1 |
| Sentiment(polarity=-0.4, subjectivity=0.6) | {'neg': 0.888, 'neu': 0.112,… | -0.836 | -1 |
| Sentiment(polarity=-0.4, subjectivity=0.6) | {'neg': 0.538, 'neu': 0.462,… | -0.5423 | -1 |
| Sentiment(polarity=-0.14583333333333334, subjectivity=0.6458333333333334) | {'neg': 0.674, 'neu': 0.326,… | -0.8839 | -1 |
| Sentiment(polarity=0.0, subjectivity=0.0) | {'neg': 0.297, 'neu': 0.703,… | -0.5859 | -1 |
| Sentiment(polarity=-0.30000000000000004, subjectivity=0.7) | {'neg': 0.415, 'neu': 0.585,… | -0.7964 | -1 |
| Sentiment(polarity=-0.4333333333333333, subjectivity=0.6166666666666667) | {'neg': 0.38, 'neu': 0.526, … | -0.8225 | -1 |
| Sentiment(polarity=-0.8, subjectivity=0.9) | {'neg': 0.32, 'neu': 0.608, … | -0.6705 | -1 |
| Sentiment(polarity=0.2857142857142857, subjectivity=0.5357142857142857) | {'neg': 0.0, 'neu': 0.839, '… | 0.3612 | 1 |
| Sentiment(polarity=0.0, subjectivity=0.0) | {'neg': 0.333, 'neu': 0.667,… | -0.128 | -1 |
| Sentiment(polarity=0.4, subjectivity=0.8) | {'neg': 0.157, 'neu': 0.597,… | 0.296 | 1 |
| Sentiment(polarity=-0.25, subjectivity=0.18333333333333332) | {'neg': 0.304, 'neu': 0.696,… | -0.5423 | -1 |
| Sentiment(polarity=0.0, subjectivity=0.0) | {'neg': 0.247, 'neu': 0.753,… | -0.3182 | -1 |
| Sentiment(polarity=0.0, subjectivity=0.0) | {'neg': 0.0, 'neu': 0.769, '… | 0.4588 | 1 |
| Sentiment(polarity=0.8, subjectivity=0.75) | {'neg': 0.0, 'neu': 0.5, 'po… | 0.7906 | 1 |
| Sentiment(polarity=1.0, subjectivity=0.3) | {'neg': 0.0, 'neu': 0.543, '… | 0.6369 | 1 |
| Sentiment(polarity=-0.5, subjectivity=1.0) | {'neg': 0.279, 'neu': 0.721,… | -0.4767 | -1 |
| Sentiment(polarity=-0.26666666666666666, subjectivity=0.7222222222222223) | {'neg': 0.358, 'neu': 0.498,… | -0.6486 | -1 |
| Sentiment(polarity=-0.1, subjectivity=0.6) | {'neg': 0.649, 'neu': 0.351,… | -0.8126 | -1 |
| Sentiment(polarity=-0.2, subjectivity=0.8) | {'neg': 0.342, 'neu': 0.658,… | -0.6369 | -1 |
| Sentiment(polarity=0.2, subjectivity=0.2) | {'neg': 0.0, 'neu': 0.857, '… | 0.3612 | 1 |

Figure 4.1.14 Sentiment polarity

The sentiment polarity of the text is analysed as shown in sentiment , sentiment analyser has further analysed the polarity of each how postive or how negative the texts are.

Compound score is used it is a metric that calculates the total of all the lexicon reating which have been normalized between -1 , that is the most negative and +1 the most postive. The compound that is more than 0.05 or equals to 0.05 indicates the positive sentiment.

Figure 4.1.15    Labels of the text

The image shows the positive, negative and neutral labels of the tweets labeled 'G6'.

## 4.2 Machine Learning Algorithm

## Naive Bayes Classifier

```
175
176
177     #Multinomial Naive Bayes
178
179     from sklearn.feature_extraction.text import CountVectorizer
180     from sklearn import metrics
181     from sklearn.naive_bayes import MultinomialNB
182
183
184
185     cv = CountVectorizer(strip_accents='ascii', token_pattern = u'(?ui)\\b\\w*[a-z]+\\w*\\b',
186                         lowercase=True, stop_words='english')
187     X_train_cv = cv.fit_transform(X_train)
188     X_test_cv = cv.transform(X_test)
189
190
191     word_freq_df = pd.DataFrame (X_train_cv.toarray(), columns = cv.get_feature_names())
192     top_words_df = pd.DataFrame(word_freq_df.sum()).sort_values(0,ascending=False)
193
194     naive_bayes = MultinomialNB()
195     naive_bayes.fit(X_train_cv, y_train)
196     predictions = naive_bayes.predict(X_test_cv)
197
198     from sklearn.metrics import accuracy_score, precision_score, recall_score
199
200     print('Accuracy score of Naive Bayes: ', accuracy_score(y_test, predictions))
201     print('Precision score Naive Bayes: ', precision_score(y_test, predictions))
202     print('Recall score Naive Bayes: ',recall_score(y_test, predictions))
203
204     print('Accuracy score of Naive Bayes in percentage: ',
205             accuracy_score(y_test, predictions)*100)
```

Figure 4.2.1 Naive Bayes classifier model

```
206
207     from sklearn.metrics import confusion_matrix
208     print (confusion_matrix(y_test, predictions))
209
210     from sklearn.metrics import classification_report
211     print(classification_report(y_test,predictions))
212
213
214
215     time.sleep(1)
216
217     end = time.time()
218
219     print(f"Runtime is {end - start}")
220
221
222
```

Figure 4.2.2 Accuracy metrics of Naive Bayes classifier

Presently we should investigate what we simply did. To begin with, we imported the Multimomia NB strategy and the accuracy_score technique. At that point we made an article naive_bayes of the sort MultinomialNB. After this, we prepared the calculation

BIS (Hons) Information Systems Engineering
Faculty of Information and Communication Technology (Kampar Campus), UTAR

on the testing data(data_train) and testing target(target_train) utilizing the fit() technique, and afterward anticipated the objectives in the test information utilizing the foresee() strategy. At long last we printed the score utilizing the accuracy_score() strategy and with this we have effectively applied the Naive-Bayes calculation to manufacture a forecast model.

**Output of Multinomial Naive Bayes Classifier**



```
Total number of positive and negative 2050
 1    1051
-1     999
Name: label, dtype: int64
Accuracy score of Naive Bayes:  0.8658536585365854
Precision score Naive Bayes:  0.8297872340425532
Recall score Naive Bayes:  0.9285714285714286
Accuracy score of Naive Bayes in percentage:  86.58536585365853
[[160  40]
 [ 15 195]]
              precision    recall  f1-score   support

          -1       0.91      0.80      0.85       200
           1       0.83      0.93      0.88       210

    accuracy                           0.87       410
   macro avg       0.87      0.86      0.86       410
weighted avg       0.87      0.87      0.87       410

Runtime is 11.466231346130371
```

Figure 4.2.3 Output of Naive Bayes classifier

Output of Naive Bayes Classifier of data set 1 shows the accuracy of 86.6%, precision score of 0.829 and recall score of 0.929

BIS (Hons) Information Systems Engineering
Faculty of Information and Communication Technology (Kampar Campus), UTAR

**K Nearest Neighbor Classifier**

KNN (K-Nearest Neighbor) is a basic administered order calculation we can use to dole out a class to new information point. It tends to be utilized for relapse also, KNN doesn't make any presumptions on the information appropriation, consequently it is non-parametric. It keeps all the preparation information to make future expectations by processing the closeness between an info test and each preparation occurrence.

KNN can be summed up as underneath: Registers the separation between the new information point with each preparation model. For registering the separation estimates, for example, Euclidean separation, Hamming separation or Manhattan separation will be utilized. Model picks K sections in the information base which are nearest to the new information point. At that point it does the greater part vote i.e the most widely recognized class/name among those K passages will be the class of the new information point

```
169
170    #K Nearest Neighbor Classifier
171
172    from sklearn.feature_extraction.text import CountVectorizer
173  ▼ cv = CountVectorizer(strip_accents='ascii', token_pattern = u'(?ui)\\b\\w*[a-z]+\\w*\\b',
174                        lowercase=True, stop_words='english')
175    X_train_cv = cv.fit_transform(X_train)
176    X_test_cv = cv.transform(X_test)
177
178
179    word_freq_df = pd.DataFrame (X_train_cv.toarray(), columns = cv.get_feature_names())
180    top_words_df = pd.DataFrame(word_freq_df.sum()).sort_values(0,ascending=False)
181
182
183    from sklearn.neighbors import KNeighborsClassifier
184    classifier = KNeighborsClassifier(n_neighbors=5)
185    classifier.fit(X_train_cv, y_train)
186
187    y_pred = classifier.predict(X_test_cv)
188
```

Figure 4.2.4 Model of KNN classifier

Scikit-learn is composed into modules, with the goal that we can import the applicable classes without any problem. Import the class 'KNeighborsClassifer' from 'neighbors' module and Instantiate the ('assessor' is scikit-learn's term for a model). We are calling model as assessor on the grounds that their essential job is to gauge obscure amounts.

BIS (Hons) Information Systems Engineering
Faculty of Information and Communication Technology (Kampar Campus), UTAR

In our model we are making an occasion ('knn' ) of the class 'KNeighborsClassifer', at the end of the day we have made an article called 'knn' which realizes how to do KNN characterization once we give the information. The boundary 'n_neighbors' is the tuning boundary/hyper boundary (k) . All different boundaries are set to default esteems. 'fit' strategy is utilized to prepare the model on preparing information (X_train,y_train) and 'anticipate' technique to do the testing on testing information (X_test). Picking the ideal estimation of K is basic, so we fit and test the model for various qualities for K (from 1 to 25) utilizing a for circle and record the KNN's trying precision in a variable (scores).

```
191
192     print('Accuracy score of K Nearest Neighbor: ', accuracy_score(y_test, y_pred))
193     print('Precision score K Nearest Neighbor: ', precision_score(y_test, y_pred))
194     print('Recall score K Nearest Neighbor: ',recall_score(y_test, y_pred))
195
196     print('Accuracy score of K Nearest Neighbor: ', accuracy_score(y_test, y_pred)*100)
197
198     from sklearn.metrics import confusion_matrix
199     print (confusion_matrix(y_test, y_pred))
200
201     from sklearn.metrics import classification_report
202     print(classification_report(y_test, y_pred))
203
204
205     time.sleep(1)
206
207     end = time.time()
208
209     print(f"Runtime is {end - start}")
```

Figure 4.2.5 Accuracy evaluation of KNN classifier

Plot the connection between the estimations of K and the comparing testing precision utilizing the matplotlib library. As should be obvious there is a raise and fall in the exactness and it is very normal while looking at the model multifaceted nature with the precision. All in all as the estimation of K increment there gives off an impression of being a raise in the exactness and again it falls.

When all is said in done the Training precision ascends as the model intricacy increments, for KNN the model multifaceted nature is dictated by the estimation of K. Bigger K esteem prompts smoother choice limit (less intricate model). Littler K prompts more mind boggling model (may prompt overfitting). Testing precision punishes models that are too complex(over fitting) or not mind boggling enough(underfit).

BIS (Hons) Information Systems Engineering
Faculty of Information and Communication Technology (Kampar Campus), UTAR

**Output of K Nearest Neighbor**

```
Name: label, dtype: int64
Accuracy score of K Nearest Neighbor:  0.7365853658536585
Precision score K Nearest Neighbor:  0.8642857142857143
Recall score K Nearest Neighbor:  0.5761904761904761
Accuracy score of K Nearest Neighbor:  73.65853658536585
[[181  19]
 [ 89 121]]
              precision    recall  f1-score   support

          -1       0.67      0.91      0.77       200
           1       0.86      0.58      0.69       210

    accuracy                           0.74       410
   macro avg       0.77      0.74      0.73       410
weighted avg       0.77      0.74      0.73       410

Runtime is 11.907018423080444
```

Figure 4.2.6 Accuracy output of KNN Classifier

Accuracy score of K Nearest Neighbor Classifier is depicted as shown, 73.7%

**Support Vector Machine**

```
165
166    #Support Vector Machine
167
168    from sklearn.feature_extraction.text import CountVectorizer
169
170  ▾ cv = CountVectorizer(strip_accents='ascii', token_pattern = u'(?ui)\\b\\w*[a-z]+\\w*\\b',
171                        lowercase=True, stop_words='english')
172    X_train_cv = cv.fit_transform(X_train)
173    X_test_cv = cv.transform(X_test)
174
175
176    word_freq_df = pd.DataFrame (X_train_cv.toarray(), columns = cv.get_feature_names())
177    top_words_df = pd.DataFrame(word_freq_df.sum()).sort_values(0,ascending=False)
178
179    from sklearn import svm
180
181    SupportVectorMachine = svm.SVC(kernel='linear')
182    SupportVectorMachine.fit(X_train_cv, y_train)
183    predictions = SupportVectorMachine.predict(X_test_cv)
184
```

Figure 4.2.7 Support Vector Machine classifier model

After we characterized the model above we have to prepare the model utilizing the information given. For this we are utilizing the fit() strategy as appeared previously.

BIS (Hons) Information Systems Engineering
Faculty of Information and Communication Technology (Kampar Campus), UTAR

Author make a model variable and launch the SVC class. After this we will prepare the model, yet before that let us examine a portion of the significant boundaries of the help vector classifier model, recorded beneath.

**Straight Kernel** is one of the most usually utilized bits. This is utilized when the information is Linearly distinct methods information can be isolated utilizing a solitary Line.

RBF portion is utilized when the information isn't directly distinct. The RBF bit of Support Vector Machine makes a non-direct mixes of given highlights and changes the given information tests to a higher dimensional element space can utilize a straight choice limit to isolate the classes.

**Regularization C**: C is a regularization boundary. The regularization is contrarily relative to C. It must be positive.

**Degree**: Degree is the level of the polynomial bit work. It is overlooked by all different portions like straight.

**Verbose**: This empowers verbose yield. Verbose is an overall programming term for produce the vast majority of logging yield. Verbose methods requesting that the program inform everything regarding what it is doing constantly.

**Random_State**: random_state is the seed utilized by the arbitrary number generator. This is utilized to guarantee reproducibility. As it were, to get a deterministic conduct during fitting, random_state must be fixed

When the model is prepared, it's prepared to make expectations. We can utilize the foresee strategy on the model and pass x_test as a boundary to get the yield as y_pred. Notice that the expectation yield is a variety of genuine numbers comparing to the information exhibit.

```
185
186    from sklearn.metrics import accuracy_score, precision_score, recall_score
187
188    print('Accuracy score of Support Vector Machine: ', accuracy_score(y_test, predictions))
189    print('Precision score Support Vector Machine: ', precision_score(y_test, predictions))
190    print('Recall score Support Vector Machine: ',recall_score(y_test, predictions))
191
192  ▼ print('Accuracy score of Support Vector Machine in percentage: ',
193         accuracy_score(y_test, predictions)*100)
194
195    from sklearn.metrics import confusion_matrix
196    print (confusion_matrix(y_test, predictions))
```

Figure 4.2.8    model testing of Support Vector Machine

```
195    from sklearn.metrics import confusion_matrix
196    print (confusion_matrix(y_test, predictions))
197
198    from sklearn.metrics import classification_report
199    print(classification_report(y_test,predictions))
200
201
202
203    time.sleep(1)
204
205    end = time.time()
206
207    print(f"executed time: {end-start}")
```

Figure 4.24 Evaluation of Support Vector Machine

When the model is prepared, it's prepared to make expectations. We can utilize the foresee strategy on the model and pass x_test as a boundary to get the yield as y_pred. Notice that the expectation yield is a variety of genuine numbers comparing to the information exhibit.

**Output of Support Vector Machine**

```
Accuracy score of Support Vector Machine:  0.9146341463414634
Precision score Support Vector Machine:  0.8854625550660793
Recall score Support Vector Machine:  0.9571428571428572
Accuracy score of Support Vector Machine in percentage:  91.46341463414635
[[174  26]
 [  9 201]]
              precision    recall  f1-score   support

         -1       0.95      0.87      0.91       200
          1       0.89      0.96      0.92       210

   accuracy                           0.91       410
  macro avg       0.92      0.91      0.91       410
weighted avg       0.92      0.91      0.91       410

executed time: 11.874178409576416
```

Figure 4.2.9 Output of Support Vector Machine

The output depicts accuracy score of Support Vector Machine Classifier 91.5%

BIS (Hons) Information Systems Engineering
Faculty of Information and Communication Technology (Kampar Campus), UTAR

## CHAPTER 5 SYSTEM TESTING

The Machine Learning algorithms were tested using two different data sets.

Dataset 1

| Dataframe | 3886 tweets |
|-----------|-------------|
| Positive | 1836 |
| Negative | 1051 |
| Neutral | 999 |

Table 5.1 Dataframe of dataset 1



| | A | B | C | D | E | F | G | H | I | J | K | L | M | N |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 7 | 623495518 | #AAPL OR | #AAPL:5 Rocket Stocks to Buy for December Gains: Apple and More...http://t.co/eG5XhXdLLS | | | | | | | | | | | |
| 8 | 623495519 | #AAPL OR | Top 3 all @Apple #tablets. Damn right! http://t.co/RJiGn2JUuB | | | | | | | | | | | |
| 9 | 623495520 | #AAPL OR | CNBCTV: #Apple's margins better than expected? #aapl http://t.co/7geVrtOGLK | | | | | | | | | | | |
| 10 | 623495521 | #AAPL OR | Apple Inc. Flash Crash: What You Need to Know http://t.co/YJIgtifdAj #AAPL | | | | | | | | | | | |
| 11 | 623495522 | #AAPL OR | #AAPL:This Presentation Shows What Makes The World's Biggest Tech Companies ...http://t.co/qlH9PqSoSd | | | | | | | | | | | |
| 12 | 623495523 | #AAPL OR | WTF MY BATTERY WAS 31% ONE SECOND AGO AND NOW IS 29% WTF IS THIS @apple | | | | | | | | | | | |
| 13 | 623495524 | #AAPL OR | Apple Watch Tops Search Engine List of Best Wearable Tech http://t.co/LTEzJzqqF8 #AAPL #iWatch #AppleWatch | | | | | | | | | | | |
| 14 | 623495525 | #AAPL OR | The Best-Designed #iPhone #Apps In the World, According to @apple: http://t.co/Razqvpxofg http://t.co/ev7uKWiEcz | | | | | | | | | | | |
| 15 | 623495526 | #AAPL OR | RT @peterpham: Bought my @AugustSmartLock at the @apple store..pretty good logo match . can't wait to install it! http://t.co/z8VKMhbnR | | | | | | | | | | | |
| 16 | 623495527 | #AAPL OR | @apple Contact sync between Yosemite and iOS8 is seriously screwed up. It used to be much more stable in the past. #icloud #isync | | | | | | | | | | | |
| 17 | 623495528 | #AAPL OR | #aapl @applenws Thanks to the non factual dumb Twitter followers stock drops 3 points in one minute. Thanks dummies. #rumors | | | | | | | | | | | |
| 18 | 623495529 | #AAPL OR | WARNING IF YOU BUY AN IPHONE 5S UNLOCKED FROM @APPLE IPHONE YOU CANNOT USE IT ON VERIZON NETWORK | | | | | | | | | | | |
| 19 | 623495530 | #AAPL OR | @Apple John Cantlie has been a prisoner of ISIS for 739 days, show you have not abandoned him. Sign https://t.co/WTn4fuiJ0P | | | | | | | | | | | |
| 20 | 623495531 | #AAPL OR | @apple- thanks for xtra checkin at upper westside store- but why are appointments running almost 50 minutes late? | | | | | | | | | | | |
| 21 | 623495532 | #AAPL OR | Why #AAPL Stock Had a Mini-Flash Crash Today: Money Morning: Nothing the analysts suggested would make a widel... http://t.co/jFGsSy2Ei | | | | | | | | | | | |
| 22 | 623495533 | #AAPL OR | $AAPL dip only momentarily....just an aberration in the equity world..#AAPL | | | | | | | | | | | |
| 23 | 623495534 | #AAPL OR | The JH Hines Staff with their newly issued @apple #ConnectED Macbook and iPad mini #txed http://t.co/82YjiCJBxH | | | | | | | | | | | |
| 24 | 623495535 | #AAPL OR | @robconeybeer: You need an IP portfolio to defend against big companies - just look at @Samsung @Apple court battles | | | | | | | | | | | |
| 25 | 623495536 | #AAPL OR | @Apple, For the love of GAWD, CENTER the '1'on the damn calendar app. You fixed it once, its back, off center, AGAIN! http://t.co/dMyAHEm | | | | | | | | | | | |
| 26 | 623495537 | #AAPL OR | i get the storage almost full notification literally every 5 minutes chill @apple | | | | | | | | | | | |
| 27 | 623495538 | #AAPL OR | I had to do made the #switch from iPhone 6 to the galaxy note edge. @apple keep up http://t.co/1Vve1htP0n | | | | | | | | | | | |
| 28 | 623495539 | #AAPL OR | @ me RT @101Baemations: Can't stand those ppl with @Apple stickers everywhere. 9/10 they prob just bought an iPod shuffle | | | | | | | | | | | |
| 29 | 623495540 | #AAPL OR | Justice Department cites 18th century federal law to get @Apple to unlock iPhones: http://t.co/Eth0QpAIom | | | | | | | | | | | |
| 30 | 623495541 | #AAPL OR | Latest Apple Products Leading in Efficiency http://t.co/KHeNlVT1FJ @apple #iPhone #iPad #plugloads | | | | | | | | | | | |

Figure 5.1 Dataset 1

The above image depicts a part of tweets of .csv file of Twitter Apple Review.

Dataset 2

| Dataframe | 7156 tweets |
|-----------|-------------|
| Positive | 3871 |
| Negative | 2293 |
| Neutral | 992 |

Table 5.2 Dataframe of dataset 2

| 7060 | 724328216 | 7032 | Trying to figure out why folks assume Apple working on a car means Apple is working on a self-driving car. Not saying wrong just confused. |
|------|-----------|------|---|
| 7061 | 724328217 | 7033 | @BenedictEvans @a16z I think a big piece of the puzzle will be creating the "App Store" of the car once they become self driving. |
| 7062 | 724328218 | 7034 | My dad's Mexican cowboy looking self driving my tiny car and not his truck is just a really funny image |
| 7063 | 724328219 | 7035 | Volvo to launch self-driving pilot program in 2017 http://t.co/mH0T7LXGK4 |
| 7064 | 724328220 | 7036 | https://t.co/9KkMurHxjC Today Volvo announced a real on-the-streets test of its self-driving cars ‰ÛÓ by 2017. |
| 7065 | 724328221 | 7037 | Thinking about how cars of the future will prevent road rage. #bigdata #analytics #robotics #selfdrivingcars |
| 7066 | 724328222 | 7041 | @semil ride sharing eventually becomes self driving cars. Roads become bandwidth. What's a good analogy? Mobile spectrum? |
| 7067 | 724328223 | 7044 | How is it that we live in a time where we're developing hoverboards and self driving cars yet we still can't send a non-shitty-looking fax? |
| 7068 | 724328224 | 7045 | Unmarked #SelfDrivingCars experiment spotted on the way to work. #SiliconValley http://t.co/D9RwFSME80 |
| 7069 | 724328225 | 7046 | Forget self-driving cars I'm still waiting for the technology that flashes headlights when you honk the horn like Homer Simpson's car. |
| 7070 | 724328226 | 7047 | Self-driving cars smart street furniture: 76 Designs of the Year nominated http://t.co/WUSHa3tp8y |
| 7071 | 724328227 | 7048 | Volvo Will Test Self-Driving Cars With Real Customers in 2017 | WIRED http://t.co/9qXPgbnzRl |
| 7072 | 724328228 | 7049 | I don't want a self-driving car I love driving. Let's build something to override human stupidity so we can avoid most of car collisions. |
| 7073 | 724328229 | 7051 | "How'd you find us?" (the courier in the self-driving car) "I just go to where the cars tell me." - Hot Tub Time Machine 2 |
| 7074 | 724328230 | 7052 | Cashmore said self-driving cars are "driving" the lifestyle in the US. #AccidentalPun #SMWMashable #SMWNYC |
| 7075 | 724328231 | 7053 | Getting spooked by this whole self driving cars and AI convo w/ @petecashmore #SMWNYC #smwmashable |
| 7076 | 724328232 | 7055 | Hmm. ‰ÛÏ@CNET: Nokia knocks net neutrality: self-driving cars "won't get the service you need" http://t.co/9axkYttaeW http://t.co/GnrfUg |
| 7077 | 724328233 | 7057 | Price? Whatever it is double it. Volvo to release self-driving cars by 2017 http://t.co/lliJNRHFCe #sopumped #volvo #designateddriver |
| 7078 | 724328234 | 7058 | Nokia knocks Net neutrality: Self-driving cars 'won't get the service you need' -‰Û_ http://t.co/sb0yUgAJSu #business http://t.co/9ha0TGRhl |
| 7079 | 724328235 | 7059 | Re: news of $NXPI acquiring $FSL from the Carlyle Group (& similar "old money VCs") new chip for self-driving cars: http://t.co/NzBIzqeNtB |
| 7080 | 724328236 | 7060 | @EdGL @jasonstory @brycebstory Volvo self-driving car to hit Swedish roads @automotive_news http://t.co/W5yAxdjoVz‰Û▣ |
| 7081 | 724328237 | 7062 | MWC: Nokia CEO cites self-driving cars and home healthcare to attack net neutrality http://t.co/GW1QWvSHkC |
| 7082 | 724328238 | 7063 | Stumbled upon a photo shoot for this self-driving mercedesbenz prototype @ Twin Peaks Vista Point https://t.co/8GXQ0P9L6U |
| 7083 | 724328239 | 7064 | Sounds like NVIDIA's new GPU Will have big implications for self-driving cars #GTC15 |

Figure 5.2 Dataset 2

The image visualizes the reviews in Twitter on self driving cars in .csv file format.

# CHAPTER 5 SYSTEM TESTING

**Tested Steps**

| Steps | Description | Result Dataset1 | Result Dataset2 | Verified by |
|---|---|---|---|---|
| **Preprocessing steps** | can utilize the read_csv strategy for Pandas to stack the information into a Pandas. have to do some pre-handling steps. | √ | √ | Aishnivya |
| Remove Pattern | Charaters, patterns and punctuations are removed. | √ | √ | Aishnivya |
| Tokenized Tweets | the way toward fragmenting text into words, sentences | √ | √ | Aishnivya |
| Stemming | lessening related words to a typical stem. | √ | √ | Aishnivya |
| Remove Stopwords | expulsion of usually utilized words probably not going to be valuable for learning. | √ | √ | Aishnivya |
| Rejoin Words | meaningful words will be concatenated to a string | √ | √ | Aishnivya |
| **Sentiment Analyser** | Sentiment Polarity was analysed using Sentiment Analyser | √ | √ | Aishnivya |

| | | | | |
|---|---|---|---|---|
| **Splitting Dataset** | Data set split into train set and test set using Scikit Learn methodology | √ | √ | Aishnivya |
| **Text Classification Algorithm** | Multinomial Naive Bayes, Gaussian Naive Bayes, K Nearest Neighbor, Support Vector Machine | √ | √ | Aishnivya |
| **Multinomial Naive Bayes** | imported the Multinomial NB strategy and the accuracy_score technique | √ | √ | Aishnivya |
| Accuracy Score | Evaluating the classification model | √ | √ | Aishnivya |
| Confusion Matrix | Summary of prediction results | √ | √ | Aishnivya |
| Classification Report | Determine Precison score, recall score and F1 Score of Multinomial Naive Bayes | √ | √ | Aishnivya |
| **K Nearest Neighbor** | Fit() technique of K Nearest Neighbor is implemented | √ | √ | Aishnivya |
| Accuracy Score | Performance of classification model is observed. | √ | √ | Aishnivya |

BIS (Hons) Information Systems Engineering
Faculty of Information and Communication Technology (Kampar Campus), UTAR

| Confusion Matrix | Summary of count values and break down to each class | √ | √ | Aishnivya |
|---|---|---|---|---|
| Classification Report | Determine Precison score, recall score and F1 Score of K Nearest Neighbor | √ | √ | Aishnivya |
| **Support Vector Machine** | Sklearn technique of Support Vector Machine is handled. | √ | √ | Aishnivya |
| Accuracy Score | Depicts the evaluation of classification model | √ | √ | Aishnivya |
| Confusion Matrix | Number of correct and incorrect estimate is viewed | √ | √ | Aishnivya |
| Classification Report | Determine Precison score, recall score and F1 Score of Support Vector Machine | √ | √ | Aishnivya |
| Execution Time | Total execution time that the model took to complete. | √ | √ | Aishnivya |

Table 5.3 Steps Verification

Steps that has been conducted and verified with both the data sets.

## CHAPTER 6 RESULTS AND DISCUSSION

### RESULTS

The results or performance of the classifiers is evaluated using accuracy score, confusion matrix and classification report that includes precision score, recall score and f1-Score

### Confusion Matrix

**Confusion Matrix** table frequently used to depict the exhibition of an arrangement model (or "classifier") on a lot of test information for which the genuine qualities are known in implementation of the classifiers. The disarray grid itself is moderately easy to see, however the related phrasing can be befuddling.

### Classification Report

### Precision

The capacity of a classifier not to name an occasion positive that is really negative. For each class it is characterized as the proportion of genuine positives to the aggregate of valid and bogus positives.

**TP – True Positives**

**FP – False Positives**

Precision – Accuracy of positive predictions.

Precision = TP/(TP + FP)

### Recall

Recall is the capacity of a classifier to locate Twitter data sets. For each class it is characterized as the proportion of genuine positives to the entirety of genuine positives and bogus negatives.

**FN – False Negatives**

Recall: Fraction of positives that were correctly identified.

Recall = TP/(TP+FN)

BIS (Hons) Information Systems Engineering
Faculty of Information and Communication Technology (Kampar Campus), UTAR

**F1 score**

The F1 score is a weighted symphonious mean of exactness and review with the end goal that the best score is 1.0 and the most noticeably awful is 0.0. As a rule, F1 scores are lower than exactness measures as they install accuracy and review into their calculation. As a general guideline, the weighted normal of F1 ought to be utilized to analyze classifier models, not worldwide precision.

F1 Score = 2*(Recall * Precision) / (Recall + Precision)

BIS (Hons) Information Systems Engineering
Faculty of Information and Communication Technology (Kampar Campus), UTAR

**Output Observed of Algorithm In terms of Accuracy score, Precision score and Recall Score**

**Dataset 1**

Multinomial Naive Bayes

| Classifier [Performance Evaluation] | Ratings |
|---|---|
| Accuracy score | 0.865 |
| Precision score | 0.829 |
| Recall score | 0.928 |
| Execution time (s) | 11.466 |

Table 6.1 Evaluation score of Naive Bayes Classifier

K Nearest Neighbor

| Classifier [Performance Evaluation] | Ratings |
|---|---|
| Accuracy score | 0.736 |
| Precision score | 0.86 |
| Recall score | 0.695 |
| Execution time (s) | 11.9 |

Table 6.2 Evaluation score of KNN

Support Vector Machine

| Classifier [Performance Evaluation] | Ratings |
| --- | --- |
| Accuracy score | 0.9146 |
| Precision score | 0.885 |
| Recall score | 0.576 |
| Execution time (s) | 11.87 |

Table 6.3 Evaluation scores of Support Vector Machine

BIS (Hons) Information Systems Engineering
Faculty of Information and Communication Technology (Kampar Campus), UTAR

**Dataset 1**

**Classification report of Algorithms**

Multinomial Naive Bayes Classifier

|  | Precision | Recall | F1-Score | Support |
|---|---|---|---|---|
| Positive | 0.83 | 0.93 | 0.88 | 210 |
| Negative | 0.91 | 0.80 | 0.85 | 200 |

Table 6.4 Classification report of Naive Bayes

K-Nearest Neighbor Classifier

|  | Precision | Recall | F1-Score | Support |
|---|---|---|---|---|
| Positive | 0.86 | 0.58 | 0.69 | 210 |
| Negative | 0.67 | 0.91 | 0.77 | 200 |

Table 6.5 Classification report of KNN

Support Vector Machine

|  | Precision | Recall | F1-Score | Support |
|---|---|---|---|---|
| Positive | 0.89 | 0.96 | 0.92 | 210 |
| Negative | 0.95 | 0.87 | 0.91 | 200 |

Table 6.6 Classification report of Support Vector Machine

**Performance Comparison of Accuracy using Percentage**

| Classifier | Accuracy [%] |
|---|---|
| Multinomial Naive Bayes | 86.6 |
| K-Nearest Neighbor | 73.7 |
| Support Vector Machine | 91.5 |

Table 6.7 Performance Comparison

BIS (Hons) Information Systems Engineering
Faculty of Information and Communication Technology (Kampar Campus), UTAR

**Confusion Matrix of the Algorithms**

**Naive Bayes Classifier**

```
Accuracy score of Naive Bayes in percentage:  86.58536585365853
[[160  40]
 [ 15 195]]
              precision    recall  f1-score   support

          -1       0.91      0.80      0.85       200
           1       0.83      0.93      0.88       210
```

Figure 6.1 Confusion Matrix of Naive Bayes Classifier

Naive Bayes Classifier confusion matrix depicts that out 200 negative tweets 160 were predicted correctly in true negative while among 210 positive tweets 195 were true positive.

**K Nearest Neighbor Classifier**

```
Accuracy score of K Nearest Neighbor:  73.65853658536585
[[181  19]
 [ 89 121]]
              precision    recall  f1-score   support

          -1       0.67      0.91      0.77       200
           1       0.86      0.58      0.69       210
```

Figure 6.2 Confusion Matrix of KNN

The image shows visualise 181 tweets were predicted correctly in -1 while 121 true positive tweets are shown as output with 121 tweets.

**Support Vector Machine**



Figure 6.3 Confusion Matrix of Support Vector Machine

Support Vector Machine shows 174 tweets out of 200 tweets correctly predicted as negative tweets and 201 tweets out 210 were correctly stated in true positive.

**Dataset 2**

**Output Observed of Algorithm In terms of Accuracy score, Precision score and Recall Score**

Mulinomial Naive Bayes Classifier

| Classifier [Performance Evaluation] | Ratings |
|---|---|
| Accuracy score | 0.7884 |
| Precision score | 0.7749 |
| Recall score | 0.9610 |
| Execution time (s) | 18.06 |

Table 6.8 Performance of Naive Bayes Classifier

K Nearest Neighbor

| Classifier [Performance Evaluation] | Ratings |
| --- | --- |
| Accuracy score | 0.69 |
| Precision score | 0.72 |
| Recall score | 0.87 |
| Execution time (s) | 18.46 |

Table 6.9 Performance of KNN Classifier

Support Vector Machine

| Classifier [Performance Evaluation] | Ratings |
| --- | --- |
| Accuracy score | 0.88 |
| Precision score | 0.72 |
| Recall score | 0.87 |
| Execution time (s) | 18.63 |

Table 6.10 Performance of Support Vector Machine classifier

**Classification report of Algorithms**

Multinomial Naive Bayes Classifier

|          | Precision | Recall | F1-Score | Support |
|----------|-----------|--------|----------|---------|
| Positive | 0.86      | 0.58   | 0.69     | 210     |
| Negative | 0.67      | 0.91   | 0.77     | 200     |

Table 6.11 Classification report of Naive Bayes classifier

K-Nearest Neighbor Classifier

|          | Precision | Recall | F1-Score | Support |
|----------|-----------|--------|----------|---------|
| Positive | 0.73      | 0.87   | 0.79     | 437     |
| Negative | 0.58      | 0.35   | 0.44     | 220     |

Table 6.12 Classification report of KNN classifier

Support Vector Machine

|          | Precision | Recall | F1-Score | Support |
|----------|-----------|--------|----------|---------|
| Positive | 0.89      | 0.95   | 0.92     | 437     |
| Negative | 0.88      | 0.77   | 0.82     | 220     |

Table 6.13 Classification report of Support Vector Machine classifier

**Performance Comparison of Accuracy using Percentage**

| Classifier | Accuracy [%] |
|---|---|
| Multinomial Naive Bayes | 78.8 |
| K-Nearest Neighbor | 69.0 |
| Support Vector Machine | 88.0 |

Table 6.14 Performance comparison data set 2

**Confusion Matrix of the algorithm**

Naive Bayes Classifier



Figure 6.4 Confusion Matrix of Naive Bayes Classifier

In Dataset 2, Naive Bayes Clasifier show 98 negative tweets that were predicted and 420 true positive tweets.

K Nearest Neighbor Classifier



Figure 6.5 Confusion Matrix of KNN Classifier

The diagram above shows 77 true negative tweets and 381 true positive tweets.

BIS (Hons) Information Systems Engineering
Faculty of Information and Communication Technology (Kampar Campus), UTAR

Support Vector Machine



```
Accuracy score of Support Vector Machine in percentage:  88.88888888888889
[[170  50]
 [ 23 414]]
              precision    recall  f1-score   support

          -1       0.88      0.77      0.82       220
           1       0.89      0.95      0.92       437
```

Figure 6.6 Confusion Matrix of Support Vector Machine Classifier

170 tweets were predicted correctly as negative while 414 tweets were true positive that was labeled '1'.

**DISCUSSION**



Figure 6.7 Data set 1 accuracy comparison

The bar graph above depicts the accuracy in percentage representation of Multinomial Naive Bayes Classifier that gives the accuracy of 86.5%, K Nearest Neighbor Classifier stated 73.6% and 91.5% of accuracy output of Support Vector Machine has been well visualised. In short, it is understood that Naive Bayes Classifier in Scikit Learn outperforms K Nearest Neighbor and Support Vector Machine performs the best among these mentioned classifiers.

BIS (Hons) Information Systems Engineering
Faculty of Information and Communication Technology (Kampar Campus), UTAR

Figure 6.8 Data set 2 accuracy comparison

The classifiers worked well within the second data set. Multinomial Naive Bayes shows 78.9% and K Nearest Neighbor with 69% and Support Vector Machine 88%. Support Vector Machine have the highest accuracy of evaluation compared to Naive Bayes Classifier and K Nearest Neighbor Classifier. K Nearest Neighbor Classifier stated lesser accuracy among these classifiers , while Naive Bayes Classifier outperformed K Nearest Neighbor classifier.

BIS (Hons) Information Systems Engineering
Faculty of Information and Communication Technology (Kampar Campus), UTAR

CHAPTER 6 RESULTS AND DISCUSSION

The results well depicts that the hypothesis that was predicted Naive Bayes Classifier would outperform K Nearest Neighbor is accepted. The results of accuracy in dataset 1 in percentage of Sklearn Multinomial Naive Bayes Classifier has 86.5% while sklearn K Nearest Neighbor Classifer visualizes 73.6%. Nevertheless, Support Vector Machine Algorithm has the best output in terms of accuracy stating 91.5% in dataset 1(appletwitter.csv).

Dataset 2 that has been used in the similar environment also depicts Scikit Learn Naive Bayes Classifier outperforms Scikit Learn K Nearest Neighbor Classifier. The results shows 78.8% and 69% respectively. Support Vector Machine classifier in Scikit Learn that has been claimed to be a good classification algorithm , once again proved the accuracy with 88%.

Naive Bayes is also said to be much faster compared to K Nearest Neighbor algorithm due to real time execution. KNN also required proper scaling when it comes to classifiication as it measures the euclidean distance as explained in Chapter 3(System Design). As Naive Bayes said to show a good performance in large dataset as well, in small scale dataset as well the good result is depicted as dataset 1 contains 2050 cleaned tweets the accuracy shows higher compared to dataset 2 which contains more than 7000 tweets.

As observed in terms of execution time though Support vector Machine has the highest accuracy compared to K Nearest Neighbor, the execution time Support Vector Machine algorithm is slightly more compared to KNN. Owing to this situation , the reason is learned that Knn is also called a lazy learner, in real execution it does not require training period as the euclidean distance is calculated, the time is executed on testing the dataset.

Support Vector Machine is visualized to show the best accuracy in these datasets, with accuracy score of 91.5%, precision score of 0.885 and recall score of 0.95 in dataset 1. Dataset 2 shows accuracy score of 88%, Precision score of 0.89 and recall score 0.94. These three evaluation metrics for Support Vector Machine is the highest

compared to other classification algorithm mentioned. Nonetheless, the execution time is higher as it is well observed in larger datasets.

Naive Bayes would be best to be used in text classification, spam filtering and sentiment analysis as it also has better accuracy and a stable execution time. Naive Bayes classifer would also be good in weather predictions. Support vector Machine as it has an advantage to provide the highest accuracy it could be implemented well in text and hypertext classification , image recognition anf handwriting recognition, it yields the better accuracy for image classification while text and hypertext classification uses different categories, this algorithm categories document into different categories and uses scores generated to be compared with threshold values. Naive Bayes classifier outperform K Nearest Neighbor algorithm in terms of accuracy and execution time. Support vector Machine is the algorithm with highest accuracy thought it takes a slightly longer time to be executed

**CHAPTER 7 CONCLUSION**

The research study sets objective to evaluate the performance for sentiment classification that takes into account on accuracy and precision. The main problem as mentioned was information overload as known, pages and users has very own perception and thoughts on their daily life. Thus is makes difficulty in identifying content based on interest and in most situation recommending unwanted items. Therefore solving such situations is a significant task by condensing data overload and identify the real need of users to be recommended. It is also important to recommend the right items at right time. The hypothesis predicted was accepted that mentioned Naive Bayes Classifier outperform K Nearest Neighbor Algorithm. Comparing Naive Bayes Classifier and K Nearest Neighbor it is proven that Naive Bayes Classifier works the best with large data sets.

The research study conducted compare Naive Bayes Algorithm , K Nearest Neighbor Algorithm and Support Vector Machine Algorithm to evaluate the performance in terms of accuracy. Personalised Recommender System implements algorithms in recommending context or a subject of classification. Scikit Learn library is used to module, train and test the data sets. Sentiment Analyser has been handled to get the positive, negative and neutral tweets, Preprocessing steps were implemented to clean the tweets which has punctuation and stop words.

The results obtained shows Support Vector Machine performs the best among the classification algorithm Naive Bayes and K Nearest Neighbor. However it has higher execution time, it is well observed in large data sets. Naive Bayes Classifier outperform K Nearest Neighbor in term of accuracy and execution time.

CHAPTER 7 CONCLUSION

**Objectives Achieved**

There were objectives to be achieved in the project and the objectives are achieved in the implementation. The objective is to condense data overload by retrieving the most suited data services from large amount of data. It has been achieved in the context the tweets that has been imported has been preprocessed to eliminate unwanted characters which are irrelevant to the message. Information overload is the issue , the data content reviews of many kind, neutral reviews in this case had been dropped as it does not help in predicting the polarity of the message being conveyed. The objective has been achieved. Besides in general context it is also very significant in the usage of market research despite on what the company thinks of the product they will be able to know the perspective of their customers. Much interestingly the company would also be able to predict on the review of other companies.

The objective is also to analyse users' conveyed information, using text classification algorithm that predicts higher accuracy and performance. Text classification has become increasingly important. Naive Bayes Classifier, K Nearest Neighbor Classifier and Support Vector Machine Classifier has been evaluated in the project. Naive Bayes algorithm is considered to have higher accuracy compared to K Nearest Neighbor algorithm and Support Vector Machine classifier depicts the highest accuracy among the text classifications methodologies. Therefore the objective is achieved.

The objective stated is to identify classifier that recommend suitable or most similar context at a specific situation. The objective is achieved in the implementation of the project. Support Vector machine classifier is most suited text classification compared to other algorithms that has been analyzed with Twitter data sets. It yields the highest accuracy among both the Twitter data sets that has been implemented. Naive Bayes classifier however is still suitable , it has outperformed K Nearest Neighbor. It is well depicted that Naive Bayes and Support Vector Machine has higher accuracy it large datasets while K Nearest Neighbor performs slightly lesser accuracy in large data sets. Sentiment Analysis been used to gauge public opinion to policy announcement and messaged ahead of even presidential elections. Brands and companies has also been

BIS (Hons) Information Systems Engineering
Faculty of Information and Communication Technology (Kampar Campus), UTAR

using sentiment analysis as recommender system to analyses user reviews. Television advertisements and songs widely uses these sentiment analysis methodologies as they are able to justify the reviews of the customers at right time and fast to be able to overcome the situation.

**Highlight of any novelties and contribution**

The project has analysed machine learning algorithm with lexicons based approach. The text classification methods that has been implemented is Naive Bayes Classifier , KNN Classifier and Support Vector Machine Classifier. Shown in chapter 6 results and discussion Support Vector Machine classifier has the highest accuracy. Therefore it is best suited in image classification applications and handwriting recognitions. It yields higher performance in terms of classification compared to Naive Bayes classifier and K Nearest Neighbor Classifier. Naive Bayes Classifier although it shows a moderate among these mentioned classifiers, it still works the best with text classification of Twitter data set. Naive Bayes classifier would also be able work well in weather prediction as it allows continuous data distribution. K Nearest Neighbor as it used Euclidean Distance in the algorithm it is simpler yet testing can be added seamlessly into data set.

**Future Work**

Thrawated expressions should be analysed. For instead the tweet 'the movie should have been interesting, the plot was nice somehow the actors should have acted much better'. Though this sentence states all positive words, the message that has been conveyed is negative. Therefor the future work should consider such analyses. Usage of symbols and punctuation is very important is meaning conveying as well. The position of punctuation should also be helpful in determining message polarity. The question mark, full stops and symbol also conveys meaning to the tweets. Therefore it should be also considered in the future work. Entity recognition. I like ice cream but I dislike the flavour. There is a need to separate out the tweet or texts about an entity

before sentiment analysis being conducted. This sentence has both positive and negative word, thus it might be considered neutral. However the sentence carries specific polarity for each entity or the object.

A-1

**APPENDIX A Final Year Project Biweekly Report**

**FINAL YEAR PROJECT WEEKLY REPORT**

(Project I / <u>Project II</u>)

| Trimester, Year: MAY,2020 | Study week no.: 1 - 2 |
|---|---|
| Student Name & ID: AISHNIVYA A/P BALAMURUGAN     1706183 | |
| Supervisor: DR. RAMESH KUMAR AYYASAMY | |
| Project Title: A study on Personalised Recommender System using social media. | |

**1. WORK DONE**

[Please write the details of the work done in the last fortnight.]

- The planning of the process, the estimated time to complete FYP 2

- Objectives verifies

- Implementation of system design

**2. WORK TO BE DONE**

- Sentiment analyse and implementation of text classification.

**3. PROBLEMS ENCOUNTERED -**    No Problems encountered

**4. SELF EVALUATION OF THE PROGRESS**

**-** need to proceed further

_____

Supervisor's signature

_____

Student's signature

A-1

BIS (Hons) Information Systems Engineering

Faculty of Information and Communication Technology (Kampar Campus), UTAR

# FINAL YEAR PROJECT WEEKLY REPORT

*(Project I / Project II)*

| | |
|---|---|
| **Trimester, Year: MAY, 2020** | **Study week no.: 3 - 4** |
| **Student Name & ID: AISHNIVYA A/P BALAMURUGAN     1706183** | |
| **Supervisor: DR. RAMESH KUMAR AYYASAMY** | |
| **Project Title: A study on Personalised Recommender System using social media.** | |

**1. WORK DONE**

[Please write the details of the work done in the last fortnight.]

- Finalize Project Scope

- Finalized problem statement and objective

- Conducted Sentiment Analysis

**2. WORK TO BE DONE**

- Implementation of text classification

**3. PROBLEMS ENCOUNTERED**

- No problem encountere

**4. SELF EVALUATION OF THE PROGRESS**

- Okay

_____                    _____

Supervisor's signature                          Student's signature

BIS (Hons) Information Systems Engineering

Faculty of Information and Communication Technology (Kampar Campus), UTAR

# FINAL YEAR PROJECT WEEKLY REPORT

*(Project I / Project II)*

| | |
|---|---|
| **Trimester, Year: MAY,2020** | **Study week no.: 5 - 6** |
| **Student Name & ID: AISHNIVYA A/P BALAMURUGAN     1706183** | |
| **Supervisor: DR. RAMESH KUMAR AYYASAMY** | |
| **Project Title: A study on Personalised Recommender System using social media.** | |

---

**1. WORK DONE**

[Please write the details of the work done in the last fortnight.]

- Text classification algorithm completed

---

**2. WORK TO BE DONE**

- Evaluation metrics on text classification

---

**3. PROBLEMS ENCOUNTERED**

- No problem encountered

---

**4. SELF EVALUATION OF THE PROGRESS**

**-** Okay

---

_____                    _____

Supervisor's signature                                                  Student's signature

BIS (Hons) Information Systems Engineering

Faculty of Information and Communication Technology (Kampar Campus), UTAR

# FINAL YEAR PROJECT WEEKLY REPORT

*(Project I / Project II)*

| | |
|---|---|
| **Trimester, Year: MAY,2020** | **Study week no.: 7 - 8** |
| **Student Name & ID: AISHNIVYA A/P BALAMURUGAN      1706183** | |
| **Supervisor: DR. RAMESH KUMAR AYYASAMY** | |
| **Project Title: A study on Personalised Recommender System using social media.** | |

**1. WORK DONE**

[Please write the details of the work done in the last fortnight.]

- Completed evaluation metrics

**2. WORK TO BE DONE**

**-**System testing

- Update the results and discussion in documentation

**3. PROBLEMS ENCOUNTERED**

- No Problems Encountered

**4. SELF EVALUATION OF THE PROGRESS**

**-**Okay

Supervisor's signature                                    Student's signature

BIS (Hons) Information Systems Engineering

Faculty of Information and Communication Technology (Kampar Campus), UTAR

# FINAL YEAR PROJECT WEEKLY REPORT

*(Project I / Project II)*

| | |
|---|---|
| **Trimester, Year: MAY,2020** | **Study week no.: 9 - 10** |
| **Student Name & ID: AISHNIVYA A/P BALAMURUGAN      1706183** | |
| **Supervisor: DR. RAMESH KUMAR AYYASAMY** | |
| **Project Title: A study on Personalised Recommender System using social media.** | |

**1. WORK DONE**

[Please write the details of the work done in the last fortnight.]

- Results and discussion been updated

- System testing completed

- Project has been concluded

**2. WORK TO BE DONE**

- Add appendix

**3. PROBLEMS ENCOUNTERED**

- No problems encountered

**4. SELF EVALUATION OF THE PROGRESS**

**-** Moderate

_____                     _____

Supervisor's signature                                     Student's signature

A-5

# FINAL YEAR PROJECT WEEKLY REPORT

*(Project I / Project II)*

| | |
|---|---|
| **Trimester, Year: MAY,2020** | **Study week no.: 11 - 12** |
| **Student Name & ID: AISHNIVYA A/P BALAMURUGAN    1706183** | |
| **Supervisor: DR. RAMESH KUMAR AYYASAMY** | |
| **Project Title: A study on Personalised Recommender System using social media.** | |

**1. WORK DONE**

[Please write the details of the work done in the last fortnight.]

-update information

-review report

- Completed report

**2. WORK TO BE DONE**

**-** Submit for turnitin check

- Submission to supervisor

- submission of report

- prepare for the presentation

**3. PROBLEMS ENCOUNTERED**

**-** Project has been completed

**4. SELF EVALUATION OF THE PROGRESS**

- Okay

_____                    _____

Supervisor's signature                                              Student's signature

A-6

BIS (Hons) Information Systems Engineering

Faculty of Information and Communication Technology (Kampar Campus), UTAR

# APPENDIX B POSTER

## APPENDIX C CODES IN PYTHON

```python
16
17    import numpy as np
18    import re
19    import pandas as pd
20    import pylab as pl
21
22
23    from nltk.tokenize import WordPunctTokenizer
24    from bs4 import BeautifulSoup
25
26    from sklearn.feature_extraction.text import TfidfVectorizer
27    from sklearn.decomposition import PCA, TruncatedSVD
28
29    from vaderSentiment.vaderSentiment import SentimentIntensityAnalyzer
30    analyser = SentimentIntensityAnalyzer()
31
32    from sklearn import metrics
33    from mpl_toolkits.mplot3d import Axes3D
34
35    from matplotlib import pyplot
36    import preprocessor as p
37    import time
38
39    import nltk
40
41
42    start = time.time()
43    df = pd.read_csv('appletwitter.csv')
44    print(len(df))
45
46    #df["_unit_id"].value_counts(ascending=False).head(2)
47
48
49    def remove_pattern(input_txt, pattern):
50        r = re.findall(pattern, input_txt)
51        for i in r:
52            input_txt = re.sub(i, '', input_txt)
53
```

Appendix C Codes in Python

```python
49      def remove_pattern(input_txt, pattern):
50          r = re.findall(pattern, input_txt)
51          for i in r:
52              input_txt = re.sub(i, '', input_txt)
53
54          return input_txt
55
56
57      df['Clean_text'] = np.vectorize(remove_pattern)(df['text'], "@[\w]*")
58
59
60      df['Clean_text'] = df ['text'].str.lower()
61      example_review = df.iloc[0]
62      print(example_review['Clean_text'])
63
64
65      print(nltk.word_tokenize(example_review['Clean_text']))
66
67
68      def identify_token(row):
69          Clean_text=row['Clean_text']
70          tokens = nltk.word_tokenize(Clean_text)
71
72          #taken only words and not punctuation
73          token_words =[w for w in tokens if w.isalpha()]
74          return token_words
75      df['words'] = df.apply(identify_token, axis=1)
76
77
78
79      from nltk.stem import PorterStemmer
80      stemming = PorterStemmer()
81      def stem_list(row):
82          my_list = row['words']
83          stemmed_list = [stemming.stem(word) for word in my_list]
84          return(stemmed_list)
85      df['stemmed_words'] = df.apply(stem_list,axis=1)
```

```python
88
89      from nltk.corpus import stopwords
90      stops = set(stopwords.words("english"))
91
92    ▼ def remove_stops(row):
93          my_list = row['stemmed_words']
94          meaningful_words = [w for w in my_list if not w in stops]
95          return (meaningful_words)
96      df['stem_meaningful'] = df.apply(remove_stops, axis=1)
97
98
99    ▼ def rejoin_words(row):
100         my_list = row['stem_meaningful']
101         joined_words = ( " ".join(my_list))
102         return joined_words
103     df['Tweets_clean'] = df.apply(rejoin_words, axis=1)
104
105     missing_values_count=df.isnull().sum()
106     missing_values_count[0:13]
107
108
109
110     #To determine length of clean text
111     df['Clean_text_length'] = df['Tweets_clean'].apply(len)
112
113
114
115     from textblob import TextBlob
116
117   ▼ def calculate_sentiment(Tweets_clean):
118         return TextBlob(Tweets_clean).sentiment
119
120   ▼ def calculate_sentiment_analyser(Tweets_clean):
121         return analyser.polarity_scores(Tweets_clean)
122
123     df['sentiment']=df.Tweets_clean.apply(calculate_sentiment)
124     df['sentiment_analyser']=df.Tweets_clean.apply(calculate_sentiment_analyser)
```

```python
125
126
127     s = pd.DataFrame(index = range(0,len(df)),columns= ['compound_score','compound_score_sentiment'])
128
129   ▼ for i in range(0,len(df)):
130       s['compound_score'][i] = df['sentiment_analyser'][i]['compound']
131
132   ▼   if (df['sentiment_analyser'][i]['compound'] <= -0.05):
133         s['compound_score_sentiment'][i] = 'Negative'
134   ▼   if (df['sentiment_analyser'][i]['compound'] >= 0.05):
135         s['compound_score_sentiment'][i] = 'Positive'
136   ▼   if ((df['sentiment_analyser'][i]['compound'] >= -0.05) & (df['sentiment_analyser'][i]['compound'] <= 0.05)):
137         s['compound_score_sentiment'][i] = 'Neutral'
138
139     df['compound_score'] = s['compound_score']
140     df['compound_score_sentiment'] = s['compound_score_sentiment']
141     df.head()
142
143
144     count=df['compound_score_sentiment'].value_counts()
145     print(count)
146
147     #drop neutral tweets
148     df = df.set_index("compound_score_sentiment")
149     df = df.drop("Neutral", axis=0)
150
151
152     df['label'] = df['compound_score'].apply(lambda x: -1 if x < 0 else 1)
153
154
155     #View total number of rows in df
156     rows=len(df.axes[0])
157     print('Total number of positive and negative', rows)
158
159     #Count numbers of positive and negative tweets
160     count=df['label'].value_counts()
161     print(count)
```
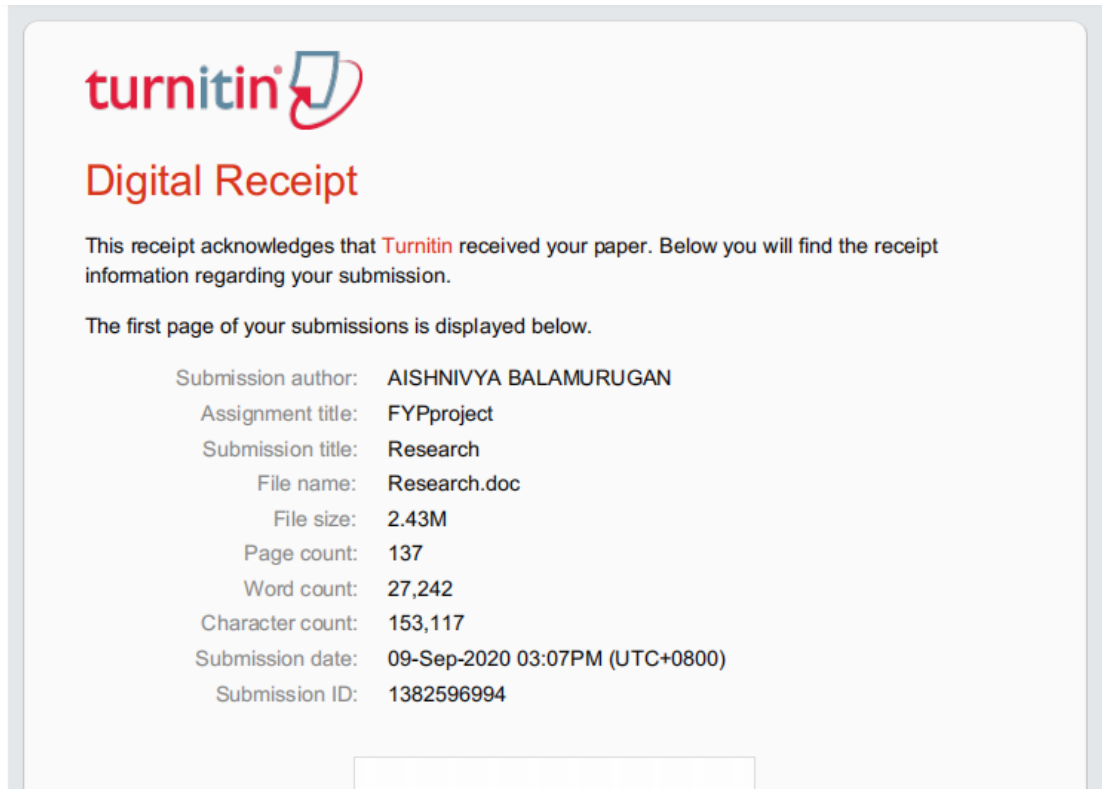
```python
159    #Count numbers of positive and negative tweets
160    count=df['label'].value_counts()
161    print(count)
162
163
164
165
166    #spliiting dataset into train set and test set
167    from sklearn.model_selection import train_test_split
168    X_train, X_test, y_train, y_test = train_test_split(df['Tweets_clean'],
169             df['label'], test_size=0.2, random_state=0)
170
171
172    #Multinomial Naive Bayes
173
174    from sklearn.feature_extraction.text import CountVectorizer
175    from sklearn import metrics
176    from sklearn.naive_bayes import MultinomialNB
177
178
179
180    cv = CountVectorizer(strip_accents='ascii', token_pattern = u'(?ui)\\b\\w*[a-z]+\\w*\\b',
181                         lowercase=True, stop_words='english')
182    X_train_cv = cv.fit_transform(X_train)
183    X_test_cv = cv.transform(X_test)
184
185
186    word_freq_df = pd.DataFrame (X_train_cv.toarray(), columns = cv.get_feature_names())
187    top_words_df = pd.DataFrame(word_freq_df.sum()).sort_values(0,ascending=False)
188
189    naive_bayes = MultinomialNB()
190    naive_bayes.fit(X_train_cv, y_train)
191    predictions = naive_bayes.predict(X_test_cv)
192
193      from sklearn.metrics import accuracy_score, precision_score, recall_score
194
195      print('Accuracy score of Naive Bayes: ', accuracy_score(y_test, predictions))
196      print('Precision score Naive Bayes: ', precision_score(y_test, predictions))
197      print('Recall score Naive Bayes: ',recall_score(y_test, predictions))
198
199      print('Accuracy score of Naive Bayes in percentage: ',
200              accuracy_score(y_test, predictions)*100)
201
202      from sklearn.metrics import confusion_matrix
203      print (confusion_matrix(y_test, predictions))
204
205      from sklearn.metrics import classification_report
206      print(classification_report(y_test,predictions))
207
208
209
210      time.sleep(1)
211
212      end = time.time()
213
214      print(f"Runtime is {end - start}")
215
216
```
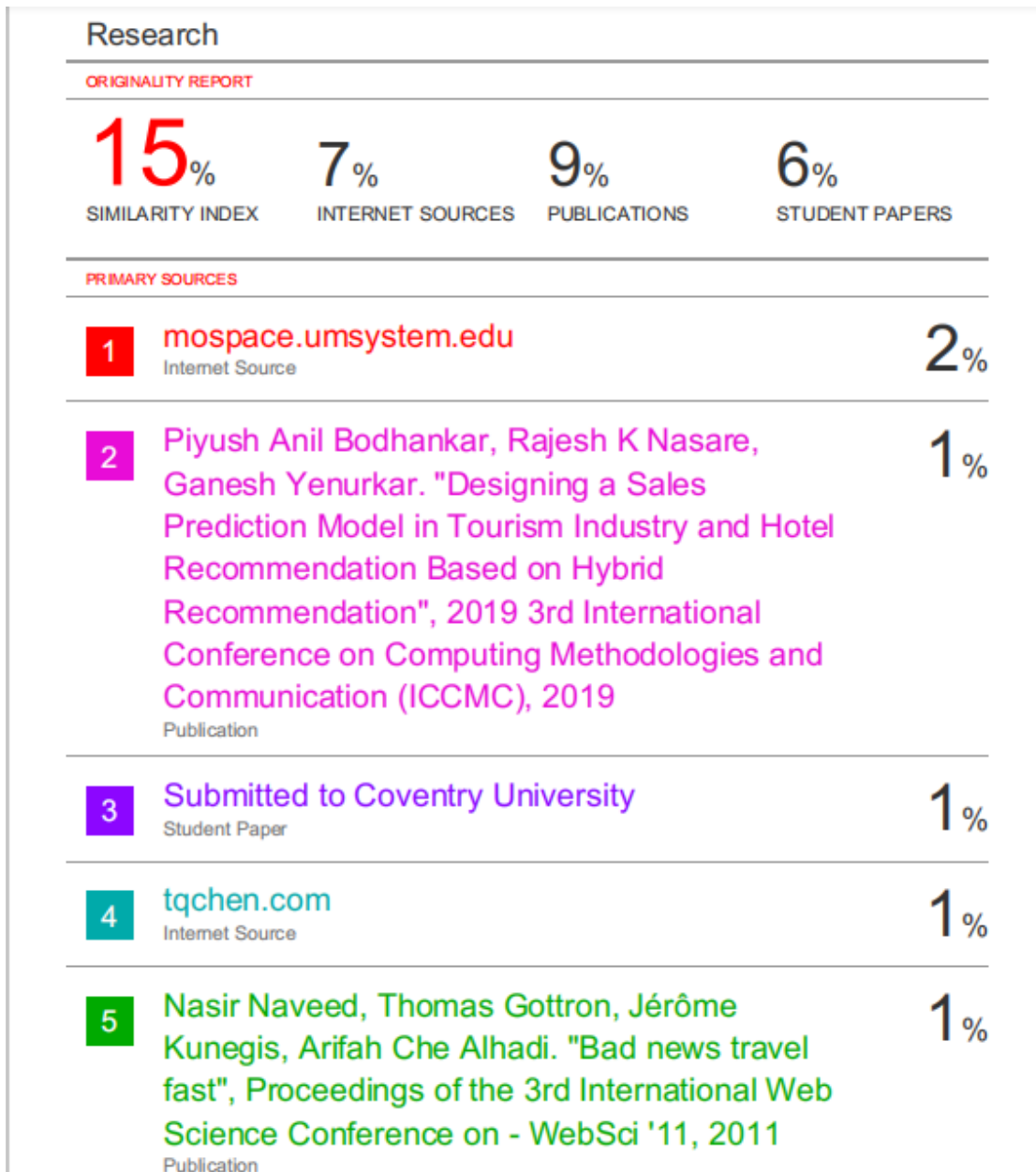
**APPENDIX D TURNITIN CHECK RESULT**



turnitin

**Digital Receipt**

This receipt acknowledges that Turnitin received your paper. Below you will find the receipt information regarding your submission.

The first page of your submissions is displayed below.

| | |
|---|---|
| Submission author: | AISHNIVYA BALAMURUGAN |
| Assignment title: | FYPproject |
| Submission title: | Research |
| File name: | Research.doc |
| File size: | 2.43M |
| Page count: | 137 |
| Word count: | 27,242 |
| Character count: | 153,117 |
| Submission date: | 09-Sep-2020 03:07PM (UTC+0800) |
| Submission ID: | 1382596994 |

## Research

ORIGINALITY REPORT

| 15% | 7% | 9% | 6% |
|---|---|---|---|
| SIMILARITY INDEX | INTERNET SOURCES | PUBLICATIONS | STUDENT PAPERS |

PRIMARY SOURCES

| 1 | mospace.umsystem.edu<br>Internet Source | 2% |
|---|---|---|
| 2 | Piyush Anil Bodhankar, Rajesh K Nasare, Ganesh Yenurkar. "Designing a Sales Prediction Model in Tourism Industry and Hotel Recommendation Based on Hybrid Recommendation", 2019 3rd International Conference on Computing Methodologies and Communication (ICCMC), 2019<br>Publication | 1% |
| 3 | Submitted to Coventry University<br>Student Paper | 1% |
| 4 | tqchen.com<br>Internet Source | 1% |
| 5 | Nasir Naveed, Thomas Gottron, Jérôme Kunegis, Arifah Che Alhadi. "Bad news travel fast", Proceedings of the 3rd International Web Science Conference on - WebSci '11, 2011<br>Publication | 1% |

Balachandran Manavalan, Jooyoung Lee.

| 6 | "SVMQA: support–vector-machine-based protein single-model quality assessment", Bioinformatics, 2017<br>Publication | 1% |
| 7 | Thomas Renault. "Sentiment analysis and machine learning in finance: a comparison of methods and models on one million messages", Digital Finance, 2019<br>Publication | 1% |
| 8 | Submitted to Ain Shams University<br>Student Paper | <1% |
| 9 | Submitted to Victorian Institute of Technology<br>Student Paper | <1% |
| 10 | www.bitdegree.org<br>Internet Source | <1% |
| 11 | Chaithra V. D. "Hybrid approach: naive bayes and sentiment VADER for analyzing sentiment of mobile unboxing video comments", International Journal of Electrical and Computer Engineering (IJECE), 2019<br>Publication | <1% |

Publication

| | | |
|---|---|---|
| 12 | link.springer.com<br>Internet Source | <1% |
| 13 | Rishabh Ahuja, Arun Solanki, Anand Nayyar. "Movie Recommender System Using K-Means Clustering AND K-Nearest Neighbor", 2019 9th | <1% |

International Conference on Cloud Computing, Data Science & Engineering (Confluence), 2019
Publication

| | | |
|---|---|---|
| 14 | Submitted to Higher Education Commission Pakistan<br>Student Paper | <1% |
| 15 | d4datascience.wordpress.com<br>Internet Source | <1% |
| 16 | Submitted to University of Wales Institute, Cardiff<br>Student Paper | <1% |
| 17 | mafiadoc.com<br>Internet Source | <1% |
| 18 | Submitted to CBA<br>Student Paper | <1% |

| | | |
|---|---|---|
| 18 | **Submitted to CBA**<br>Student Paper | <1% |
| 19 | **Submitted to University of Leeds**<br>Student Paper | <1% |
| 20 | **Submitted to Harrisburg University of Science and Technology**<br>Student Paper | <1% |
| 21 | **ojs3.unpatti.ac.id**<br>Internet Source | <1% |
| 22 | **Submitted to University of Northumbria at Newcastle**<br>Student Paper | <1% |
| 23 | **"Data Science and Big Data Analytics", Springer Science and Business Media LLC, 2019**<br>Publication | <1% |
| 24 | **Y Findawati, I R Indra Astutik, A S Fitroni, I Indrawati, N Yuniasih. "Comparative analysis of Naïve Bayes, K Nearest Neighbor and C.45 method in weather forecast", Journal of Physics: Conference Series, 2019**<br>Publication | <1% |

| | | |
|---|---|---|
| 25 | **Submitted to Banaras Hindu University** <br> Student Paper | <1% |
| 26 | **www.dataquest.io** <br> Internet Source | <1% |
| 27 | **Submitted to Savitribai Phule Pune University** <br> Student Paper | <1% |
| 28 | **"ICDSMLA 2019", Springer Science and Business Media LLC, 2020** <br> Publication | <1% |
| 29 | **matheo.uliege.be** <br> Internet Source | <1% |
| 30 | **Submitted to Kookmin University** <br> Student Paper | <1% |
| 31 | **Rungroj Maipradit, Hideaki Hata, Kenichi Matsumoto. "Sentiment Classification Using N-Gram Inverse Document Frequency and** | <1% |

**Automated Machine Learning", IEEE Software, 2019**
Publication

Lecture Notes in Computer Science, 2015

| 32 | Lecture Notes in Computer Science, 2015.<br>Publication | <1% |
| 33 | "Enhancing Classifier Accuracy in Ayurvedic Medicinal Plants using WO-DNN", International Journal of Engineering and Advanced Technology, 2019<br>Publication | <1% |
| 34 | Submitted to University of Ulster<br>Student Paper | <1% |
| 35 | Tarek Mahmoud, Tarek Abd-El-Hafeez, Doha El-Deen. "A Design of an Automatic Web Page Classification System", British Journal of Applied Science & Technology, 2016<br>Publication | <1% |
| 36 | M. S. Neethu, R. Rajasree. "Sentiment analysis in twitter using machine learning techniques", 2013 Fourth International Conference on Computing, Communications and Networking Technologies (ICCCNT), 2013<br>Publication | <1% |
| 37 | csce.unl.edu<br>Internet Source | <1% |
| 38 | Submitted to London School of Commerce - | |

| 38 | Submitted to London School of Commerce - | |
|----|------------------------------------------|---|

| | Dhaka<br>Student Paper | <1% |
|----|---|---|
| 39 | K. Deivanai, V. Vijayakumar, Priyanka. "Chapter 32 Automated Workload Management Using Machine Learning", Springer Science and Business Media LLC, 2019<br>Publication | <1% |
| 40 | Submitted to University of Florida<br>Student Paper | <1% |
| 41 | "Linguistic-Based and user-Based Recommending Posts using Two-Level Clustering Methods", International Journal of Innovative Technology and Exploring Engineering, 2019<br>Publication | <1% |
| 42 | Eslam Omara, Mervat Mosa, Nabil Ismail. "Deep Convolutional Arabic Sentiment Analysis with Imbalanced Data", 2019 15th International Computer Engineering Conference (ICENCO), 2019<br>Publication | <1% |

| | | |
|---|---|---|
| 43 | Submitted to Deakin University<br>Student Paper | <1% |
| 44 | Gianni Barlacchi, Christos Perentis, Abhinav Mehrotra, Mirco Musolesi, Bruno Lepri. "Are you getting sick? Predicting influenza-like symptoms | <1% |

| | | |
|---|---|---|
| | using human mobility behaviors", EPJ Data Science, 2017<br>Publication | |
| 45 | Submitted to Indian Institute of Information Technology, Allahabad<br>Student Paper | <1% |
| 46 | docplayer.net<br>Internet Source | <1% |
| 47 | Submitted to Visvesvaraya Technological University, Belagavi<br>Student Paper | <1% |
| 48 | Submitted to LearnBook<br>Student Paper | <1% |
| 49 | Submitted to University of Central Lancashire<br>Student Paper | <1% |

| 50 | Abdullah Alsaeedi, Mohammad Zubair. "A Study on Sentiment Analysis Techniques of Twitter Data", International Journal of Advanced Computer Science and Applications, 2019 Publication | <1% |
| 51 | muthu.co Internet Source | <1% |
| 52 | Submitted to Sardar Vallabhbhai National Inst. of Tech.Surat Student Paper | <1% |
| 53 | www.mdpi.com Internet Source | <1% |
| 54 | Kumar, Hemantha, T.A. Ranjit Kumar, M. Amarnath, and V. Sugumaran. "Fault diagnosis of bearings through vibration signal using Bayes classifiers", International Journal of Computer Aided Engineering and Technology, 2014. Publication | <1% |
| 55 | "Artificial Intelligence", Springer Science and Business Media LLC, 2019 Publication | <1% |

| | | |
|---|---|---|
| 61 | www.ijbhtnet.com<br>Internet Source | <1% |
| 62 | www.emerald.com<br>Internet Source | <1% |
| 63 | www.sersc.org<br>Internet Source | <1% |
| 64 | Lecture Notes in Computer Science, 2014.<br>Publication | <1% |
| 65 | export.arxiv.org<br>Internet Source | <1% |
| 66 | pdfs.semanticscholar.org<br>Internet Source | <1% |
| 67 | docs.oracle.com<br>Internet Source | <1% |
| 68 | openjournals.wu.ac.at<br>Internet Source | <1% |
| 69 | worldwidescience.org<br>Internet Source | <1% |
| 70 | research.aalto.fi | |

Internet Source

**APPENDIX E TURNITIN FORM**

| Universiti Tunku Abdul Rahman | | | |
|---|---|---|---|
| **Form Title : Supervisor's Comments on Originality Report Generated by Turnitin for Submission of Final Year Project Report (for Undergraduate Programmes)** | | | |
| Form Number: FM-IAD-005 | Rev No.: 0 | Effective   Date: 01/10/2013 | Page No.: 1of 1 |

**FACULTY OF INFORMATION AND COMMUNICATION TECHNOLOGY**

| Full Name(s) of Candidate(s) | AISHNIVYA A/P BALAMURUGAN |
|---|---|
| ID Number(s) | 17ACB06183 |
| Programme / Course | IA |
| Title of Final Year Project | A Study on Personalised Recommender System using Social Media |

| Similarity | Supervisor's Comments (Compulsory if parameters of originality exceeds   the limits approved by UTAR) |
|---|---|
| **Overall similarity index:_ 15%_____ %** **Similarity by source** Internet Sources:7_____% Publications:          9      % Student Papers:            6     % | |
| **Number of individual sources listed** of more than 3% similarity:  0_____ | |

**Parameters of originality required and limits approved by UTAR are as Follows:**
   **(i)       Overall similarity index is 20% and below, and**
   **(ii)      Matching of individual sources listed must be less than 3% each, and**
   **(iii) Matching texts in continuous block must not exceed 8 words**
*Note: Parameters (i) – (ii) shall exclude quotes, bibliography and text matches which are*

Note   Supervisor/Candidate(s) is/are required to provide softcopy of full set of the originality report to Faculty/Institute

*Based on the above results, I hereby declare that I am satisfied with the originality of the Final Year Project Report submitted by my student(s) as named above.*

_____          _____

  Signature of Supervisor                              Signature of Co-Supervisor

  Name:      DR.Ramesh        Kumar            Name: _____
Ayyasamy_____

  Date:    09/09/2020__                              Date: _____

**APPENDIX F Final Year Project Check List**



**UNIVERSITI TUNKU ABDUL RAHMAN**
FACULTY OF INFORMATION & COMMUNICATION TECHNOLOGY
(KAMPAR CAMPUS)
**CHECKLIST FOR FYP2 THESIS SUBMISSION**

| Student Id | 17ACB06183 |
|---|---|
| Student Name | AISHNIVYA A/P BALAMURUGAN |
| Supervisor Name | DR.RAMESH KUMAR AYYASAMY |

| TICK (√) | DOCUMENT ITEMS |
|---|---|
| | Your report must include all the items below. Put a tick on the left column after you have checked your report with respect to the corresponding item. |
| √ | Front Cover |
| √ | Signed Report Status Declaration Form |
| √ | Title Page |
| √ | Signed form of the Declaration of Originality |
| √ | Acknowledgement |
| √ | Abstract |
| √ | Table of Contents |
| √ | List of Figures (if applicable) |
| √ | List of Tables (if applicable) |
| √ | List of Symbols (if applicable) |
| √ | List of Abbreviations (if applicable) |
| √ | Chapters / Content |
| √ | Bibliography (or References) |
| √ | All references in bibliography are cited in the thesis, especially in the chapter of literature review |
| √ | Appendices (if applicable) |
| √ | Poster |
| √ | Signed Turnitin Report (Plagiarism Check Result - Form Number: FM-IAD-005) |

*Include this form (checklist) in the thesis (Bind together as the last page)

| I, the author, have checked and confirmed all the items listed in the table are included in my report. | Supervisor verification. Report with incorrect format can get 5 mark (1 grade) reduction. |
|---|---|
| ___ _____ (Signature of Student) Date: 09/09/2020 | _____ _____ (Signature of Supervisor) Date: 09/09/2020 |