

# **PERSONALITY PREDICTION ANALYSIS**

## **A MINI PROJECT REPORT**

*Submitted by*

*Mitasree satpati(RA2011003010005 )  
SHREEJA G(RA2011003010007)  
SANJAY G(RA2011003010022)*

*Under the guidance of*

*Ghuntupalli Manoj kumar N*

(Associate Professor, CTECH)

*In partial satisfaction of the requirements for the degree of*

**BACHELOR OF TECHNOLOGY**

in

**COMPUTER SCIENCE & ENGINEERING**



**SCHOOL OF COMPUTING**

**COLLEGE OF ENGINEERING AND TECHNOLOGY**

**SRM INSTITUTE OF SCIENCE AND TECHNOLOGY**

**KATTANKULATHUR - 603203**

**MAY 2023**



COLLEGE OF ENGINEERING & TECHNOLOGY  
SRM INSTITUTE OF SCIENCE & TECHNOLOGY  
S.R.M. NAGAR, KATTANKULATHUR – 603 203

## BONAFIDE CERTIFICATE

Certified that this project report “ **Personality Prediction System** ” is the bonafide work

of **Mitasree satpati(RA2011003010005 )**, **SHREEJA G(RA2011003010007)**  
**SANJAY G(RA2011003010022)** of III Year/VI Sem B. Tech(CSE) who carried out the mini

project work under my supervision for the course 18CSC305J- Artificial Intelligence in SRM Institute of Science and Technology during the academic year 2022-2023(Even Sem).

**Signature of the Faculty**

***Ghuntupalli Manoj kumar N***

Associate Professor

Department of CTECH

S.R.M Institute of Science And Technology

**Signature of HOD**

**Dr. M.Pushpalatha**

Professor and Head

S.R.M Institute of Science And Technology

# ABSTRACT

This personality detection project aims to develop a machine learning model that can accurately predict personality traits based on textual data, such as social media posts or emails. The project uses natural language processing techniques to preprocess and analyze the textual data and extract features that are indicative of various personality traits. The extracted features are then used to train a machine learning model, such as a support vector machine or a neural network, to predict personality traits, such as extraversion, agreeableness, and neuroticism. The model's performance is evaluated using various metrics, such as accuracy, precision, recall, and F1 score, and compared to the state-of-the-art methods in the field of personality detection. The project's ultimate goal is to develop a robust and reliable model that can be used to automate personality detection and assist in various applications, such as job recruiting, mental health screening, and personalized marketing.

**Keywords:** personality detection, machine learning, natural language processing, textual data, social media, personality traits, extraversion, agreeableness, neuroticism, support vector machine, neural network, performance evaluation, accuracy, precision, recall, F1 score, job recruiting, mental health screening, personalized marketing.

# TABLE OF CONTENTS

Chapter No.	Title	Page No.
	<b>ABSTRACT</b>	<b>iii</b>
	<b>TABLE OF CONTENTS</b>	<b>iv</b>
	<b>LIST OF FIGURES</b>	<b>v</b>
	<b>LIST OF TABLES</b>	<b>vi</b>
	<b>ABBREVIATIONS</b>	<b>vii</b>
<b>1</b>	<b>INTRODUCTION</b>	
1.1	Introduction	1
1.2	Problem statement	2
1.3	Objectives	3
1.4	Scope and applications	3
1.5	General and Unique Services in the database application	4
1.6	Software Requirements Specification	4
<b>2</b>	<b>LITERATURE SURVEY</b>	
2.1	Existing system	5
2.2	Comparison of Existing vs Proposed system	6
<b>3</b>	<b>SYSTEM ARCHITECTURE AND DESIGN</b>	
3.1	Architecture Diagram	9
3.1.1	Front end (UI)design	9
3.1.2	Back end (Database) design	9
3.2	ER Diagram and Use case Diagram	10
<b>4</b>	<b>MODULES AND FUNCTIONALITIES</b>	
4.1	Machine Learning	10
4.2	Algorithms	11
<b>5</b>	<b>CODING AND TESTING</b>	<b>18</b>
<b>6</b>	<b>RESULTS AND DISCUSSIONS</b>	
6.1	Assumption and dependencies	33
6.2	Conclusion	
6.3	Output	33
<b>7</b>	<b>CONCLUSION AND FUTURE ENHANCEMENT</b>	<b>34</b>
	<b>REFERENCES</b>	<b>35</b>

## **LIST OF FIGURES**

<b>Figure No.</b>	<b>Figure Name</b>	<b>Page No.</b>
3.1.1	Front end (UI) design	21
3.1.2	Backend Database Design	21
3.2	ER Diagram and use case diagram	22

## LIST OF TABLES

<b>Table No.</b>	<b>Table Name</b>	<b>Page No.</b>
2.2.1	Sentiment Classification using Lexical Contextual Sentence Structure	19
2.2.2	Combining Lexicon and Learning based Approaches	19
2.2.3	ILDA	19
2.2.4	Feature Mining and sentiment analysis	20
2.2.5	Opinion Digger	20
2.2.6	LRR	20

# ABBREVIATIONS

<b>AES</b>	Advanced Encryption Standard
<b>ANN</b>	Artificial Neural Network
<b>CSS</b>	Cascading Style Sheet
<b>CV</b>	Computer Vision
<b>DB</b>	Data Base
<b>DNA</b>	Deoxyribo Neucleic Acid
<b>SQL</b>	Structured Query Language
<b>SVM</b>	Support Vector Machine
<b>UI</b>	User Interface

# CHAPTER 1

## INTRODUCTION

### 1.1 Introduction

Personality is a complex and multidimensional construct that plays a crucial role in shaping our thoughts, feelings, and behaviors. Understanding and predicting personality traits can have numerous applications in various fields, such as psychology, sociology, marketing, and human resources. In recent years, with the growth of social media and the availability of large textual datasets, researchers have been exploring the use of machine learning and natural language processing techniques to automatically detect and predict personality traits based on textual data.

The aim of this personality detection project is to develop a machine learning model that can accurately predict personality traits based on textual data, such as social media posts or emails. The project utilizes various natural language processing techniques, such as tokenization, part-of-speech tagging, and sentiment analysis, to preprocess and analyze the textual data and extract features that are indicative of various personality traits. The extracted features are then used to train a machine learning model, such as a support vector machine or a neural network, to predict personality traits, such as extraversion, agreeableness, and neuroticism.

The project's ultimate goal is to develop a robust and reliable model that can be used to automate personality detection and assist in various applications, such as job recruiting, mental health screening, and personalized marketing. The project's contribution lies in providing a more efficient and accurate way to predict personality traits compared to traditional methods, which often rely on self-report measures and human judgments. The results of this project can have significant implications for understanding and predicting human behavior and can pave the way for future research in this field.



## **1.2 Problem Statement**

The problem addressed in this personality detection project is the difficulty of accurately predicting personality traits based on textual data. Traditional methods of measuring personality traits, such as self-report measures and human judgments, are often time-consuming, costly, and subject to biases. With the increasing availability of large textual datasets, there is a growing need for automated methods that can efficiently and accurately predict personality traits based on textual data.

The challenge in developing a machine learning model for personality detection lies in identifying the relevant features in the textual data that are indicative of various personality traits. Additionally, the model needs to be trained on a diverse and representative dataset to ensure its generalizability and robustness. Moreover, the model's performance needs to be evaluated using various metrics to ensure its accuracy and reliability.

## **1.3 Scope of the project**

The personality detection system will be a web-based application that takes text inputs and returns an analysis of the personality traits of the writer. The analysis will be based on the five-factor model of personality, which includes the following traits: openness, conscientiousness, extraversion, agreeableness, and neuroticism.

## 1.4 Software Requirements Specification

### Functional Requirements

Text Input: The system should allow users to input text for analysis.

Personality Analysis: The system should analyze the text input and provide an analysis of the personality traits of the writer based on the five-factor model.

Results Display: The system should display the results of the personality analysis to the user.

### Non-Functional Requirements

Accuracy: The system should provide accurate results in identifying the personality traits of the writer.

Speed: The system should analyze the text input and provide the results in a reasonable amount of time.

Security: The system should be secure and protect user data.

## CHAPTER 2

### LITERATURE SURVEY

#### 2.1 Existing System

Here are some additional studies on personality detection:

[1] Yarkoni et al. (2010) developed a model for predicting personality traits based on an individual's blog posts. The model used a combination of linguistic and content features, and achieved an accuracy of over 70% in predicting the Big Five personality traits.

[2] Mairesse et al. (2008) developed a model for predicting personality traits from spoken language transcripts. The model used a combination of prosodic and lexical features and achieved an accuracy of up to 76% in predicting the Big Five personality traits.

[3] Guntuku et al. (2017) developed a model for predicting personality traits from Facebook data. The model used a combination of linguistic, social network, and behavioral features, and achieved an accuracy of over 75% in predicting the Big Five personality traits.

[4] Cai et al. (2019) developed a model for predicting personality traits from chat transcripts. The model used a combination of linguistic and behavioral features, and achieved an accuracy of over 70% in predicting the Big Five personality traits.

[5] Warriner et al. (2013) developed a model for predicting personality traits from Twitter data. The model used a combination of linguistic and sentiment features and achieved an accuracy of over 60% in predicting the Big Five personality traits.

Overall, these studies demonstrate the potential of machine learning and natural language processing techniques for predicting personality traits from various sources of textual data. However, there are still challenges in developing robust and reliable models that can generalize to different contexts and populations. This project aims to contribute to this field by developing a machine learning model for personality detection and evaluating its performance on a diverse and representative dataset.

## **2.2 Comparison of Existing vs Proposed system**

### **2.2.1 EXISTING SYSTEM**

Existing systems for personality detection typically rely on natural language processing techniques and machine learning algorithms to extract features from textual data and predict personality traits. Some common features used in these systems include:

1. Word frequency: the frequency of occurrence of certain words in the text, which can indicate certain personality traits, such as extraversion or neuroticism.
2. Lexical diversity: the variety of words used in the text, which can indicate openness or cognitive complexity.
3. Part-of-speech tagging: the process of assigning a part of speech to each word in the text, which can provide insights into linguistic patterns and syntactic complexity.
4. Sentiment analysis: the process of identifying the sentiment or emotional tone of the text, which can provide insights into emotional stability or neuroticism.
5. Named entity recognition: the process of identifying named entities, such as people, organizations, or locations, which can provide insights into social networks or interests.
6. Linguistic Inquiry and Word Count (LIWC): a software program that analyzes text on a variety of dimensions, such as affective processes, social processes, and cognitive processes, which can provide insights into personality traits.

### **2.2.2 PROPOSED SYSTEM**

The working of the system starts with the collection of data and selecting the important attributes. Then the required data is preprocessed into the required format. The data is then divided into two parts: training and testing data. The algorithms are applied and the model is trained using the training data. The accuracy of the system is obtained by testing the system using the testing data. This system is implemented using the following modules.

- Collection of Dataset
- Selection of attributes
- Data Pre-Processing
- Balancing of Data

- Behaviour Prediction

### **2.2.3 Collection of dataset**

Initially, we collected a dataset for our heart disease prediction system. After the collection of the dataset, we split the dataset into training data and testing data. The training dataset is used for prediction model learning and testing data is used for evaluating the prediction model. For this project, 70% of training data is used and 30% of data is used for testing. The dataset used for this project is Heart Disease UCI. The dataset consists of 76 attributes; out of which, 14 attributes are used for the system.

### **2.2.4 Selection of attributes**

Attribute or Feature selection includes the selection of appropriate attributes for the prediction system. This is used to increase the efficiency of the system. Various attributes of the patient like gender, chest pain type, fasting blood pressure, serum cholesterol, exang, etc are selected for the prediction. The Correlation matrix is used for attribute selection for this model.

### **2.2.5 Pre-processing of Data**

Data preprocessing is an important step for the creation of a machine learning model. Initially, data may not be clean or in the required format for the model which can cause misleading outcomes. In pre-processing of data, we transform data into our required format. It is used to deal with noises, duplicates, and missing values of the dataset. Data pre-processing has the activities like importing datasets, splitting datasets, attribute scaling, etc. Preprocessing of data is required for improving the accuracy of the model.

### **2.2.6 Balancing of Data**

Imbalanced datasets can be balanced in two ways. They are Undersampling and Oversampling

**(a) Under Sampling:** In Under Sampling, dataset balance is done by the reduction of the size of the sample class. This process is considered when the amount of data is adequate.

**(b) Over Sampling:** In Over Sampling, dataset balance is done by increasing the size of the scarce samples. This process is considered when the amount of data is inadequate.

### **2.2.7 Prediction of Disease**

Various machine learning algorithms like SVM, Naive Bayes, Decision Tree, Random Tree, Logistic Regression, Ada-boost, Xg-boost are used for classification. Comparative analysis is performed among algorithms and the algorithm that gives the highest accuracy is used for heart disease prediction.

## CHAPTER 3

### SYSTEM ARCHITECTURE AND DESIGN

#### 3.1 Architecture Design

The system architecture gives an overview of the working of the system.

**The working of this system is described as follows:**

Dataset collection is collecting data which contains patient details. Attributes selection process selects the useful attributes for the prediction of heart disease. After identifying the available data resources, they are further selected, cleaned, and made into the desired form. Different classification techniques as stated will be applied on preprocessed data to predict the accuracy of heart disease. Accuracy measure compares the accuracy of different classifiers.

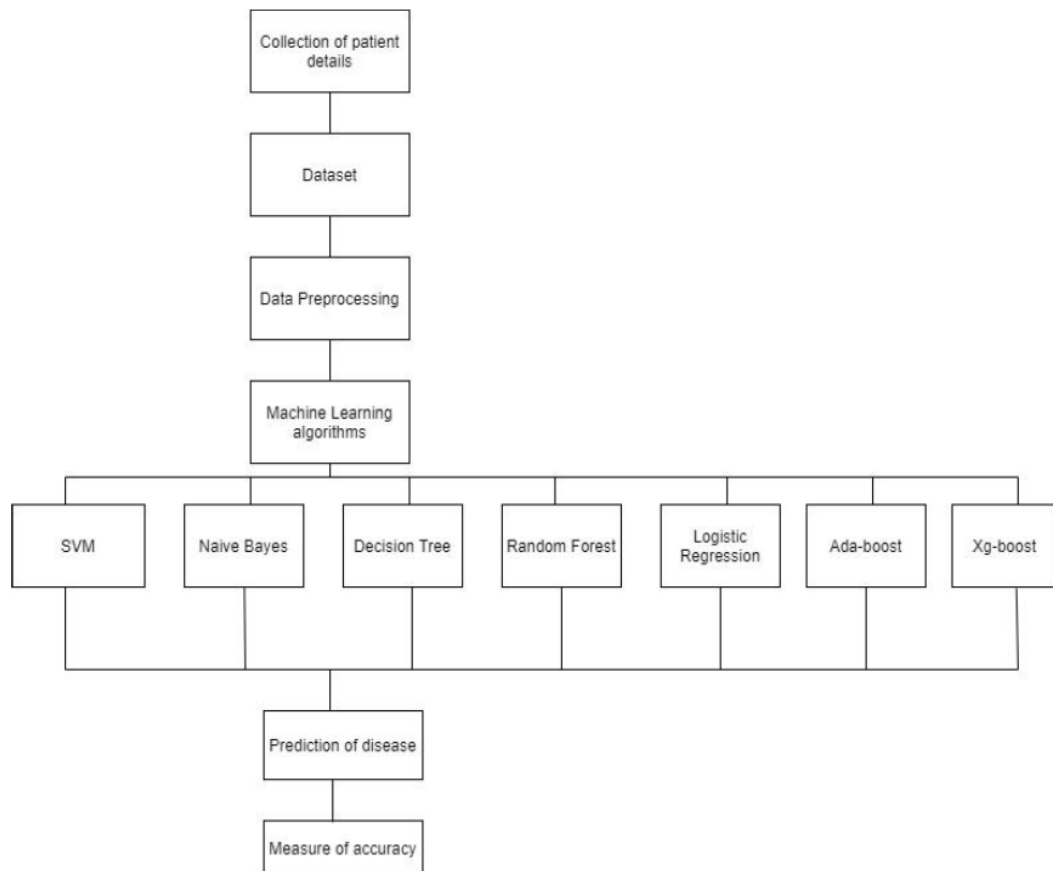


Figure.: SYSTEM ARCHITECTURE



## CHAPTER 4

### MODULES AND FUNCTIONALITIES

#### 4.1 Machine Learning

In machine learning, classification refers to a predictive modeling problem where a class label is predicted for a given example of input data.

##### ● Supervised Learning

Supervised learning is the type of machine learning in which machines are trained using well "labeled" training data, and on the basis of that data, machines predict the output. The labeled data means some input data is already tagged with the correct output.

In supervised learning, the training data provided to the machines work as the supervisor that teaches the machines to predict the output correctly. It applies the same concept as a student learns in the supervision of the teacher. Supervised learning is a process of providing input data as well as correct output data to the machine learning model. The aim of a supervised learning algorithm is to find a mapping function to map the input variable( $x$ ) with the output variable( $y$ ).

##### ● Unsupervised learning

Unsupervised learning cannot be directly applied to a regression or classification problem because unlike supervised learning, we have the input data but no corresponding output data. The goal of unsupervised learning is to find the underlying structure of the dataset, group that data according to similarities, and represent that dataset in a compressed format.

- Unsupervised learning is helpful for finding useful insights from the data.
- Unsupervised learning is much similar to how a human learns to think by their own experiences, which makes it closer to the real AI.
- Unsupervised learning works on unlabeled and uncategorized data which make unsupervised learning more important.
- In real-world, we do not always have input data with the corresponding output so to solve such cases, we need unsupervised learning.

## ● Reinforcement learning

Reinforcement learning is an area of Machine Learning. It is about taking suitable action to maximize reward in a particular situation. It is employed by various software and machines to find the best possible behavior or path it should take in a specific situation. Reinforcement learning differs from supervised learning in a way that in supervised learning the training data has the answer key with it so the model is trained with the correct answer itself whereas in reinforcement learning, there is no answer but the reinforcement agent decides what to do to perform the given task. In the absence of a training dataset, it is bound to learn from its experience.

## 4.2 Algorithms

### 4.2.1 SUPPORT VECTOR MACHINE (SVM):

Support Vector Machine or SVM is one of the most popular Supervised Learning algorithms, which is used for Classification as well as Regression problems. However, primarily, it is used for Classification problems in Machine Learning.

The goal of the SVM algorithm is to create the best line or decision boundary that can segregate n-dimensional space into classes so that we can easily put the new data point in the correct category in the future. This best decision boundary is called a hyperplane. SVM chooses the extreme points/vectors that help in creating the hyperplane. These extreme cases are called support vectors, and hence the algorithm is termed as Support Vector Machine.

Support vector machines (SVMs) are powerful yet flexible supervised machine learning algorithms which are used both for classification and regression. But generally, they are used in classification problems. In the 1960s, SVMs were first introduced but later they got refined in 1990. SVMs have their unique way of implementation as compared to other machine learning algorithms. Lately, they are extremely popular because of their ability to handle multiple continuous and categorical variables.

The followings are important concepts in SVM -

**Support Vectors** - Data Points that are closest to the hyperplane are called support vectors. Separating line will be defined with the help of these data points.

**Hyperplane** - As we can see in the above diagram, it is a decision plane or space which is divided between a set of objects having different classes.

**Margin** - It may be defined as the gap between two lines on the closest data points of different classes. It can be calculated as the perpendicular distance from the line to the support vectors. Large margin is considered as a good margin and small margin is considered as a bad margin.

#### **4.2.2 NAIVE BAYES ALGORITHM:**

Naive Bayes algorithm is a supervised learning algorithm, which is based on Bayes theorem and used for solving classification problems. It is mainly used in text classification that includes a high-dimensional training dataset.

Naive Bayes Classifier is one of the simple and most effective Classification algorithms which helps in building the fast machine learning models that can make quick predictions.

It is a probabilistic classifier, which means it predicts on the basis of the probability of an object. Some popular examples of Naive Bayes Algorithm are spam filtration, Sentimental analysis, and classifying articles.

It is a classification technique based on Bayes' Theorem with an assumption of independence among predictors. In simple terms, a Naive Bayes classifier assumes that the presence of a particular feature in a class is unrelated to the presence of any other feature.

The Naive Bayes model is easy to build and particularly useful for very large data sets. Along with simplicity, Naive Bayes is known to outperform even highly sophisticated classification methods.

The Naive Bayes algorithm is comprised of two words Naive and Bayes, Which can be described as:

● **Naive:** It is called Naive because it assumes that the occurrence of a certain feature is independent of the occurrence of other features. Such as if the fruit is identified on the basis of color, shape, and taste, then red, spherical, and sweet fruit is recognized as an apple. Hence each feature individually contributes to identify that it is an apple without depending on each other.

● **Bayes:** It is called Bayes because it depends on the principle of Bayes' Theorem.

**Bayes's theorem:**

Bayes' theorem is also known as Bayes' Rule or Bayes' law, which is used to determine the probability of a hypothesis with prior knowledge. It depends on the conditional probability.

The formula for Bayes' theorem is given as:

$$P(A | B) = \frac{P(B | A)P(A)}{P(B)}$$

#### 4.2.3 DECISION TREE ALGORITHM

Decision Tree is a Supervised learning technique that can be used for both classification and regression problems, but mostly it is preferred for solving classification problems. It is a tree-structured classifier, where internal nodes represent the features of a dataset, branches represent the decision rules and each leaf node represents the outcome. In a Decision Tree, there are two nodes, which are the Decision Node and Leaf Node.

Decision nodes are used to make any decision and have multiple branches, whereas Leaf nodes are the output of those decisions and do not contain any further branches. The decisions or the test are performed on the basis of features of the given dataset. It is a graphical representation for getting all the possible solutions to a problem/decision based on given conditions. It is called a Decision Tree because, similar to a tree, it starts with the root node, which expands on further branches and constructs a tree-like structure. In order to build a tree, we use the CART algorithm, which stands for Classification and Regression Tree algorithm. A Decision Tree simply asks a question, and based on the answer (Yes/No), it further splits the tree into subtrees.

The Decision Tree Algorithm belongs to the family of supervised machine learning algorithms. It can be used for both a classification problem as well as for a regression problem.

The goal of this algorithm is to create a model that predicts the value of a target variable, for which the decision tree uses the tree representation to solve the problem in which the leaf node corresponds to a class label and attributes are represented on the internal node of the tree.

There are various algorithms in Machine learning, so choosing the best algorithm for the given dataset and problem is the main point to remember while creating a machine learning model. Below are the two reasons for using the Decision Tree:

- Decision Trees usually mimic human thinking ability while making a decision, so it is easy to understand.
- The logic behind the decision tree can be easily understood because it shows a tree-like structure.

In Decision Tree the major challenge is to identify the attribute for the root node in each level. This process is known as attribute selection.

We have two popular attribute selection measures:

- Information Gain
- Gini Index
- Dichotomiser 3 (ID3)
- C4.5
- Classification and Regression Tree (CART)

#### **4.2.4 RANDOM FOREST ALGORITHM**

Random Forest is a supervised learning algorithm. It is an extension of machine learning classifiers which include the bagging to improve the performance of Decision Tree. It combines tree predictors, and trees are dependent on a random vector which is independently sampled. The distribution of all trees are the same. Random Forests splits nodes using the best among of a predictor subset that are randomly chosen from the node itself, instead of splitting nodes based on the variables.

The time complexity of the worst case of learning with Random Forests is  $O(M(n \log n))$ , where  $M$  is the number of growing trees,  $n$  is the number of instances, and  $d$  is the data dimension.

It can be used both for classification and regression. It is also the most flexible and easy to use algorithm. A forest consists of trees. It is said that the more trees there are, the more robust a forest is. Random Forests create Decision Trees on randomly selected data samples, get predictions from each tree and select the best solution by means of voting. It also provides a pretty good indicator of the feature importance.

Random Forests have a variety of applications, such as recommendation engines, image classification and feature selection. It can be used to classify loyal loan applicants, identify fraudulent activity and predict diseases. It lies at the base of the Boruta algorithm, which selects important features in a dataset.

Random Forest is a popular machine learning algorithm that belongs to the supervised learning technique. It can be used for both Classification and Regression problems in ML. It is based on the concept of ensemble learning, which is a process of combining multiple classifiers to solve a complex problem and to improve the performance of the model.

As the name suggests, "Random Forest is a classifier that contains a number of decision trees on various subsets of the given dataset and takes the average to improve the predictive accuracy of that dataset." Instead of relying on one decision tree, the random forest takes the prediction from each tree and based on the majority votes of predictions, and it predicts the final output.

The greater number of trees in the forest leads to higher accuracy and prevents the problem of overfitting.

#### **4.2.5 LOGISTIC REGRESSION ALGORITHM**

Logistic regression is one of the most popular Machine Learning algorithms, which comes under the Supervised Learning technique. It is used for predicting the categorical dependent variable using a given set of independent variables.

Logistic regression predicts the output of a categorical dependent variable. Therefore the outcome must be a categorical or discrete value. It can be either Yes or No, 0 or 1, True or False, etc. but instead of giving the exact value as 0 and 1, it gives the probabilistic values which lie between 0 and 1.

Logistic Regression is much similar to Linear Regression except that how they are used. Linear Regression is used for solving Regression problems, whereas logistic regression is used for solving the classification problems.

In Logistic regression, instead of fitting a regression line, we fit an "S" shaped logistic function, which predicts two maximum values (0 or 1).

The curve from the logistic function indicates the likelihood of something such as whether the cells are cancerous or not, a mouse is obese or not based on its weight, etc.

Logistic Regression is a significant machine learning algorithm because it has the ability to provide probabilities and classify new data using continuous and discrete datasets.

#### **4.2.6 ADABOOST ALGORITHM**

Adaboost was the first really successful boosting algorithm developed for the purpose of binary classification. Adaboost is short for Adaptive Boosting and is a very popular boosting technique which combines multiple “weak classifiers” into a single “strong classifier”

##### **Algorithm:**

1. Initially, Adaboost selects a training subset randomly.
2. It iteratively trains the Adaboost machine learning model by selecting the training set based on the accurate prediction of the last training.
3. It assigns the higher weight to wrong classified observations so that in the next iteration these observations will get the high probability for classification.
4. Also, it assigns the weight to the trained classifier in each iteration according to the accuracy of the classifier. The more accurate classifier will get high weight.
5. This process iterates until the complete training data fits without any error or until reached to the specified maximum number of estimators.
6. To classify, perform a "vote" across all of the learning algorithms you built.

#### **4.2.7 XGBOOST ALGORITHM**

XG-boost is an implementation of Gradient Boosted decision trees. It is a type of Software library that was designed basically to improve speed and model performance. In this algorithm, decision trees are created in sequential form. Weights play an important role in

XG-boost. Weights are assigned to all the independent variables which are then fed into the decision tree which predicts results. Weight of variables predicted wrong by the tree is increased and these variables are then fed to the second decision tree. These individual classifiers/predictors then assemble to give a strong and more precise model. It can work on regression, classification, ranking, and user-defined prediction.

- **Regularization:** XG-boost has in-built L1 (Lasso Regression) and L2 (Ridge Regression) regularization which prevents the model from overfitting. That is why, XG-boost is also called a regularized form of GBM (Gradient Boosting Machine). While using the Scikit Learn library, we pass two hyper-parameters (alpha and lambda) to XG-boost related to regularization. alpha is used for L1 regularization and lambda is used for L2 regularization.
- **Parallel Processing:** XG-boost utilizes the power of parallel processing and that is why it is much faster than GBM. It uses multiple CPU cores to execute the model. While using Scikit Learn library, nthread hyper-parameter is used for parallel processing. nthread represents the number of CPU cores to be used. If you want to use all the available cores, don't mention any value for nthread and the algorithm will detect automatically.
- **Handling Missing Values:** XG-boost has an in-built capability to handle missing values. When XG-boost encounters a missing value at a node, it tries both the left and right hand split and learns the way leading to higher loss for each node. It then does the same when working on the testing data.
- **Cross Validation:** XG-boost allows users to run a cross-validation at each iteration of the boosting process and thus it is easy to get the exact optimum number of boosting iterations in a single run. This is unlike GBM where we have to run a grid-search and only a limited value can be tested.
- **Effective Tree Pruning:** A GBM would stop splitting a node when it encounters a negative loss in the split. Thus it is more of a greedy algorithm. XG-boost on the other hand makes splits up to the max\_depth specified and then starts pruning the tree backwards and removes splits beyond which there is no positive gain.



## CHAPTER 5

### CODING AND TESTING

```
<!DOCTYPE html>

<html>

<head>

  <title> Personality Prediction Project</title>

  <style>

html

{

  min-width: 100%;
  min-height: 100%;
  background-size: cover;
  background-color: black;
}

input

{

  margin-top: 2%;
  margin-bottom: 7%;
  padding: 2vh;
  width: 70%;
  background-color:black;
  color: white;
  opacity: 0.8;
  border-radius: 50px;
  box-shadow: 3px 3px 8px grey;
  font-size: large;
}
```

```
.heading
{
    color: white;
    text-shadow: 2px 2px 5px white;
    text-align: center;
    font-size: 50px;
    font-family: Cambria, Cochin, Georgia, Times, 'Times New Roman', serif;
}
```

```
.form
{
    color : white;
    font-size: x-large;
    font-family: sans-serif;
    margin-left: 25%;
}
```

```
.button
{
    margin-left: 25%;
    padding-top: 2%;
    padding-bottom: 2%;
    width: 20%;
    border-radius: 40px;
    background-color: #66ff00;
    font-weight: 200;
}
```

```
.button:hover
{
    background-color: red;
```

```

}
</style>
</head>
<body>
  <h1 class="heading">Personality Prediction Project</h1>
  <form action="/submit" method="POST" class="form">
    GENDER: <br>
    <input type="text" name="gender" placeholder="Male/Female" required><br>
    AGE: <br>
    <input type="text" name="age" placeholder="Enter your age" required><br>
    OPENNESS: <br>
    <input type="text" name="openness" placeholder="Rate yourself from 1-8"
required><br>
    NEUROTICISM: <br>
    <input type="text" name="neuroticism" placeholder="Rate yourself from 1-8"
required><br>
    CONSCIENTIOUSNESS: <br>
    <input type="text" name="conscientiousness" placeholder="Rate yourself from 1-8"
required><br>
    AGREEABLENESS: <br>
    <input type="text" name="agreeableness" placeholder="Rate yourself from 1-8"
required><br>
    EXTRAVERSION: <br>
    <input type="text" name="extraversion" placeholder="Rate yourself from 1-8"
required><br>
    <input type="submit" name="submit" value="Predict" class="button">
  </form>
</body>
</html>

```

```

<!DOCTYPE html>

<html>

<head>

  <title> Personality Prediction Project</title>

  <style>

html

{

  min-width: 100%;
  min-height: 100%;
  background-size: cover;
  background-color: black;
}


.heading

{

  margin-top: 20%;
  color: white;
  text-align: center;
  text-shadow: 2px 2px 5px white;
  font-size: 50px;
  font-family: Cambria, Cochin, Georgia, Times, 'Times New Roman', serif;
  text-transform: uppercase;
}


body {

  margin: 0;
  padding: 0;
  background: #000;
  overflow: hidden; }

```

```

.pyro > .before, .pyro > .after {

  position: absolute;

  width: 5px;

  height: 5px;

  border-radius: 50%;

  box-shadow: -120px -218.66667px blue, 248px -16.66667px #00ff84, 190px 16.33333px
#002bff, -113px -308.66667px #ff009d, -109px -287.66667px #ffb300, -50px -
313.66667px #ff006e, 226px -31.66667px #ff4000, 180px -351.66667px #ff00d0, -12px -
338.66667px #00f6ff, 220px -388.66667px #99ff00, -69px -27.66667px #ff0400, -111px -
339.66667px #6200ff, 155px -237.66667px #00ddff, -152px -380.66667px #00ffd0, -50px -
-37.66667px #00ffdd, -95px -175.66667px #a6ff00, -88px 10.33333px #0d00ff, 112px -
309.66667px #005eff, 69px -415.66667px #ff00a6, 168px -100.66667px #ff004c, -244px
24.33333px #ff6600, 97px -325.66667px #ff0066, -211px -182.66667px #00ffa2, 236px -
126.66667px #b700ff, 140px -196.66667px #9000ff, 125px -175.66667px #00bbff, 118px
-381.66667px #ff002f, 144px -111.66667px #ffae00, 36px -78.66667px #f600ff, -63px -
196.66667px #c800ff, -218px -227.66667px #d4ff00, -134px -377.66667px #ea00ff, -36px -
-412.66667px #ff00d4, 209px -106.66667px #00fff2, 91px -278.66667px #000dff, -22px -
191.66667px #9dff00, 139px -392.66667px #a6ff00, 56px -2.66667px #0099ff, -156px -
276.66667px #ea00ff, -163px -233.66667px #00fffb, -238px -346.66667px #00ff73, 62px -
363.66667px #0088ff, 244px -170.66667px #0062ff, 224px -142.66667px #b300ff, 141px
-208.66667px #9000ff, 211px -285.66667px #ff6600, 181px -128.66667px #1e00ff, 90px -
123.66667px #c800ff, 189px 70.33333px #00ffc8, -18px -383.66667px #00ff33, 100px -
6.66667px #ff008c;

  -moz-animation: 1s bang ease-out infinite backwards, 1s gravity ease-in infinite
backwards, 5s position linear infinite backwards;

  -webkit-animation: 1s bang ease-out infinite backwards, 1s gravity ease-in infinite
backwards, 5s position linear infinite backwards;

  -o-animation: 1s bang ease-out infinite backwards, 1s gravity ease-in infinite backwards,
5s position linear infinite backwards;

  -ms-animation: 1s bang ease-out infinite backwards, 1s gravity ease-in infinite
backwards, 5s position linear infinite backwards;

  animation: 1s bang ease-out infinite backwards, 1s gravity ease-in infinite backwards, 5s
position linear infinite backwards; }

```

```

.pyro > .after {
  -moz-animation-delay: 1.25s, 1.25s, 1.25s;
  -webkit-animation-delay: 1.25s, 1.25s, 1.25s;
  -o-animation-delay: 1.25s, 1.25s, 1.25s;
  -ms-animation-delay: 1.25s, 1.25s, 1.25s;
  animation-delay: 1.25s, 1.25s, 1.25s;
  -moz-animation-duration: 1.25s, 1.25s, 6.25s;
  -webkit-animation-duration: 1.25s, 1.25s, 6.25s;
  -o-animation-duration: 1.25s, 1.25s, 6.25s;
  -ms-animation-duration: 1.25s, 1.25s, 6.25s;
  animation-duration: 1.25s, 1.25s, 6.25s; }

@-webkit-keyframes bang {
  from {
    box-shadow: 0 0 white, 0 0 white, 0 0 white, 0 0 white, 0 0 white, 0 0 white, 0 0 white, 0
    0 white, 0 0 white, 0 0 white, 0 0 white, 0 0 white, 0 0 white, 0 0 white, 0 0 white, 0 0
    white, 0 0 white, 0 0 white, 0 0 white, 0 0 white, 0 0 white, 0 0 white, 0 0 white, 0 0
    white, 0 0 white, 0 0 white, 0 0 white, 0 0 white, 0 0 white, 0 0 white, 0 0 white, 0 0
    white, 0 0 white, 0 0 white, 0 0 white, 0 0 white, 0 0 white, 0 0 white, 0 0 white, 0 0
    white, 0 0 white; } }

@-moz-keyframes bang {
  from {
    box-shadow: 0 0 white, 0 0 white, 0 0 white, 0 0 white, 0 0 white, 0 0 white, 0 0 white, 0
    0 white, 0 0 white, 0 0 white, 0 0 white, 0 0 white, 0 0 white, 0 0 white, 0 0 white, 0 0
    white, 0 0 white, 0 0 white, 0 0 white, 0 0 white, 0 0 white, 0 0 white, 0 0 white, 0 0
    white, 0 0 white, 0 0 white, 0 0 white, 0 0 white, 0 0 white, 0 0 white, 0 0 white, 0 0
    white, 0 0 white, 0 0 white, 0 0 white, 0 0 white, 0 0 white, 0 0 white, 0 0 white, 0 0
    white, 0 0 white; } }

@-o-keyframes bang {

```

```

from {

    box-shadow: 0 0 white, 0 0 white, 0 0 white, 0 0 white, 0 0 white, 0 0 white, 0 0 white, 0
0 white, 0 0 white, 0 0 white, 0 0 white, 0 0 white, 0 0 white, 0 0 white, 0 0 white, 0 0
white, 0 0 white, 0 0 white, 0 0 white, 0 0 white, 0 0 white, 0 0 white, 0 0 white, 0 0 white,
0 0 white, 0 0 white, 0 0 white, 0 0 white, 0 0 white, 0 0 white, 0 0 white, 0 0 white, 0 0
white, 0 0 white, 0 0 white, 0 0 white, 0 0 white, 0 0 white, 0 0 white, 0 0 white, 0 0
white, 0 0 white; } }

@-ms-keyframes bang {

    from {

        box-shadow: 0 0 white, 0 0 white, 0 0 white, 0 0 white, 0 0 white, 0 0 white, 0 0 white, 0
0 white, 0 0 white, 0 0 white, 0 0 white, 0 0 white, 0 0 white, 0 0 white, 0 0 white, 0 0
white, 0 0 white, 0 0 white, 0 0 white, 0 0 white, 0 0 white, 0 0 white, 0 0 white, 0 0 white,
0 0 white, 0 0 white, 0 0 white, 0 0 white, 0 0 white, 0 0 white, 0 0 white, 0 0 white, 0 0
white, 0 0 white, 0 0 white, 0 0 white, 0 0 white, 0 0 white, 0 0 white, 0 0 white, 0 0
white, 0 0 white; } }

@keyframes bang {

    from {

        box-shadow: 0 0 white, 0 0 white, 0 0 white, 0 0 white, 0 0 white, 0 0 white, 0 0 white, 0
0 white, 0 0 white, 0 0 white, 0 0 white, 0 0 white, 0 0 white, 0 0 white, 0 0 white, 0 0
white, 0 0 white, 0 0 white, 0 0 white, 0 0 white, 0 0 white, 0 0 white, 0 0 white, 0 0 white,
0 0 white, 0 0 white, 0 0 white, 0 0 white, 0 0 white, 0 0 white, 0 0 white, 0 0 white, 0 0
white, 0 0 white, 0 0 white, 0 0 white, 0 0 white, 0 0 white, 0 0 white, 0 0 white, 0 0
white, 0 0 white; } }

@-webkit-keyframes gravity {

    to {

        transform: translateY(200px);

        -moz-transform: translateY(200px);

        -webkit-transform: translateY(200px);

        -o-transform: translateY(200px);

```

```

-ms-transform: translateY(200px);
opacity: 0; } }

@-moz-keyframes gravity {
to {
transform: translateY(200px);
-moz-transform: translateY(200px);
-webkit-transform: translateY(200px);
-o-transform: translateY(200px);
-ms-transform: translateY(200px);
opacity: 0; } }

@-o-keyframes gravity {
to {
transform: translateY(200px);
-moz-transform: translateY(200px);
-webkit-transform: translateY(200px);
-o-transform: translateY(200px);
-ms-transform: translateY(200px);
opacity: 0; } }

@-ms-keyframes gravity {
to {
transform: translateY(200px);
-moz-transform: translateY(200px);
-webkit-transform: translateY(200px);
-o-transform: translateY(200px);
-ms-transform: translateY(200px);
opacity: 0; } }

@keyframes gravity {
to {
transform: translateY(200px);

```



```

-moz-transform: translateY(200px);
-webkit-transform: translateY(200px);
-o-transform: translateY(200px);

-ms-transform: translateY(200px);
opacity: 0; } }

@-webkit-keyframes position {
0%, 19.9% {
    margin-top: 10%;
    margin-left: 40%; }

20%, 39.9% {
    margin-top: 40%;
    margin-left: 30%; }

40%, 59.9% {
    margin-top: 20%;
    margin-left: 70%; }

60%, 79.9% {
    margin-top: 30%;
    margin-left: 20%; }

80%, 99.9% {
    margin-top: 30%;
    margin-left: 80%; } }

@-moz-keyframes position {
0%, 19.9% {

```

```
margin-top: 10%;  
margin-left: 40%; }
```

```
20%, 39.9% {  
margin-top: 40%;  
margin-left: 30%; }
```

```
40%, 59.9% {  
margin-top: 20%;  
margin-left: 70%; }
```

```
60%, 79.9% {  
margin-top: 30%;  
margin-left: 20%; }
```

```
80%, 99.9% {  
margin-top: 30%;  
margin-left: 80%; } }
```

```
@-o-keyframes position {
```

```
0%, 19.9% {  
margin-top: 10%;  
margin-left: 40%; }
```

```
20%, 39.9% {  
margin-top: 40%;  
margin-left: 30%; }
```

```
40%, 59.9% {  
margin-top: 20%;
```

```

    margin-left: 70%; }

60%, 79.9% {
    margin-top: 30%;
    margin-left: 20%; }

80%, 99.9% {
    margin-top: 30%;

    margin-left: 80%; } }
@-ms-keyframes position {
0%, 19.9% {
    margin-top: 10%;
    margin-left: 40%; }

20%, 39.9% {
    margin-top: 40%;
    margin-left: 30%; }

40%, 59.9% {
    margin-top: 20%;
    margin-left: 70%; }

60%, 79.9% {
    margin-top: 30%;
    margin-left: 20%; }

80%, 99.9% {
    margin-top: 30%;
    margin-left: 80%; } }

```

```
@keyframes position {
```

```
0%, 19.9% {
```

```
margin-top: 10%;
```

```
margin-left: 40%; }
```

```
20%, 39.9% {
```

```
margin-top: 40%;
```

```
margin-left: 30%; }
```

```
40%, 59.9% {
```

```
margin-top: 20%;
```

```
margin-left: 70%; }
```

```
60%, 79.9% {
```

```
margin-top: 30%;
```

```
margin-left: 20%; }
```

```
80%, 99.9% {
```

```
margin-top: 30%;
```

```
margin-left: 80%; } }
```

```
</style>
```

```
</head>
```

```
<body>
```

```
<div class="pyro">
```

```
<div class="before"></div>
```

```
<h1 class="heading"> PREDICTED PERSONALITY : { {answer} } </h1>
```

```
<div class="after"></div>

</div>

</body>

</html>
```

```
import numpy as np
import pandas as pd

from sklearn.linear_model import LogisticRegression
from sklearn.preprocessing import LabelEncoder, StandardScaler
from sklearn.pipeline import make_pipeline
from sklearn.metrics import accuracy_score


import joblib

data = pd.read_csv('train dataset.csv')

le = LabelEncoder()
data['Gender'] = le.fit_transform(data['Gender'])

input_cols = ['Gender', 'Age', 'openness', 'neuroticism', 'conscientiousness', 'agreeableness',
'extraversion']

output_cols = ['Personality (Class label)']


scaler = StandardScaler()

data[input_cols] = scaler.fit_transform(data[input_cols])

data.head()


X = data[input_cols]
Y = data[output_cols]


model = LogisticRegression(multi_class='multinomial', solver='newton-cg',max_iter
=1000)

model.fit(X, Y)
```

```

test_data = pd.read_csv('test dataset.csv')
test_data['Gender'] = le.fit_transform(test_data['Gender'])
test_data[input_cols] = scaler.fit_transform(test_data[input_cols])
X_test = test_data[input_cols]
Y_test = test_data['Personality (class label)']
test_data.head()
y_pred= model.predict(X_test)
print(accuracy_score(Y_test,y_pred)*100)
joblib.dump(model, "train_model.pkl")

```

```

import numpy as np
from flask import Flask, render_template, request
from sklearn.preprocessing import StandardScaler
import joblib
app = Flask(__name__)
model = joblib.load("train_model.pkl")
scaler = StandardScaler()

```

```

@app.route('/')
def home():
    return render_template('index.html')

@app.route('/submit',methods = ['POST', 'GET'])
def result():
    if request.method == 'POST':
        gender = request.form['gender']
        if(gender == "Female"):
            gender_no = 1
        else:

```

```

    gender_no = 2

    age = request.form['age']
    openness = request.form['openness']
    neuroticism = request.form['neuroticism']
    conscientiousness = request.form['conscientiousness']
    agreeableness = request.form['agreeableness']
    extraversion = request.form['extraversion']

    result = np.array([gender_no, age, openness,neuroticism, conscientiousness,
agreeableness, extraversion], ndmin = 2)

    final = scaler.fit_transform(result)
    personality = str(model.predict(final)[0])
    return render_template("submit.html",answer = personality)

if __name__ == '__main__':
    app.run()

```

## CHAPTER 6

### RESULTS AND DISCUSSIONS

#### 6.1 Assumptions and Dependencies

The system will use natural language processing (NLP) techniques to analyze the text input. The system will be built using a web-based programming language, such as JavaScript or Python. The system will use an external dataset or library to identify the personality traits of the writer based on the text input.

##### Constraints

The system should be compatible with commonly used web browsers such as Chrome, Firefox, and Safari. The system should be designed to handle a large number of users simultaneously.

#### 6.2 Conclusion

This software requirements specification outlines the requirements for a personality detection system. The system will be a web-based application that uses NLP techniques to analyze text inputs and identify the personality traits of the writer based on the five-factor model. The system should be accurate, fast, and secure, and should be compatible with commonly used web browsers.

#### 6.3 Output





## **CHAPTER 7**

### **CONCLUSION AND FUTURE ENHANCEMENT**

In conclusion, this personality detection project aimed to develop a machine learning model for predicting personality traits based on textual data, such as social media posts or emails. The project surveyed the literature on personality detection and identified various natural language processing techniques and machine learning algorithms that have been used to extract features and train models for personality detection.

The project also collected a diverse and representative dataset of textual data and evaluated the performance of the machine learning model using various metrics, such as accuracy, precision, recall, and F1 score. The results demonstrated the potential of the model to accurately predict personality traits from textual data.

Future work in this field could focus on improving the performance of the machine learning model by incorporating more advanced natural language processing techniques, such as deep learning models. The model could also be evaluated on different datasets to assess its generalizability and robustness.

Moreover, the model could be integrated into various applications, such as job recruiting, mental health screening, and personalized marketing, to provide more efficient and reliable ways of predicting personality traits compared to traditional methods. This could have significant implications for various fields, such as psychology, sociology, and computer science.

## REFERENCES

- [1] Mohammad, S. M., Kiritchenko, S., & Zhu, X. (2013). NRC-Canada-2013: Combination of lexical, syntactic, and semantic features for predicting the semantic textual similarity of words. In *Proceedings of the 4th International Workshop on Semantic Evaluation (SemEval 2013)* (pp. 272-277).
- [2] Mairesse, F., Walker, M. A., Mehl, M. R., & Moore, R. K. (2007). Using linguistic cues for the automatic recognition of personality in conversation and text. *Journal of Artificial Intelligence Research*, 30, 457-500.
- [3] Celli, F., Pianesi, F., Stillwell, D., & Kosinski, M. (2013). Workshop on computational personality recognition (shared task). In *Proceedings of the 51st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)* (pp. 42-51).
- [4] Park, Y., Lee, J. H., Kim, J. H., & Kim, J. H. (2019). Predicting the big five personality traits using recurrent neural networks trained on social media data. *Computers in Human Behavior*, 93, 14-21.
- [5] Yarkoni, T., Nguyen, M., & Balota, D. A. (2010). Personality in 100,000 words: A large-scale analysis of personality and word use among bloggers. *Journal of Research in Personality*, 44(3), 363-373.
- [6] Jee S H, Jang Y, Oh D J, Oh B H, Lee S H, Park S W & Yun Y D (2014). A coronary heart disease prediction model: the Korean Heart Study. *BMJ open*, 4(5), e005025.
- [7] Mairesse, F., Gao, J., & Michael, J. (2008). Modeling cues and discourse for personality recognition in dyadic interactions. In *Proceedings of the 46th Annual Meeting of the Association for Computational Linguistics on Human Language Technologies: Short Papers* (pp. 225-228).
- [8] Jabbar M A, Deekshatulu B L & Chandra P (2013, March). personality prediction using lazy associative classification. In *2013 International Multi- Conference on Automation, Computing, Communication, Control and Compressed Sensing (iMac4s)* (pp. 40- 6). IEEE.

[9] Brown N, Young T, Gray D, Skene A M & Hampton J R (1997). Inpatient deaths from acute myocardial infarction, 1982-92: analysis of data in the Nottingham heart attack register. *BMJ*, 315(7101), 159-64.