

Machine Learning Model for Schizophrenia diagnosis using Electroencephalogram (EEG)

By

Korir Ruto Emmanuel

072350

**A Thesis submitted to the School of Computing and Engineering Sciences in Partial
Fulfilment of the Requirement of degree of Master of Science in Information Technology at
Strathmore University**

Strathmore University 2024

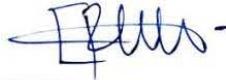
Declaration/ Approval

This dissertation is my original work and has not been presented for an award of any degree in any university.

Name: Korir Ruto Emmanuel

Admission Number :072350

Signature:



Date: 8th April 2024

This dissertation has been submitted to the School of Computing and Engineering Sciences for Examination with my approval as the University Supervisor

Name: Dr. Joseph Onderi Orero

Signature:



Date:8th April 2024.....

Abstract

Schizophrenia is a chronic mental condition that affects a significant number of the population, the disease shows psychotic symptoms which manifest as delusions, hallucinations and cognitive deficits. The signs and severity of the symptoms often progresses quite rapidly though individual patients may experience varying levels of progressions. Early identification of the onset and advanced stages of the mental illness is crucial for prompt and efficient therapy to stop or lessen the disease's progression. Currently the main means of diagnosis of different mental illnesses (including Schizophrenia) can be found in the Diagnostic and Statistical Manual of Mental Disorders (DSM-5). The patient is questioned about the disease's symptoms and prognosis as part of the conventional diagnostic procedure. Clinical assessments, including the Positive and Negative Syndrome Scale (PANSS) is often used to gauge how severe the symptoms are. Even seasoned psychiatrists may find it challenging to regularly diagnose patients accurately, largely because of the erroneous information given by patients. To increase the accuracy of the diagnosis, it is crucial to create an objective technique that measures the intensity of the symptoms using quantitative biomarkers. EEG signals can be a helpful supplementary diagnostic tool for people with schizophrenia to be used by psychiatrists. It has been observed that time-frequency modification of EEG signals obtained from electrodes is enough for the detection of schizophrenia.

This study proposes the use of machine learning methodology to build a Schizophrenia Detection model through ERP (Event Related Potential) data obtained from EEG (Electroencephalograph).

The data was sourced from an experiment conducted to determine a lack of neural and motor suppression from internal stimulus in patients with Schizophrenia. The data was used to train various machine learning algorithms and KNN (K Nearest Neighbor) provided the highest accuracy with 93%.

Keywords: *Schizophrenia detection, Schizophrenia prediction, Electroencephalography, EEG, ERP, Event Related Potential, Machine Learning, mental disorder*

Table of Contents

Declaration/ Approval.....	i
Abstract.....	ii
Table of Contents.....	iii
Table of Figures.....	vii
List of Abbreviations.....	ix
Acknowledgements.....	x
CHAPTER 1: INTRODUCTION.....	1
1.1 Background of the study.....	1
1.2 Problem Statement.....	3
1.3 Objectives.....	3
1.3.1 General Objectives.....	3
1.3.2 Specific Objectives.....	3
1.4 Research Questions.....	4
1.5 Justification.....	4
1.6 Scope and Limitation.....	4
CHAPTER 2: LITERATURE REVIEW.....	5
2.1 Introduction.....	5
2.2 Schizophrenia.....	5
2.3 Efference Copy /Corollary Discharge.....	6
2.4 EEG (Electroencephalography) and ERP (Event Related Potentials).....	8
2.5 ERP -Event related Potential.....	9
2.6 Machine Learning.....	10
2.6.1 Machine Learning Techniques.....	11

2.6.1.1	K Nearest Neighbor	11
2.6.1.2	Decision Trees Classifier	12
2.6.1.3	Random Forest Classifier	15
2.6.1.4	Convolutional Neural Networks	16
2.6.2	Integration of EEG and Machine Learning.....	18
2.7	Using EEG in Schizophrenia Detection	19
2.8	Related works.....	20
2.8.1	Schizophrenia Detection Using Machine Learning Approach from Social Media Content	20
2.8.2	Diagnosing Schizophrenia with Network Analysis and a Machine learning method	21
2.8.3	Application of principal component analysis to separate Schizophrenic Patients from Healthy Controls using fractional anisotropy images.....	22
2.8.4	Differentiation of Schizophrenia by Combining the Spatial EEG Brain Network Patterns of Rest and Task P300	23
2.9	Research Gap.....	27
2.10	Conceptual Framework	27
CHAPTER 3: RESEARCH METHODOLOGY		28
3.1	Introduction	28
3.2	Research Design.....	28
3.3	Datasets and Data Collection	28
3.4	Data Preprocessing.....	29
3.5	System Development Methodology	30
3.6	Research Quality /Validation & Reliability	31
3.7	Model Validation.....	31
3.8	Ethical Considerations.....	31

CHAPTER 4: SYSTEM ANALYSIS, DESIGN AND ARCHITECTURE	32
4.1 Introduction	32
4.2 Requirements Analysis.....	32
4.2.1 Functional Requirements	32
4.2.2 Non-Functional Requirements	32
4.3 System Architecture	33
4.4 System Diagrams.....	34
4.4.1 Use Case Diagram.....	34
4.4.2 System Sequence Diagram	38
4.4.3 Entity Relationship Diagram.....	39
4.4.4 Wireframes.....	41
CHAPTER 5: SYSTEM IMPLEMENTATION AND TESTING	43
5.1 Introduction	43
5.2 System Development.....	43
5.2.1 Hardware and Software Environment.....	43
5.3 System Implementation.....	45
5.3.1 Data Collection	45
5.3.2 Data preprocessing.....	45
5.3.3 Training the model.....	47
5.3.4 Validating the Model	50
CHAPTER 6: DISCUSSIONS	52
6.1 Introduction	52
6.2 Model Validation.....	52
6.3 Model Performance	52
6.4 Contribution to Research.....	53

CHAPTER 7: CONCLUSION AND RECOMMENDATION	54
7.1 Conclusion.....	54
7.2 Recommendation.....	55
7.3 Future work	55
REFERENCES.....	56

Table of Figures

Figure 2 1 Disconnect of the corollary discharge in schizophrenic patients (Ford & Mathalon, 2005)).....	7
Figure 2 2 Illustration of EEG Data Recording (Adapted from Brightbraincenter.co.uk)	8
Figure 2 3 Difference between Supervised and Unsupervised Learning (Adapted from IBM) ...	11
Figure 2 4 K Nearest Neighbor Classification	12
Figure 2 5 Decision Tree Classifier	13
Figure 2 6 Random Forest Classifier	16
Figure 2 7 Artificial Neural Networks	17
Figure 2 8 Convolutional Neural Networks.....	18
Figure 2 9 Illustration of N1 Suppression in Schizophrenic Patients (Adapted from (Ford et al., 2014)).....	19
Figure 2 10 Classification Performance of Machine Learning Models (Adapted from (Bae et al., 2021)).....	21
Figure 2 11 Performance of ML Models (Adapted from (Jo et al., 2020))	22
Figure 2 12 Comparison of PCA, DPCA and FLD (Adapted from (Caprihan et al., 2008))	23
Figure 2 13 Performance Differentiation (Adapted from (F. Li et al., 2019)).....	24
Figure 2 14 Conceptual Framework	27
Figure 3 1 Position of Electrodes.....	29
Figure 3 2 Rapid Application Development	30
Figure 4 1 System Architecture	33
Figure 4 2 Use Case Diagram	35
Figure 4 3 Sequence Diagram.....	39
Figure 4 4 Entity Relationship Diagram	40
Figure 4 5 Wireframe for Log in Page.....	41
Figure 4 6 Wireframe for Uploading of Patient Data	42
Figure 5 1 Importing the data and viewing it.....	46
Figure 5 2 Checking data for missing values.....	46
Figure 5 3 Checking for zero values	47
Figure 5 4 Model with KNN Classifier.....	48
Figure 5 5 Model with Decision Tree Classifier.....	48

Figure 5 6 Model with Random Forest Classifier..... 49

List of Abbreviations

EEG – Electroencephalography

ERP – Event Related Potential

ANN – Artificial Neural Network

CNN – Convolutional Neural Network

ML – Machine Learning

SZ -Patients with Schizophrenia

HC – Healthy Controls

Acknowledgements

I would like to sincerely thank and acknowledge the guidance offered to me by my supervisor Dr. Joseph Orero throughout the thesis proposal and implementation

CHAPTER 1 : INTRODUCTION

1.1 Background of the study

Schizophrenia is a chronic mental condition that affects approximately 1% of the population and is characterized by a variety of symptoms. These symptoms include positive manifestations like hallucinations and delusions, as well as negative symptoms like emotional flatness (difficulty in speech and affected individuals having less words) and avolition (complete lack of motivation)(Whiteford et al., 2013).

The Diagnostic and Statistical Manual of Mental Disorders (DSM-5) serves as the primary tool for diagnosing schizophrenia. Patients are asked a series of questions to elicit information such as the length of disease and clinical symptoms (American Psychiatric Association, 2013).

However, diagnosing the condition can be challenging due to its coexisting or co-occurring conditions, potentially leading to ineffective patient management(Bae et al., 2021). In addition, these traditional clinical diagnoses may lack accuracy due to patients concealing symptoms or overlapping symptoms with other disorders.(Lindström et al., 1994)

Psychiatrists and clinical psychologists can diagnose schizophrenia with the aid of a variety of diagnostic tools; however, traditional clinical diagnoses may not always be accurate because patients with schizophrenia may purposefully hide their symptoms, and even professionals may occasionally find it difficult to distinguish schizophrenia from other mental illnesses due to symptoms that are similar.(Lindström et al., 1994) , With the use of neuroimaging technologies, numerous researchers have worked to create objective, quantitative biomarkers that can improve overall diagnosis accuracy. Because of its low cost and great temporal resolution, electroencephalography (EEG) is considered one of the most valuable neuroimaging modalities. Numerous researches document impaired brain information processing in schizophrenia, accompanied by modified event-related potential (ERP) waveforms.(Shim et al., 2016)

The study of statistical patterns of interdependencies between EEG (electroencephalogram), MEG (magnetoencephalography), and, more recently, MRI has been conducted extensively to learn more about the brain's "functional connectivity." The real anatomical connections of the

underlying neural networks are a prerequisite for functional connectivity, but they are not the same. Regarding functional connectivity, one inquiry is whether it represents particular temporal and spatial characteristics associated with an ideal state for information processing(Stam, 2004).

Schizophrenia patients frequently misinterpret feelings and experiences, which may be a factor in the psychotic symptoms that define the disorder. These misunderstandings and misinterpretations may stem from a fundamental incapacity to forecast predicted feelings and experiences with any degree of accuracy. Normal, healthy individuals benefit from brain systems that enable them to anticipate events without conscious thought. This process helps them analyze experiences and discern between what is expected and what is unexpected(Ford et al., 2014). Sensations that ought to have been anticipated but weren't could have an improper significance if predictive processes are malfunctioning.

One hypothesis suggests that individuals with Schizophrenia struggle to distinguish between internal and external stimuli due to issues with the nervous systems corollary discharge process. This difficulty in predicting and interpreting sensations may contribute to psychotic symptoms. Normal individuals utilize predictive processes to anticipate events, a function that may be impaired in patients with Schizophrenia.(Ford et al., 2014)

In this study, we will assess the neural electrical activity of individuals both, those who have schizophrenia and healthy controls using machine learning. The study of machine learning in computer science evaluates the application of algorithms to carry out a given task without explicit human guidance(Mohri, M., Rostamizadeh, A., 2018). It has recently been widely used in many disciplines, such as neuroscience and biology. In this regard, a variety of machine learning algorithms have demonstrated tremendous usefulness, especially when it comes to classifying items by numerous attributes. These have recently been used in studies on schizophrenia and have demonstrated excellent accuracy in classifying the illness based on a variety of non-pathognomonic characteristics.

1.2 Problem Statement

The ability to point out schizophrenia in patients can sometimes pose several challenges due to the complexity of the disorder. Traditional methods of diagnosis often rely on assessment by clinicians or physicians, which may lack impartiality and reliability.(Oh et al., 2019)

The integration of event related potential (ERP) data attained from electroencephalography (EEG)with machine learning models can help detect schizophrenia more accurately. ERPs provide important insights into the neural activity linked to schizophrenia because they show the brains electrical activity in response to a particular stimulus.(Ford et al., 2014)

The development of machine learning algorithms that can efficiently integrate and evaluate ERP data from EEG is required to address these issues. These models can extract important information from ERP signals and distinguish between patients with schizophrenia and healthy controls with a considerable amount of accuracy. Furthermore, the detection of intricate patterns and interactions in EEG data is made possible by machine learning techniques, which improve the diagnostic accuracy of schizophrenia detection.(Oh et al., 2019)

Therefore, the primary objective of this study is to investigate the value of machine learning models utilizing ERP data from EEG in enhancing schizophrenia detection. Through the development and validation of a machine learning model, this research aims to identify unique neurophysiological signs of schizophrenia and come up with a credible criterion for diagnostic decision making

1.3 Objectives

1.3.1 General Objectives

The general objective would be to develop a machine learning model that detects schizophrenia in patients by use of electroencephalogram EEG pulses.

1.3.2 Specific Objectives

- i. To investigate the problems associated with schizophrenia diagnosis.
- ii. To review existing methods that have been put in place to detect Schizophrenia
- iii. Develop and train a model that detects Schizophrenia in individuals using EEG data
- iv. To validate the model

1.4 Research Questions

- i. What are the problems associated with schizophrenia diagnosis?
- ii. What are the methods currently in place to detect Schizophrenia?
- iii. How to develop a model that detects Schizophrenia in individuals using EEG data?
- iv. How will the functionality of the data be tested?

1.5 Justification

Schizophrenia is a critical illness of the brain that interferes with normal cognitive functions, verbal dialogue and the behavioral attributes of a person(Chatterjee et al., 2018), According to the National Institute of Mental Health, 2.4 million Americans over the age of 18 are affected by Schizophrenia, making it a major contributor to the disease burden.(Oh et al., 2019)

The disorder has a great impact on those who suffer from it, their quality of life is often compromised and most of them are unable to function properly in their designated workplaces, with 20-40% resorting to suicide attempts, whilst 5-10% succumbing to it,(Tibbets, 2013) that is why an accurate and timely prognosis is necessary for the patient's recovery for better treatment. At the moment there is no recognized clinical test for Schizophrenia, instead professional observation of behavioral signs is used for diagnosis. These evaluations lack objectivity and accuracy and are unable to identify underlying anomalies in the operation of the brain(Oh et al., 2019)

The use of ERP data from EEG and incorporating it into a machine learning model would enable us to extract the patterns and signals that are unique to schizophrenia patients, enabling us to pick out the identifiers of the disorder relatively early and aid in early detection.

1.6 Scope and Limitation

The study is constrained to using secondary data in order to develop a model for Schizophrenia detection in patients, the model will not be used as a replacement for diagnosis but as a tool for early detection

CHAPTER 2 : LITERATURE REVIEW

2.1 Introduction

This chapter will examine the factors that contribute towards Schizophrenia worldwide, which individuals are more likely to suffer from Schizophrenia and how it affects the livelihoods of those who suffer from it, we will also explore Electroencephalography (EEG) and Event Related Potential (ERP) and their significance in the recording of brain activity and neural responses to stimuli. A conceptual framework will also be shown in order to guide the forthcoming chapters.

2.2 Schizophrenia

In the past, mental and substance use disorders were not considered a priority in global health circles, especially when put up against communicable diseases and non-communicable diseases such as heart diseases or cancer.(Whiteford et al., 2013), but in recent times these mental diseases have begun receiving quite a significant amount of attention. Schizophrenia is one of these mental disorders. The ailment is a syndrome, characterized by a group of symptoms and signs with no apparent cause, primarily psychotic symptoms. Schizophrenia typically manifests as auditory hallucinations and paranoid delusions in late adolescence or early adulthood. Over the previous century, there hasn't been much of a shift in these symptoms of the illness.(Insel, 2010) A common symptom of schizophrenia that affects about 75% of patients is auditory hallucinations.(Nayani & David, 1996). These audible signals are experienced as voices despite the fact that there isn't a single person talking.

The precise root of Schizophrenia has not been completely comprehended, though it is hypothesized to come about from a series of neurobiological, environmental and genetic contributors. An in-depth look into it points to anomalies in the cranial make up and neurotransmitter function, specifically glutamate and dopamine may be a contributing factor in the progression of the disorder. When it comes to the treatment of Schizophrenia, antipsychotic drugs, psychological therapies and support therapies are used in conjunction with one another(Meltzer, 2004). Despite the difficulties that stem from Schizophrenia many people often lead happy productive lives when given the right care and assistance. For individuals impacted

by this intricate and incapacitating illness, early intervention and all-encompassing treatment are crucial to improving their quality of life.

2.3 Efference Copy /Corollary Discharge

This section highlights Efference copy and corollary discharge which are neural signals involved in movement control and perception, the two terms both point to the process of the brain passing motor commands to a number of sensory destinations.

Efference Copy is a signal that the brain's motor control sends to the sensory areas prior to a movement being performed. In essence, it is a duplicate of the motor commands that are to be used. The brain can predict the sensory feedback that the movement will produce all regards to the particular copy.(Pynn & DeSouza, 2013), when the efference copy is produced, it is passed to sensory areas where it is normally used to forecast the sensory consequences of the motion. The signal that enables this predictive capability is known as corollary discharge, it aids the brain tell the difference between sensory signals that are self-generated and those that are brought about by an external stimuli.(Ford & Mathalon, 2005)

When a motor action is followed by an efference duplicate of the action, the sensory cortex receives a "corollary discharge" signal indicating that the upcoming sensations are self-initiated or self-generated. This component in the visual system may help to preserve visuospatial consistency by stabilizing the visual image during eye movements. (Sperry, 1950)The efference copy/corollary discharge mechanism, in its most basic version, suppresses perception of events that follow a self-generated action. Thus, it may enable an automatic differentiation between internally and externally generated perceptions in addition to acting as a tool for learning and optimizing our activities.

A comparable operation may be in existence in the auditory system, corollary discharges emanating from motor speech making areas in the frontal lobes prepare the auditory cortex for understanding the voice produced as self-induced.

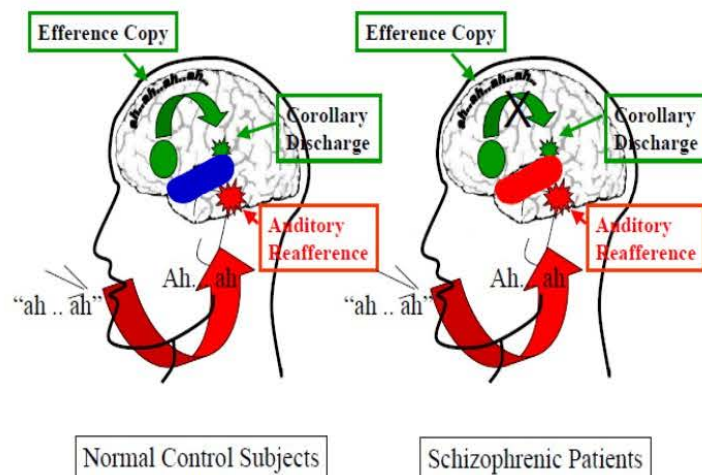


Figure 2 | Disconnect of the corollary discharge in schizophrenic patients (Ford & Mathalon, 2005)

Two schematics depicting the potential breakdown of the efference copy/corollary discharge mechanism in schizophrenia (right) and its normal functioning during speech (left). The speech plan is shown by a green circle next to Broca's area and begins in the frontal lobes. The idea or intended noises are sent as an efference copy (green ribbon) to the auditory cortex, where they are converted into a corollary discharge (green splash). Talking starts at the same moment, or maybe a few milliseconds later, and the speech sounds arrive at the auditory cortex as the auditory refferent (red splash) and form a red ribbon. Should the corollary discharge correspond with the auditory refferent, the impact of the sensory experience is either eliminated or diminished. Blue (left) in the auditory cortex represents typical repressed reactivity to the self-generated audio when it matches the corollary discharge. The diagram on the right shows that the Schizophrenic patient has a mark of X in the efference copy and the auditory cortex is marked in red to illustrate that the activity is not suppressed during speech. (Ford & Mathalon, 2005)

Failures in corollary discharge have been proposed as a possible cause of positive symptoms in schizophrenia (Feinberg & Guazzelli, 1999). Patients may experience auditory verbal hallucinations or passivity experiences if an efference copy of an intended action (or thought) does not result in a corollary discharge of the expected experience. This is because patients may be unable to distinguish between their own thoughts and voices generated by outside sources.

2.4 EEG (Electroencephalography) and ERP (Event Related Potentials)

One unique method of trying to capture auditory hallucinations in Schizophrenic patients is called “symptom capture”, this method aims to image the brain while patients are experiencing hallucinations using Electroencephalography (EEG), functional Magnetic Resonance Imaging (fMRI) or Positron emission tomography (PET). Although these methods appear to be straightforward conceptually, they are very challenging to implement in practice because it depends not only on the timely occurrence of an elusive subject experience but also on the patients consistent reporting of the event’s beginning and end (Ford & Mathalon, 2005). Quite a number of researchers have endeavored to come up with objective signs that can improve the general accuracy of analysis and detection by the help of neuroimaging technological methodologies. In the wide scope of the different modes that are considered in neuroimaging technologies, EEG is looked upon as the most advantageous, this is mostly attributed to its ability to capture and represent fast changes in the brain activity over time. (Shim et al., 2016)

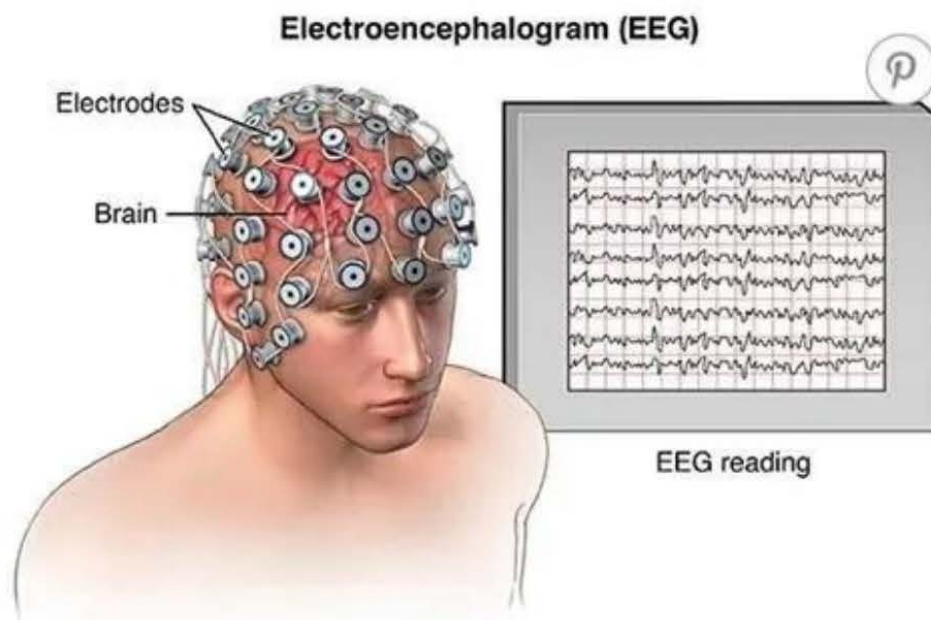


Figure 2 2 Illustration of EEG Data Recording (Adapted from Brightbraincenter.co.uk)

An EEG (electroencephalogram), is a non-invasive method used to capture electrical activity in the brain. Let us explore how this works

Placement of Electrodes: the process starts with the placement of small metal electrodes on the scalp of the individual, the electrodes are placed in particular locations according to an

international standard known as the 10-20 system(Herwig et al., 2003), which ensures consistent positioning. The electrodes are organized in a grid like pattern to cover a number of regions in the scalp.

Signal Acquisition: Once the electrodes have been positioned, they detect electrical signals generated by neurons in the brain. Neurons communicate with one another using electrical impulses(Barros et al., 2021), an external stimuli can be introduced to the individual in order to trigger electrical activity in the brain, the electrodes are connected to an amplifier, which strengthens the weak electrical signals that are picked up by the electrodes.

Processing of EEG Signals: The amplified signals are then processed by a computer, the process involves filtering of the waveforms to remove noise and other variables, Filtering help separate the specific bands of interest, which are associated with different brain states and activities.

After the processing of EEG signals is complete, they are recorded as a series of waveforms named as EEG trace, the electrical activity of the brain throughout time is represented by these waveforms, different patterns of brain activity can be seen. EEG recordings are often interpreted by professionals, such as neurologists. These electrical activity patterns can be analyzed to identify abnormalities such as sleep disorders, epileptic seizures or neurological conditions.

2.5 ERP -Event related Potential

ERP or Event Related Potentials refers to electrical measurement of brain activity which is obtained from the Electroencephalogram (EEG), they depict the electrical activity in the brain when it is reacting to particular mental or sensory inputs(Light et al., 2010). ERP's are produced by subjecting a participant to a series of stimuli, the stimuli can be, audio, visual or related to touch.

When a stimulus is introduced, it activates some neural activity in the brain associated to the processing of that stimulus, this activity leads to the production of electric potentials which in turn can be recorded by the electrodes positioned on the scalp. ERP's are often described by particular mechanisms, or peaks in the EEG waveforms that parallels to various steps in the processing of information.(Baess et al., 2011)

In terms of differentiating between ERP's (Event Related Potentials) and EEG, ERP is the neurological reply of the brain to an event be it internal or caused by the external environment, ERP values are deduced from the captured EEG signals. The activity and reactions of different

brain processors to the ongoing task are summed up in EEG signals, ERP signals on the other hand, show how the brain reacts to a particular stimulus. (Golnaz Baghdadi, Farzad Towhidkhan, 2021)

2.6 Machine Learning

Machine learning can be explained as the area of study that enables computers to learn without overt programming, by using Machine Learning, machines can be trained to handle data more efficiently. Usually, after viewing the data and using our own mental abilities we are unable to decipher the information that has been extracted, In such cases we then use machine learning methods that are able to read further into the data.(Mahesh, 2020). The need for machine learning is growing due to the number of the data sets that are available. Machine learning is used by many sectors of the economy to retrieve pertinent data. Learning from a given dataset is the primary aim of machine learning and numerous works of research have been conducted on the subject of teaching computers to learn on their own without explicit programming Nowadays, a number of AI system engineers understand that, depending on the application, it can sometimes be far simpler to train a system by providing it with instances of desired input-output behavior than to manually build it by predicting the appropriate response for every conceivable input(M. I. Jordan & T. M. Mitchell, 2015). Machine learning has also had a significant impact on computer science and a variety of other industries that deal with data-intensive problems, like consumer services, fault diagnosis in intricate systems, and the management of supply chains.

Supervised learning involves using sample input-output pairs to train a function that maps an input to an output. From labelled training data, which consists of a collection of training instances, it infers a function. The machine learning algorithms that require external aid are known as supervised algorithms(Mahesh, 2020). The train and test datasets are separated from the input dataset. An output variable from the train dataset needs to be categorized or forecasted. Every algorithm picks up patterns from the training dataset and uses them to make predictions or classify data from the test dataset. Direct feedback is used by supervised learning systems to make predictions. Regression and classification techniques are two categories under which supervised learning falls. Several supervised learning algorithms are widely used, including KNN, DT, SVM, LR, Artificial Neural Network (ANN), and Naïve Bayes (NB). The primary

goal of classification techniques is to accurately identify and forecast the likelihood of identifying the target value in the data.(Chauhan et al., 2021)

Unsupervised Machine learning on the other hand does not take response for the prediction, this kind of algorithm learns by discovering patterns in the data that are not so obvious, the concepts that are dominantly applied in unsupervised machine learning are clustering analysis and component analysis(Chauhan et al., 2021). Some of the common methods of unsupervised machine learning techniques are K Means clustering and Principal Component Analysis.

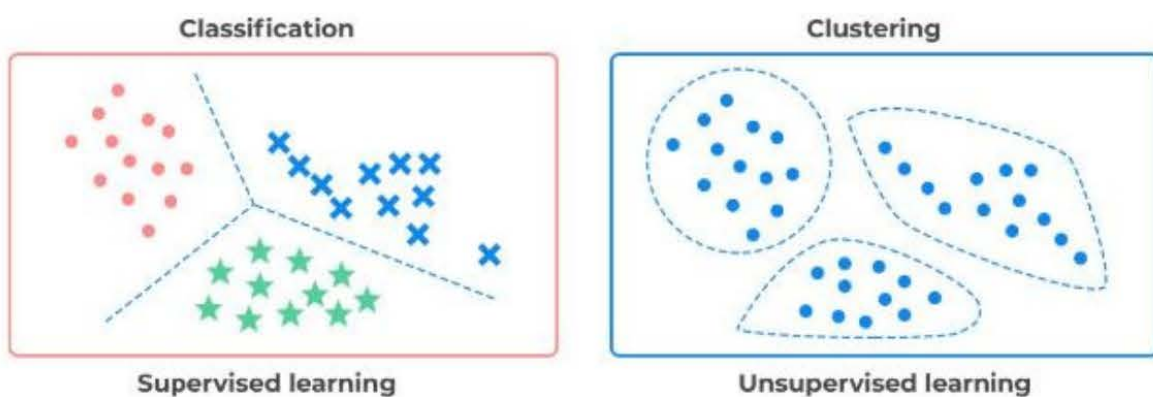


Figure 2 3 Difference between Supervised and Unsupervised Learning (Adapted from IBM)

2.6.1 Machine Learning Techniques

Machine learning methods can include quite a number of algorithms that are designed to process and interpret various datasets, some of these methods are K Nearest Neighbor, decision trees, random forest, and support vector machines (SVM). These methods offer varying approaches in solving problems, each of these classifiers has its own advantages and particular application. We are going to explore each of these techniques

2.6.1.1 K Nearest Neighbor

One of the Supervised learning machine learning models is the K-Nearest Neighbors classifier, the KNN algorithm uses proximity to classify or predict how a single data point will be grouped,

it is among the most widely used and straightforward regression and classification classifiers in machine learning.

K-nearest neighbors (KNN), also known as nearest neighbor classification is underpinned on the concept that the closest patterns to a target pattern x for which we see. KNN allocates a class label to the most of the K-Nearest forms in the data space.(Kramer, 2013)

It is important to note that the KNN method belongs to a family of models known as “lazy learning” which means that instead of the model going through a training phase, it merely saves a training dataset. This implies that whenever a classification or prediction is being produced the computing takes place. It is also called an instance-based or memory based learning approach because it mostly depends on memory to retain all its training data.(IBM, n.d.)

KNN mostly uses the Euclidean distance to allocate its distance but there are other distances to be considered such as the Manhattan Distance and the Minkowski Distance

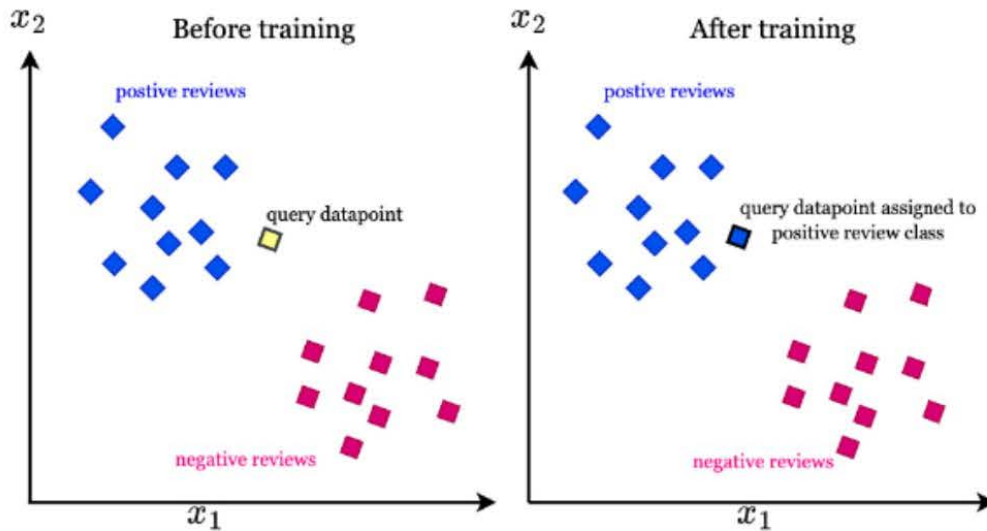


Figure 2.4 K Nearest Neighbor Classification

2.6.1.2 Decision Trees Classifier

Decision Trees is among the most common supervised machine learning classifier, using a greedy search to find the best split points inside a tree, decision tree learning uses a divide and conquer tactic. Once all, or most, of the records have been assigned a class label, this dividing procedure is then repeated top-down and recursively. The decision tree's complexity plays a major role in classifying all the data points as homogenous sets or not. Pure leaf nodes, or data

points in a single class, are easier to obtain for smaller trees. But maintaining this purity gets harder as a tree gets bigger, and this usually means that too little data falls inside a particular subtree. Data fragmentation is the term used to describe this situation.

Decision Trees are one of the most potent techniques that are used in a number of fields such as identification of patterns and image processing, Decision Trees is a sequential model that effectively and cohesively combines a number of fundamental tests in which a numerical feature is compared to a threshold parameter in each test.(Charbuty & Abdulazeez, 2021)

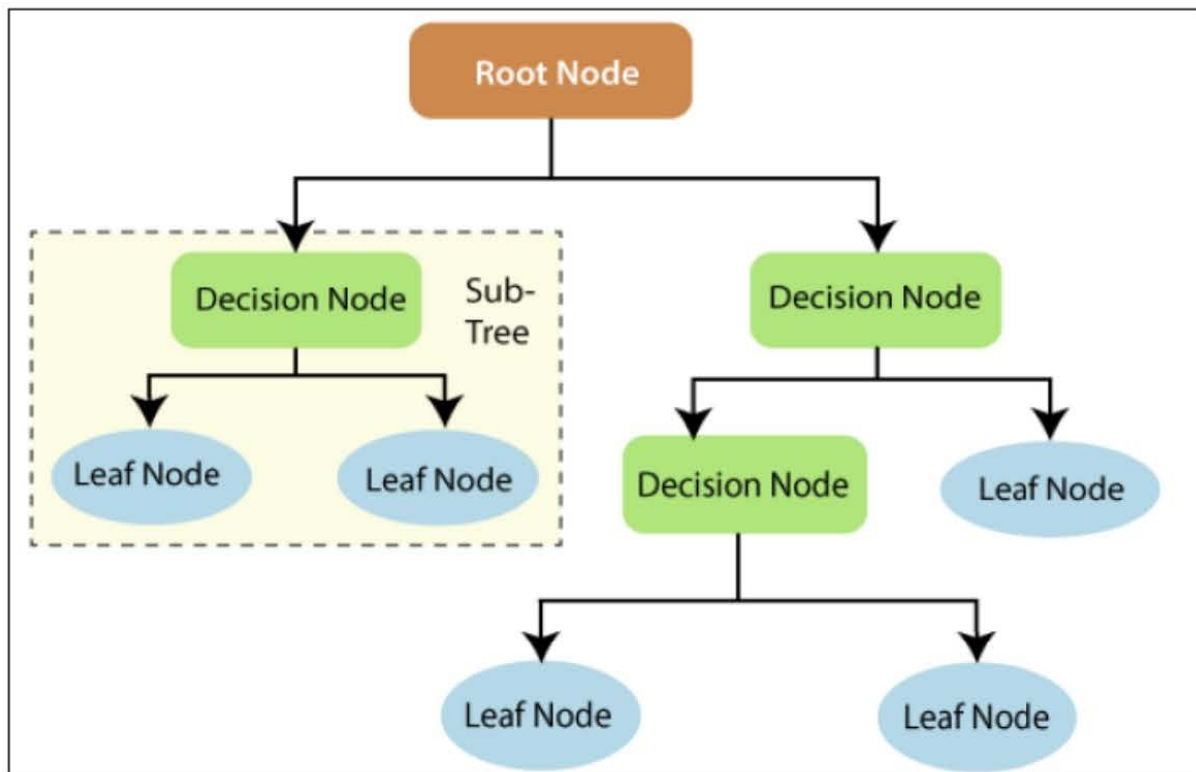


Figure 2 5 Decision Tree Classifier

In regard to picking a node for the decision tree, there are several techniques to choose the best feature at each node, however information gain and Gini impurity are the two approaches that are mostly used as the splitting criterion for decision tree models. They aid in the assessment of each test condition quality and its capacity to categorize samples into classes.(Charbuty & Abdulazeez, 2021)

In order to explain information, gain properly we must highlight entropy, Information theory gave rise to the idea of entropy, which quantifies the degree of impurity in sample values. The formula that follows defines it

$$\text{Entropy (S)} = -\sum_{i=1}^n P_i \log_2 P_i$$

- S represents the data set that entropy is calculated
- i represents the classes in set, S
- p(i) represents the proportion of data points that belong to class i to the number of total data points in set, S

Values of entropy can range from 0 to 1. When every sample in the data set S is a member of the same class, entropy is equal to zero. Entropy will peak at 1 if half of the samples are categorized into one class and the other half into a different class. The characteristic with the least level of entropy should be utilized to determine which feature is best to divide on and to identify the optimum decision tree. The difference in entropy before and after a split on a particular attribute is known as information gain. Since it is performing the greatest job of classifying the training data in accordance with its target classification, the attribute with the maximum information gain will result in the best split.

$$\text{Gain(S ,A)} = \sum_{v \in V(A)} \frac{|S_v|}{|S|} \text{Entropy(S_v)}$$

- A represents a specific attribute or class label
- *Entropy (S_v)* is the entropy of dataset, S_v
- $|S_v|/|S|$ represents the proportion of the values in S_v to the number of values in dataset, S
- *Entropy(S_v)* is the entropy of dataset, S_v

Gini Impurity

A dataset's Gini impurity is the likelihood that a random data point would be wrongly classified if its labeling were dependent on the dataset's class distribution. Comparable to entropy, the impurity of a set, S , is zero if it is pure, that is, belongs to a single class (Grabmeier & Lambe, 2007). The formula that follows indicates this:

$$Gini = 1 - \sum_{i=1}^C (p_i)^2$$

Comparing decision tree classifiers to other classification techniques, they achieve comparable and occasionally even higher accuracy. The decision tree algorithm can be applied in either a serial or parallel form, depending on the amount of data, the computing resource's memory capacity, and the algorithm's scalability. (Anyanwu & Shiva, 2009)

2.6.1.3 Random Forest Classifier

Random forest being another supervised learning classifier aggregates the output of several decision trees to get a single outcome. Its versatility and ease of use, combined with its ability to handle both regression and classification problems, have driven its popularity. Despite being a popular supervised learning algorithm, decision trees have many drawbacks, including bias and overfitting. In the random forest algorithm, however, when several decision trees come together to form an ensemble, they forecast outcomes that are more accurate, especially when the individual trees have no correlation with one another (Belgiu & Drăgu, 2016).

The three primary hyperparameters of random forest algorithms must be set prior to training. These consist of the size of the nodes, the count of trees, and the quantity of characteristics sampled. Regression and classification issues can then be resolved using the random forest classifier. Each decision tree in the ensemble of decision trees used in the random forest technique is made up of a bootstrap sample, which is a sample of data taken from a training set with replacement. One-third of the training sample is designated as test data; this is referred to as the out-of-bag sample. Feature bagging is then used to introduce yet another randomization, increasing dataset variety and decreasing decision tree correlation. The prediction's determination will change depending on the kind of problem. The individual decision trees in a regression job will be averaged, and in a classification work, the predicted class will be determined by a majority vote, or the most common categorical variable. (Petkovic et al., 2018)

Random Forest Classifier

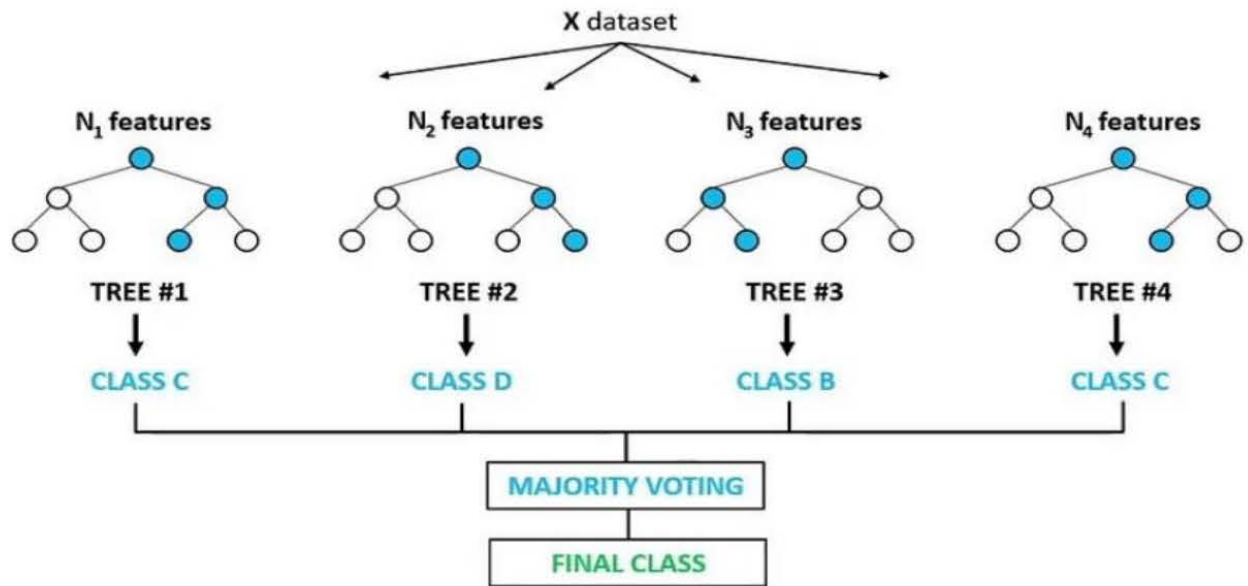


Figure 2.6 Random Forest Classifier

2.6.1.4 Convolutional Neural Networks

Convolutional Neural Networks borrow heavily from Artificial Neural Networks, so before we delve into what CNN is, let us have a look at what ANN's are. Artificial Neural Networks (ANNs) are computational processing systems that draw significant inspiration from the functioning of organic nerve systems, including the human brain. The basic component of artificial neural networks (ANNs) is a large number of networked computing nodes, also called neurons. These neurons collaborate in a dispersed manner to learn from the input and optimize the output. (Oh et al., 2019)

An ANN's fundamental structure can be modeled as follows. The input layer receives the input, which is typically loaded as a multidimensional vector and then distributed to the hidden layers. The process of learning is carried out by the hidden layers, which use the decisions made by the previous layer to determine whether a stochastic change will ultimately improve or worsen the output (Oh et al., 2019). Deep learning is the term used to describe systems with several hidden layers layered on top of each other.

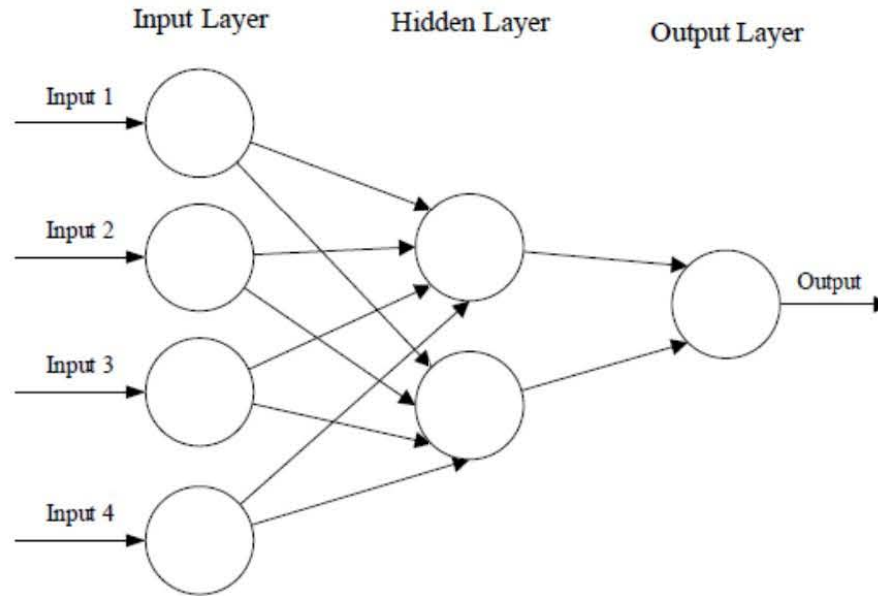


Figure 2 7 Artificial Neural Networks

The figure above shows a general structure of an Artificial Neural Network, it is comprised of a simple three-layered feedforward neural network (FNN), the layers are an input layer a hidden layer and an output layer. The structure is also common in Recurrent Neural Networks. Similar to conventional ANNs, convolutional neural networks (CNNs) are made up of neurons that learn to optimize on their own. ANNs are based on the fact that every neuron will continue to receive an input and carry out an operation. (Z. Li et al., 2022) The entire network will still express a single perceptive score function (the weight), starting with the raw picture vectors in the input and ending with the class score as the final output. The loss functions linked to the classes will be found in the final layer, and all of the standard techniques created for conventional ANNs still holds true.

One of the more significant differences between CNNs and orthodox ANN is that CNNs are majorly utilized in the area of pattern identification in images. The fact that conventional ANNs often have trouble with the computational complexity needed to compute image data is one of their biggest drawbacks.

CNN Architecture

Convolutional Neural Networks are made up of 3 distinct layers, the layers are convolutional layers, pooling layers and fully connected layers, the basic functionality of the layers are

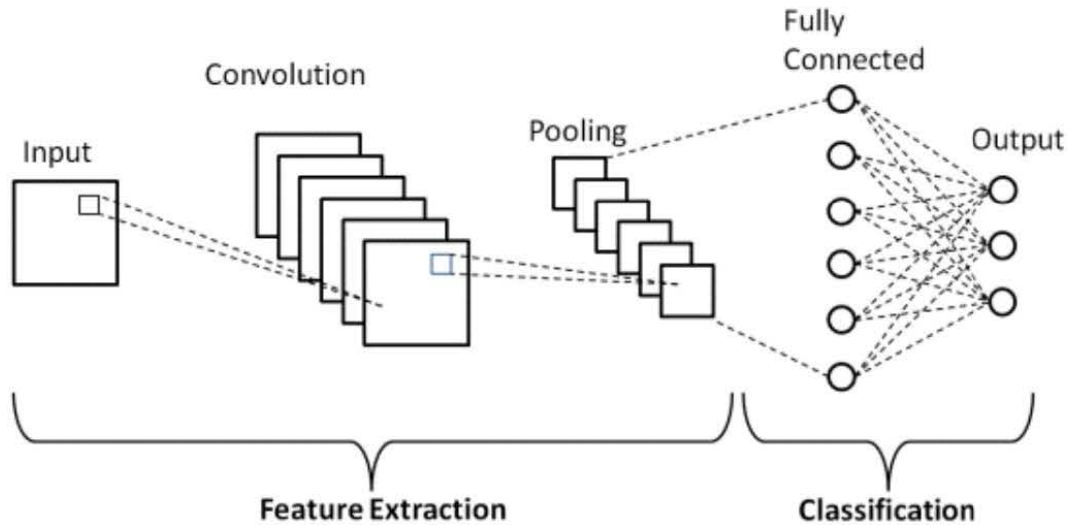


Figure 2 8 Convolutional Neural Networks

- i. The input layer will hold the values
- ii. By computing the scalar product between the weights of the neurons and the region linked to the input volume, the convolutional layer will ascertain the output of neurons that are associated with specific local parts of the input.
- iii. Afterwards, the pooling layer will merely carry out down sampling along the input's spatial dimensions, so further decreasing the number of parameters in that activation.
- iv. Next, the fully-connected layers will carry out the tasks present in conventional ANNs and try to generate class scores from the activations that will be utilized for classification.

2.6.2 Integration of EEG and Machine Learning

The Diagnosis of Schizophrenia is done mainly by the use of the Diagnostic and Statistical Manual of Mental Disorders (DSM-5), this is mostly done by taking patients through a number of questions, this helps the practitioner to gain information about the duration of the ailment and other symptoms(Shim et al., 2016).

Recently, a growing number of studies have tried to use machine learning (ML) techniques with EEG biomarkers to distinguish between patients with schizophrenia and healthy controls. Sensor-level biomarkers, such ERP amplitude and latency, were employed as features for categorization in some of these studies.(Shim et al., 2016)

2.7 Using EEG in Schizophrenia Detection

This particular segment will analyze existing contexts and systems, weighing the various techniques that have been used in Schizophrenia detection whilst using a number of technological models.

(Ford et al., 2014) conducted a study to determine whether patients with schizophrenia show less evidence of suppression in motor response as a result of a self-generated stimuli as compared to healthy controls. The study was carried out on 26 schizophrenic patients and 22 healthy controls, each subject carried out 3 conditions. The conditions were (1) press a button and automatically hear a tone (2) passively hear the same tone, and (3) press a button without a tone.

J. M. Ford et al

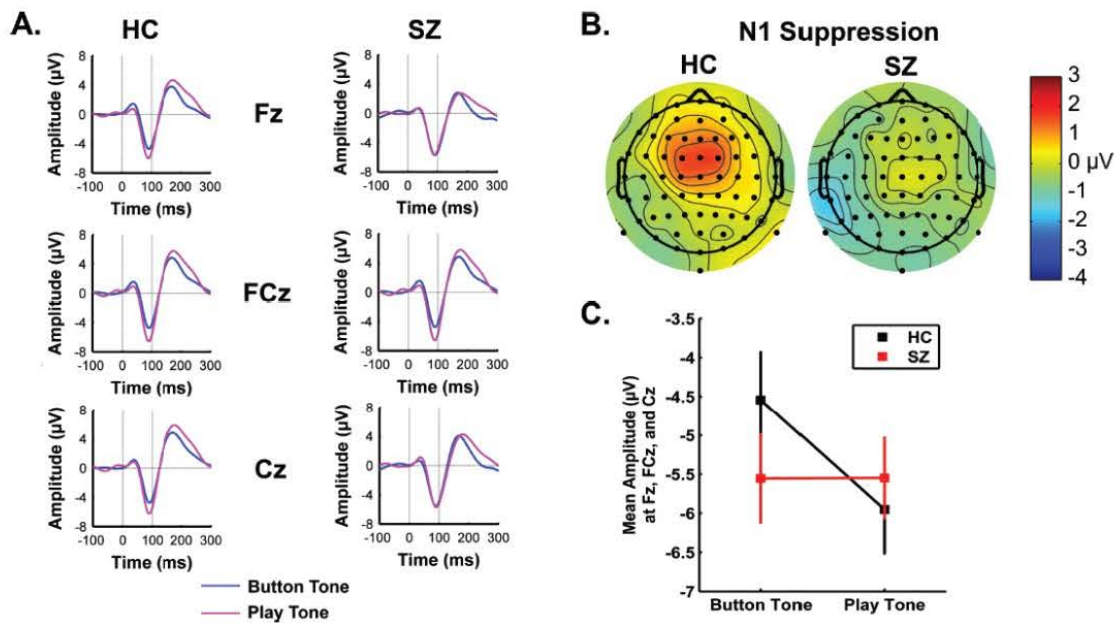


Figure 2.9 Illustration of N1 Suppression in Schizophrenic Patients (Adapted from (Ford et al., 2014))

Event related potentials to tone onset are overlaid for Button Tone and Play Tone for Healthy controls (HC) and Schizophrenia (SZ), B Shows Topography maps for suppression of N1 amplitude C. N1 amplitudes averaged during Button Tone and Play Tone for HC (black) and SZ(red.)

The findings were enough to validate the hypothesis that when a button is pressed to release a tone, patients with schizophrenia exhibit less suppression than healthy controls and less activity beforehand. The statistics indicated that this was as a result of failures of the efference copy and corollary discharge. The conclusion was that individuals with schizophrenia exhibit deficiencies in both the efference copy and corollary discharge even when the efference copy is tangentially related to the accompanying sound and does not originate from self. This was also evident in the (ERP) Event Related Potentials.(Ford et al., 2014)

2.8 Related works

2.8.1 Schizophrenia Detection Using Machine Learning Approach from Social Media Content

By examining social media messages, the study attempted to ascertain whether machine learning can be utilized to identify symptoms of schizophrenia in users. For the control group, posts were gathered from the social media site Reddit in addition to posts about fitness, parenting, relationships and teaching that had nothing to do with mental health. From the posts language characteristics and content subjects were retrieved. Linguistic markers of schizophrenia were identified by using supervised machine learning to classify posts as belonging to schizophrenia and interpreting significant aspects. The features were subjected to unsupervised clustering in order to identify a consistent semantic representation of words in schizophrenia. The accuracy rating in differentiating schizophrenia postings from controls was 96% (Bae et al., 2021)

Among the classifier algorithms, the classifier with the highest level of accuracy was the random forest algorithm with an accuracy of 96%, which goes above all the other accuracy results, the table below shows the classification performance.

Table 4. Classification performance of the machine learning classifiers.

Model	Recall	Precision	F1-Score	Accuracy	AUC
Random forest (RF)	0.94	0.98	0.96	0.96	0.97
Support vector machine (SVM)	0.91	0.90	0.91	0.91	0.91
Logistic regression (LR)	0.87	0.91	0.89	0.89	0.90
Naive Bayes (NB)	0.87	0.82	0.93	0.86	0.87

Figure 2 10 Classification Performance of Machine Learning Models (Adapted from (Bae et al., 2021))

One limitation that was observed in this study was that, there was no evidence showing that the users of the social media group (Redditors /Schizophrenia) were clinically diagnosed with schizophrenia. The other limitation was the lack of use of a Neural Network as a classifier given the remarkable accuracy levels produced by these models (ANN, CNN or RNN) when it comes to text classification, especially semantic classification.

2.8.2 Diagnosing Schizophrenia with Network Analysis and a Machine learning method

The study sought out to find out if there is an abnormality in brain connectivity as a pathophysiology of schizophrenia patients. Network analysis having become recently popular in the field of Schizophrenia detection, it was used in collaboration with machine learning in order to diagnose Schizophrenia.(Jo et al., 2020)

48 cases of Schizophrenia and 24 healthy controls were used in this particular study, using probabilistic brain tractography, graphs were reconstructed to estimate various global and nodal parameters. Following a comparison of these network properties across groups, machine learning was employed to categorize schizophrenia patients from healthy controls.

Support Vector Machine (SVM), random forest, Naïve Bayes and gradient boosting machine learning models demonstrated a considerable degree of performance in categorizing schizophrenia patients and healthy controls via network features. The model results and accuracies are listed in the table below.

TABLE 2 Performance of the machine learning models: Global

Model	Accuracy (%)	AUC
SVM	58.2 [17.5]	0.631 [0.202]
Multinomial NB	66.9 [4.0]	0.638 [0.233]
RF	68.6 [16.0]	0.680 [0.229]
XGBoost	66.3 [14.5]	0.633 [0.232]

Figure 2 11 Performance of ML Models (Adapted from (Jo et al., 2020))

The study above had some limitations, for instance the network properties that were picked for the machine learning models were intercorrelated, the characteristics selected to represent each network are not independent of each other, some are derived from one another and some are related, this affects the performance of the model during feature selection. Another limitation of this study is the use of a relatively small number of participants, this makes it hard to correctly evaluate the model and may lead to overfitting.

2.8.3 Application of principal component analysis to separate Schizophrenic Patients from Healthy Controls using fractional anisotropy images.

A Study was carried out by (Caprihan et al., 2008) to determine whether the use of Discriminant Principal Component Analysis can be utilized to obtain better results in distinguishing between patients with Schizophrenia and healthy controls. After observing that the traditional PCA is unable to produce the best results, DPCA was introduced. DPCA was applied to images generated from diffusion tensor data to differentiate between the target classes, the main goal was to identify specific brain regions that are associated with the condition. How effective the model was, was measured by a technique called “one-leave-out”, which tests the model performance by leaving out one data point and testing on it.

Data was obtained from 45 patients with Schizophrenia by use of a standard Siemens Scanner The results from the study were that DPCA produced larger separations between the two groups compared to using traditional PCA based on the given data.

In the dataset that was used, they were able to reduce the image dimension to 60 whilst maintaining performance accuracy similar to when the full dataset was used, this illustrated that the data could be significantly simplified without losing accuracy.

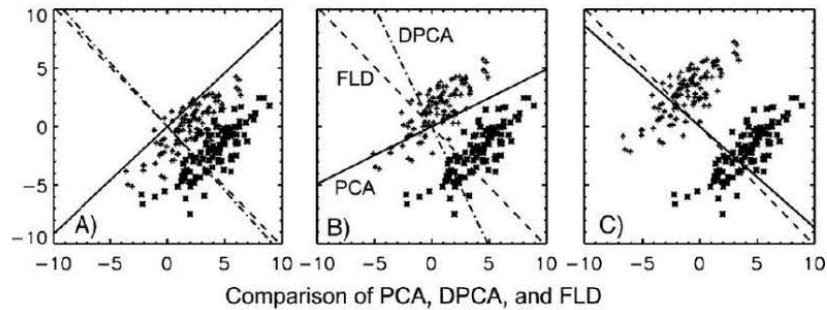


Figure 2.12 Comparison of PCA, DPCA and FLD (Adapted from (Caprihan et al., 2008))

2.8.4 Differentiation of Schizophrenia by Combining the Spatial EEG Brain Network Patterns of Rest and Task P300

The study focuses on improving the classification of patients with Schizophrenia (SZ) from healthy controls (HZ) by the use of EEG data focusing on resting state and P300 related brain activity. The approach of the study was to extract spatial pattern of network attributes for both brain states, this was done to enhance discrimination between SZ and HC.(F. Li et al., 2019)

According to the study's findings, the combined SPN features have the best accuracy of 90.48% when it comes to correctly classifying SZs from HCs, with a sensitivity of 89.47% for SZ identification and a specificity of 91.30% for HC identification.

The outcomes imply that useful information in differentiating between Healthy Controls and Schizophrenic patients can be found in both tasks related and EEG data, they hold good promise as clinical diagnostic tools for schizophrenia

Reliable distinction between SZs and HCs was achieved by combining SPN characteristics from multiple brain states (task and rest. This suggests that considering brain activity related to tasks as well as resting states improves the ability to discriminate between the two groups.

DIFFERENTIATING PERFORMANCE OF MULTIPLE MEASUREMENTS UNDER DIFFERENT BRAIN STATES

Brain state	Features	LDA			SVM		
		ACC(%)	SEN(%)	SPE(%)	ACC(%)	SEN(%)	SPE(%)
Task	Amplitude	52.38	57.89	47.83	76.19	68.42	82.61
	Network property	42.86	47.37	39.13	54.76	0.00	100.00
	SPN: One pair	61.90	68.42	56.52	73.81	57.89	86.96
	SPN: Two pairs	85.71	89.47	82.61	88.10	84.21	91.30
	SPN: Three pairs	83.33	89.47	78.26	85.71	84.21	86.96
Rest	Network property	71.43	63.16	78.26	54.76	0.00	100.00
	SPN: One pair	73.81	63.16	82.61	76.19	68.42	82.61
	SPN: Two pairs	83.33	73.68	91.30	88.10	78.95	95.65
	SPN: Three pairs	80.95	73.68	86.96	83.33	78.95	86.96
	Network property	64.29	57.89	69.57	54.76	0.00	100.00
States fusion	SPN: One pair	76.19	73.68	78.26	80.95	73.68	86.96
	SPN: Two pairs	85.71	78.95	91.30	90.48	89.47	91.30
	SPN: Three pairs	85.71	89.47	82.61	90.48	89.47	91.30

Figure 2 13 Performance Differentiation (Adapted from (F. Li et al., 2019))

STUDY	SCOPE	TECHNIQUES	RESULTS	HIGHLIGHTS /GAPS
Schizophrenia Detection Using Machine Learning Approach from Social Media Content (Bae et al., 2021)	Detection of Schizophrenia from text messages on the social media platform Reddit/Schizophrenia	Text Classification RF- Random Forest SVM -Random Forest LR-Logistic Regression NB-Naïve Bayes	RF- 96% SVM -91% LR-89% NB-86%	Users on the platform were not clinically diagnosed as Schizophrenia patients
Diagnosing Schizophrenia with Network Analysis and a Machine learning method (Jo et al., 2020)	Diagnosing of Schizophrenia patients through brain neural connectivity	SVM -Random Forest NB-Naïve Bayes RF -Random Forest XG Boost	SVM -58% NB-66.9% RF -68% XG Boost - 66.3%	Some of the network properties were intercorrelated, affecting model performance. The use of small no. of participants
Application of principal component analysis to separate Schizophrenic Patients from Healthy Controls using fractional anisotropy images. (Caprihan 2008)	Use of Diffusion Tensor Imaging to distinguish SZ patients from HC, by identifying specific brain regions associated with SZ	Discriminant Principal Component Analysis and Fishers Linear Discriminant (FLD)		DPCA was a better technique compared to FLD and traditional PCA
Differentiation of Schizophrenia by Combining the Spatial EEG Brain Network Patterns of Rest and Task P300	Distinguishing between patients with SZ and HC using SPN for brain networks	Spatial Pattern of Network	SPN -90.48%	Combination of both SPN features of rest and task resulted in a higher accuracy.

2.9 Research Gap

With the diagnosis or detection of Schizophrenia being limited to interviews and subjective behavioral observations by psychiatrists, a gap exists in leveraging the use of machine learning based EEG /ERP analysis in the detection of Schizophrenia. This study seeks to address this gap by the use of a trained machine learning model that can detect Schizophrenia in patients by the use of Event Related Potentials through Electroencephalography data.

2.10 Conceptual Framework

The following illustrates the conceptual framework

The first step involves obtaining ERP data from clinically diagnosed Schizophrenic patients and Healthy controls, the data then undergoes cleaning, to eliminate unnecessary data and impute missing values. Once the data has been processed, the data is divided into training and testing sets. The training of the data is fed into the machine learning algorithm, resulting in a prediction model. The model is then examined for its performance and thereafter positioned for use. Tests for the model can be done by pulling ERP data for an individual and testing it against the model.

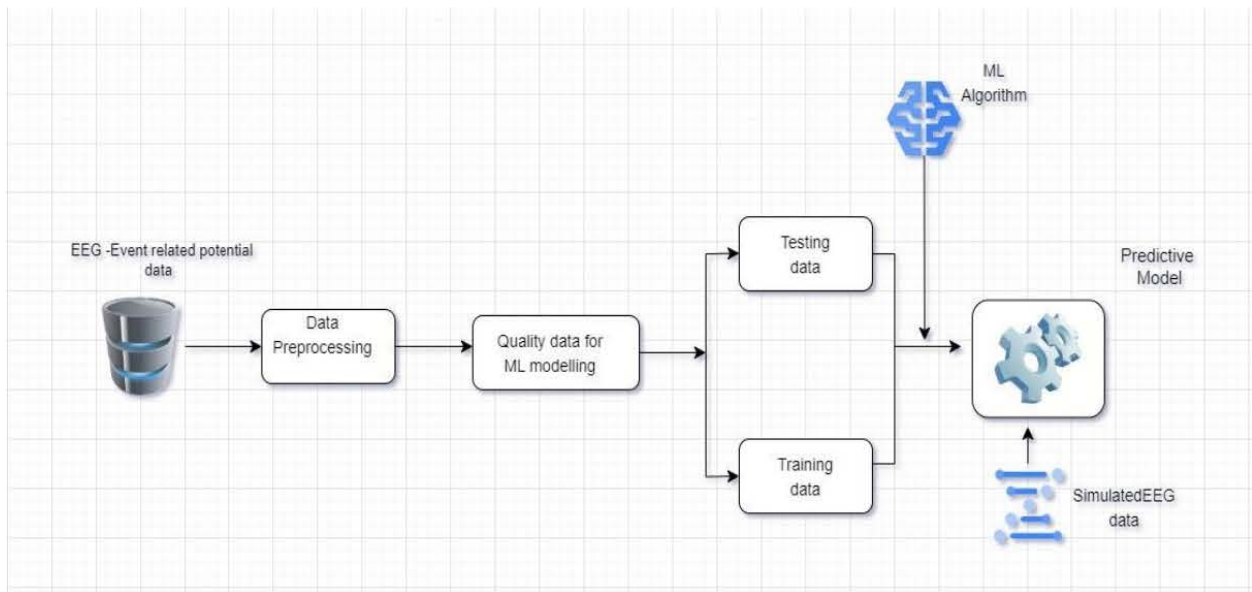


Figure 2 14 Conceptual Framework

CHAPTER 3 : RESEARCH METHODOLOGY

3.1 Introduction

This chapter outlines a methodical inquiry that was conducted to examine current issues and lay the groundwork for the best possible resolution. In order to draw findings and make defensible decisions, the study entails problem definition and revision as well as data collecting, organization, and evaluation. It covered every strategy, process, and approach used in conducting research. The use of Event Related Potential Data that has been classified for both patients with Schizophrenia and healthy patients will be used for training the model. Several tests will be created to verify the model and choose the best course of action.

3.2 Research Design

Research Design is a planned layout that has been specifically been formed to explore the questions that are the center of the research, which are the main research questions. Being able to give answers to the research question and exploring the probability of the hypothesis is the core function of the research(Dulock, 1993)

This can be considered an applied research study as it is required to help physicians in the early detection of Schizophrenia by analyzing EEG signals from brain activity in response to an external stimulus.

3.3 Datasets and Data Collection

This study's EEG dataset, which included 81 human subjects—32 healthy controls and 49 schizophrenia patients—was sourced from an online platform(Roach, n.d.). The collection includes event-related potential (ERP) EEG signals that were obtained during simple sensory activities that required pressing buttons or toning an audio system. Although the ERP averages for nine chosen electrodes (Fz, FCz, Cz, FC3, FC4, C3, C4, CP3, CP4) are given, the raw data were gathered from 64 EEG channels. The nine channels' topological placements are shown in the illustration below.

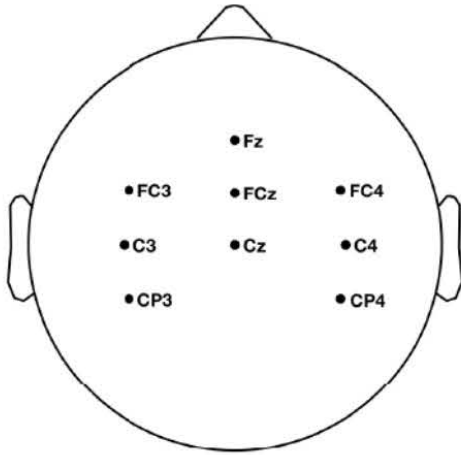


Figure 3 1 Position of Electrodes

The EEG data were digitally bandpass filtered between 0.5 Hz and 15 Hz after being captured at 1024 Hz and referred off-line to averaged earlobe electrodes. To examine the event-related EEG dynamics, data epochs are taken out of the continuous EEG signals. A data epoch in ERP data is time-locked to the start of an event or stimulus. The EEG data in this dataset were baseline adjusted at -0.6 to 0.5 seconds and divided into 3-second epochs starting 1.5 seconds before and ending 1.5 seconds after the button was pressed.(Roach, n.d.)

Three distinct conditions were applied to the subjects as they completed sensory tasks:

- 1) Generate an audio tone by pressing a button;
- 2) Listen to the same tone passively;
- 3) Press a button without producing a tone.

The final dataset obtained is a CSV file that contains, the 81 subjects, each one of them subjected to the 3 different conditions and the values recorded in microvolts for the 9 electrode positions.

3.4 Data Preprocessing

The data that was received had to go through checks for any anomalies that may interfere with the performance of the model, this to make sure that the model can learn from high-quality and reliable data. Checks that were made on the data were identifying missing values and removal of

any duplicate records. An extra column was added to the dataset to indicate a diagnosis of whether the individual had schizophrenia or was a healthy patient. This would serve as a target class for the classifier.

3.5 System Development Methodology

Rapid Application Development (RAD) system development methodology is used to create the prototype that this research suggests. Software is developed utilizing the Rapid Application Development (RAD) lifecycle, which yields higher-quality software more quickly than using the standard software development lifecycle (Beynon-Davies et al., 1999). Rapid Application Development (RAD) lets companies create software more quickly while lowering development costs and preserving software quality.

Rapid Application Development (RAD)

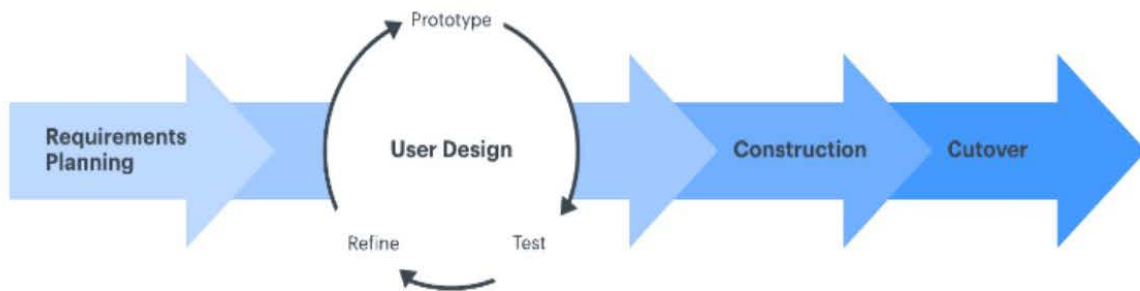


Figure 3 2 Rapid Application Development

The various segments of the Rapid Application development are shown as follows, In Regards to Requirements Planning, the stages are , definition of the problem, understanding of methodologies that have been put in place to solve these problems and Identifying gaps in existing solutions.

User Design refers to the constant improvement of the prototype in terms of usability and performance. Construction phase of RAD involves the actual development of the product, in this phase the system is implemented based on the requirements gathered during RAD planning and requirements gathering phases. The cutover phase comes in after the construction phase is completed, this mainly involves deploying the product into the production environment.it is usually critical to ensure a smooth transition from development to production.

3.6 Research Quality /Validation & Reliability

The quality of the research was based on the accuracy of the model in classifying ERP data. The testing set of the data possessed labeled ERP data that was utilized to appraise the model and give results for the confusion matrix

Model Accuracy was obtained using the following formula

$$accuracy = \frac{TP + TN}{TP + FN + TN + FP}$$

The precision of the model was obtained using the following formula

$$precision = \frac{TP}{TP + FP}$$

3.7 Model Validation

In order to validate the model, the error rate was calculated using the Mean Squared Error (MSE), Root Mean Squared Error (RMSE), and Mean Absolute Error (MAE). The average squared difference between the expected and actual numbers is known as the mean square error, or MSE. Model performance is measured by the Mean Squared Error (MSE), with values nearer zero indicating a better model.

3.8 Ethical Considerations

This research seeks to use data that is available publicly from a Kaggle repository, the author has consented to provide the data research purposes. The data will be pre-processed in order to make the process of analysis and model development easier. Furthermore, a research ethics review committee issued a certificate of ethical approval for the study.

CHAPTER 4 : SYSTEM ANALYSIS, DESIGN AND ARCHITECTURE

4.1 Introduction

This chapter describes the architecture used to construct the schizophrenia detection model, which is based on the conceptual model outlined earlier. This section also includes a description of the non-functional and functional system needs. Use case diagrams, system sequence diagrams, and flow charts are used to capture user engagement with the system as well as interactions between system components.

4.2 Requirements Analysis

This particular segment goes over the requirements that were extracted from the research objectives. The particular requirements were broken down into functional and nonfunctional requirements.

4.2.1 Functional Requirements

- i. The system should ensure that only authorized users can log in.
- ii. The system should be able to take in Event Related Potential Data (ERP Data)
- iii. The system should be able to detect Schizophrenia in patients through its trained model
- iv. The system should be able to display whether there is a likelihood of Schizophrenia in patients or not.

4.2.2 Non-Functional Requirements

- i. The system should be able to have a web interface for input of data and output of results.
- ii. The system should make sure that the data that is to be used is done in a secure way, to ensure that it is not accessible to those who do not have the right to use it.
- iii. Availability of the system is also very crucial; the system should be accessible when required.
- iv. Within a considerable amount of time the system should be able to predict and report the outcomes.

- v. In addition to being robust and simple to maintain the systems should have minimal downtimes.

4.3 System Architecture

The suggested Schizophrenia detection model's system architecture provides a high-level overview of the system's structure and design, as well as how different system components work together to achieve the system's goals.

The various components of the model and how they are linked are illustrated in the figure below.

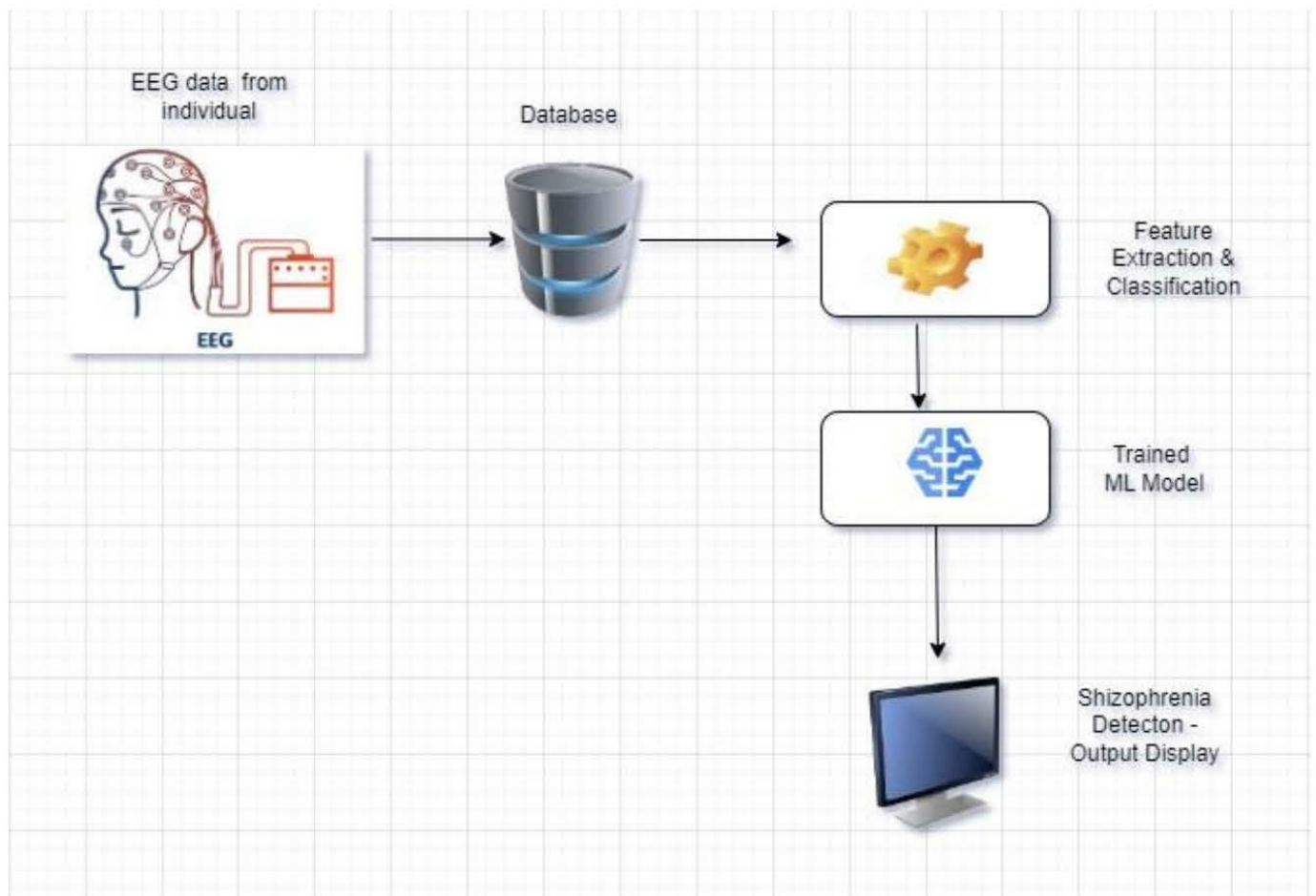


Figure 4 1 System Architecture

4.4 System Diagrams

We are going to explore some system diagrams that will help us illustrate how the various components of the Schizophrenia Prediction System Works

4.4.1 Use Case Diagram

UML's use case modeling is a widely used text-based method for system analysis and design.(Gemino & Parker, 2009) . The system is composed of these actors. The database and the administrator of the system. The administrator collects a large sample of data, trains the model and obtains a desirable accuracy, the potential patient wears an EEG contraption that records EEG signals which are then stored in a database. The administrator or physician prompts the model to fetch the data from the database and make a prediction based on model training. The system produces a prediction output determining whether the patient has schizophrenia or not.

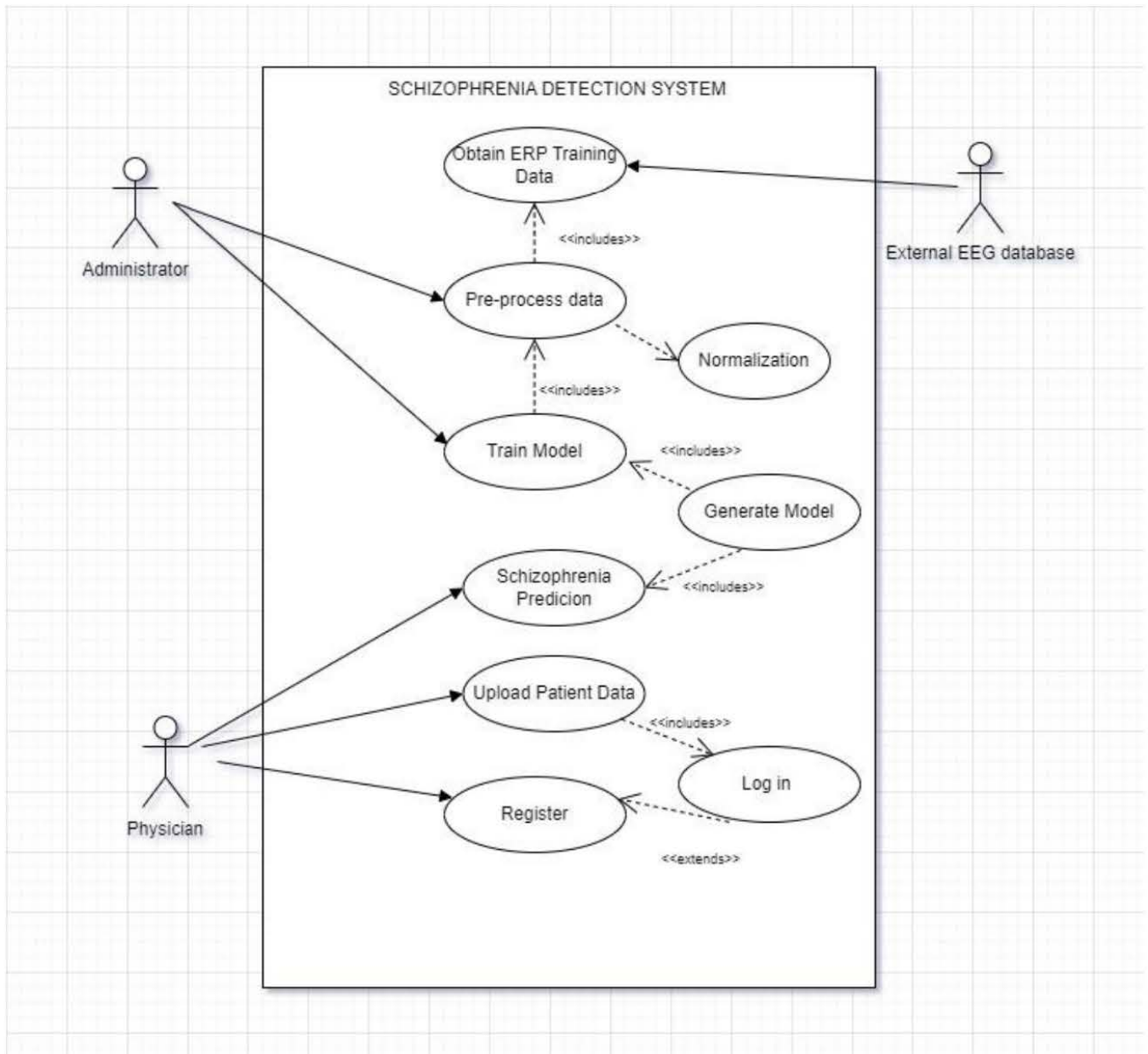


Figure 4 2 Use Case Diagram

Detailed Use Case Scenarios

Use Case	Pre- process data
Primary Actors	Administrator
Pre- Condition	Availability of Electroencephalography data set
Post - Condition	Data that has been preprocessed and ready for the model

--	--

Actor	System
Administrator obtains data in CSV format and places in folder	
Machine learning parameters are defined by the administrator	
	Extraction of features from the dataset is done by the system
	The extracted features are saved.

Use Case	Model Training
Primary Actors	Administrator
Pre-Condition	The data has been pre-processed and the machine learning algorithms is ready
Post- Condition	A trained model that can predict Schizophrenia
Main Success Scenarios	
Actor	System
The administrator initiates the training of preprocessed data	
	The system splits the dataset into train and test groups

	The system fits the data on the algorithm to be used for learning
	The system validates the model using test data.
	A model is generated by the system and is saved

Use Case	Schizophrenia Prediction
Primary Actors	Physician
Pre-Condition	Physician is a registered user Physician logged into the system
Post - Condition	Predicted Schizophrenia
Main Success Scenarios	
Actor	System
Physician Logs into the system	
Physician uploads patient data from database and prompts for prediction	
	System predicts for schizophrenia
	System generates output of Schizophrenia or healthy patient.
Physician receives potential diagnosis	

Use Case	Registration
Primary Actors	Physician
Pre-Condition	Physician has Log in Credentials
Post- Condition	Physician has access to the system
Main Success Scenario	
Actor	
Physician attempt to access	
	Physician successfully accesses the system

4.4.2 System Sequence Diagram

The sequence diagram highlights the order in which messages are sent and received as well as the structural arrangement of the sending and receiving objects.(X. Li et al., 2004), The user, being the physician in this case will enter their log in credentials. Verification of the credentials will be done and if they are correct, this will grant the user access to the system. After the physician has been granted access, they will fetch the patient data from the database and load the ERP data to the model. The model having being trained will give an output determining whether the patient has schizophrenia or not

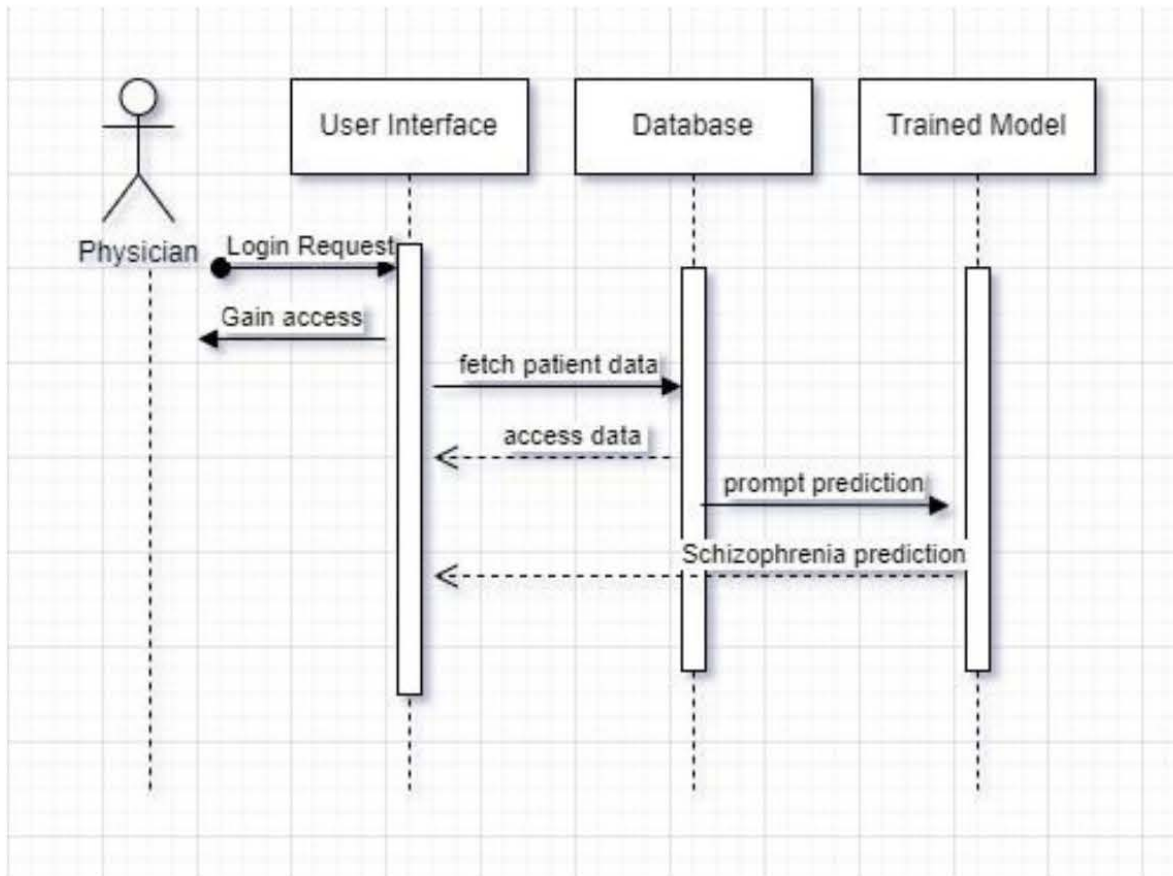


Figure 4 3 Sequence Diagram

4.4.3 Entity Relationship Diagram

The system database's current entities, their properties and how they connect to one another are illustrated by the Entity relationship diagram.(Song et al., 1995). Here are the entities and their respective relationships

Entities

- Physician
- Patient
- ERP Recording
- Schizophrenia Diagnosis

Relationships

- Physician can have multiple engagements with the system
- Patient and ERP Recording are linked as one to one
- One patient can have one Schizophrenia Diagnosis

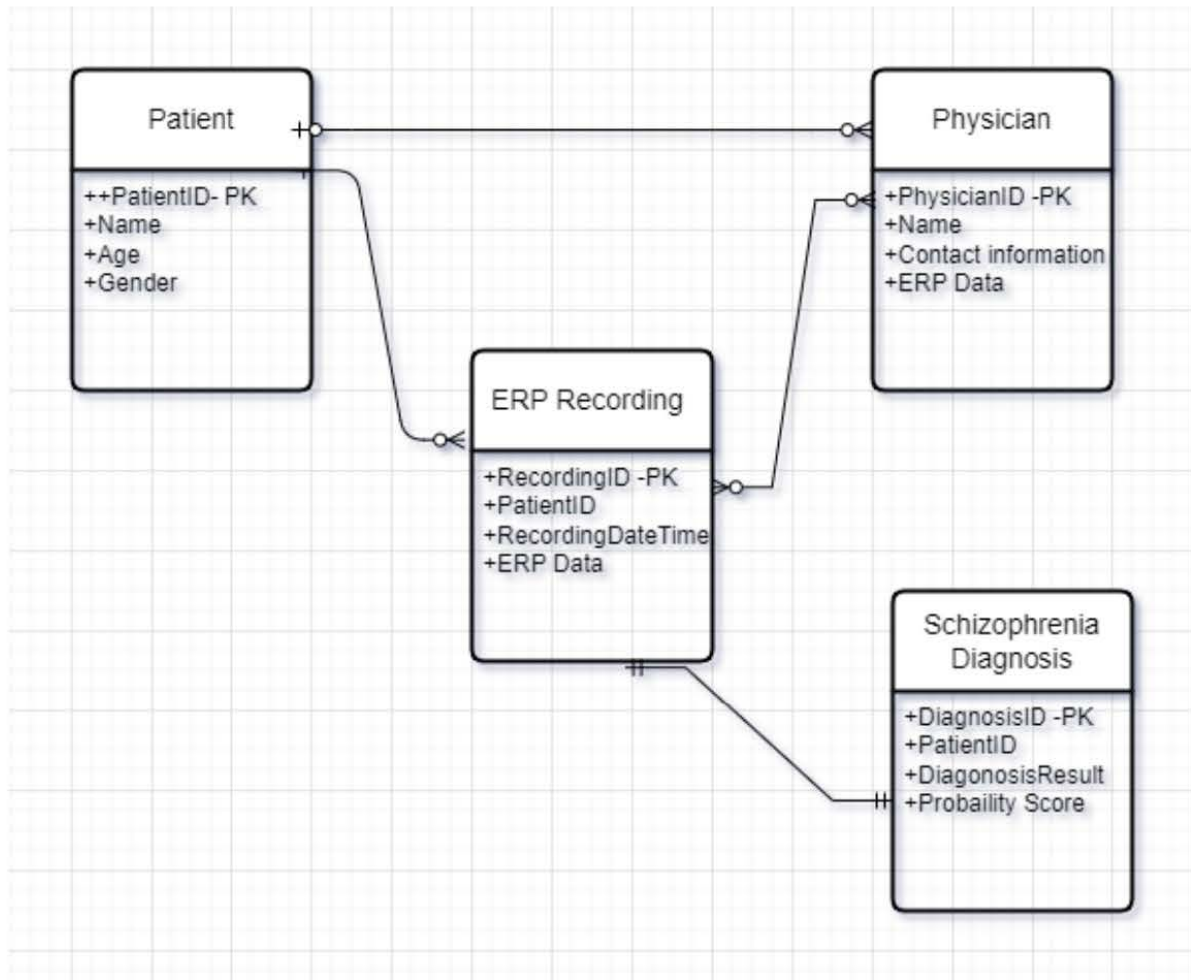


Figure 4 4 Entity Relationship Diagram

4.4.4 Wireframes

Wireframes are simplified visual representations of a user interface or the functioning of a system, they outline the basic structure, layout and content of a webpage or system interface without including much details. The following wireframe shows the conceptual parameters of the proposed system.

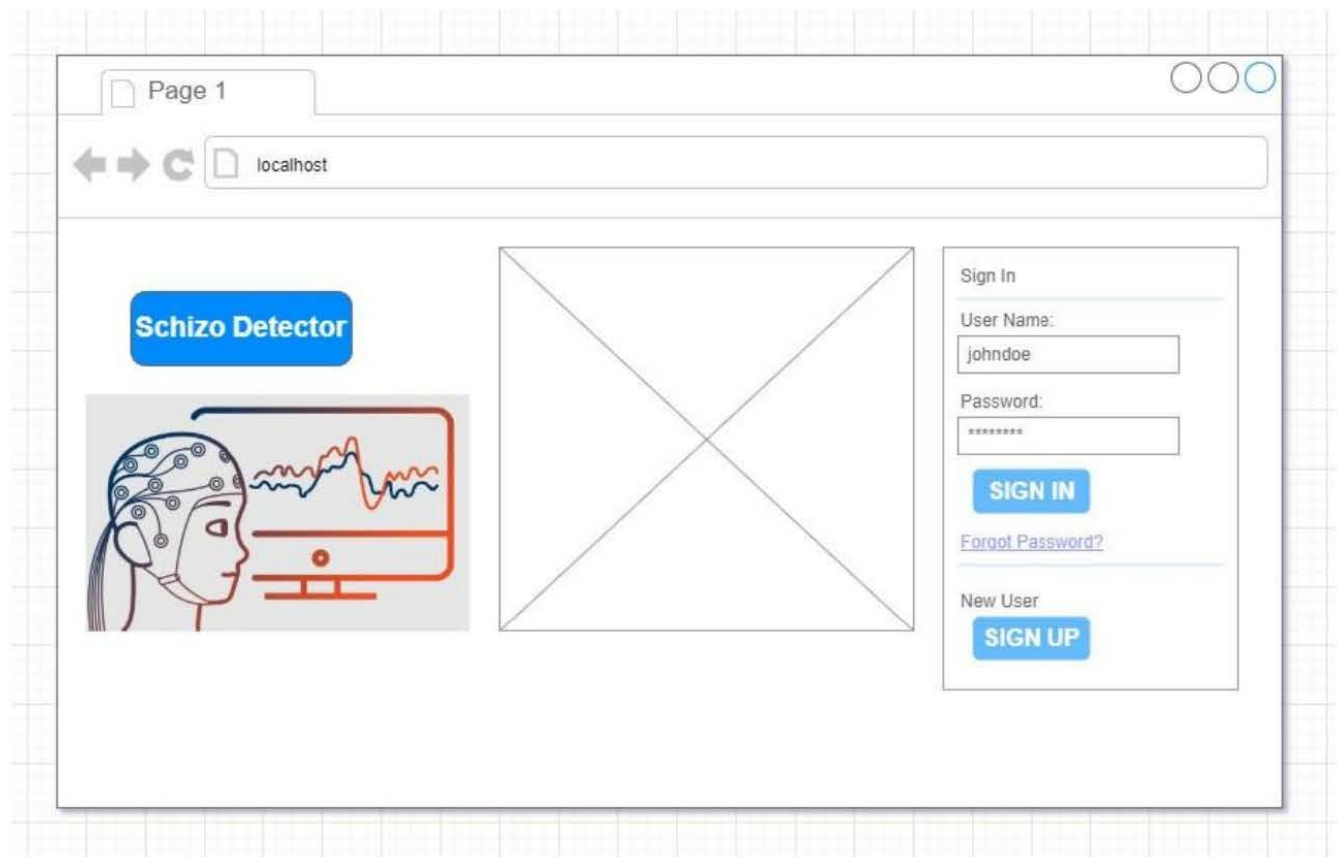


Figure 4 5 Wireframe for Log in Page

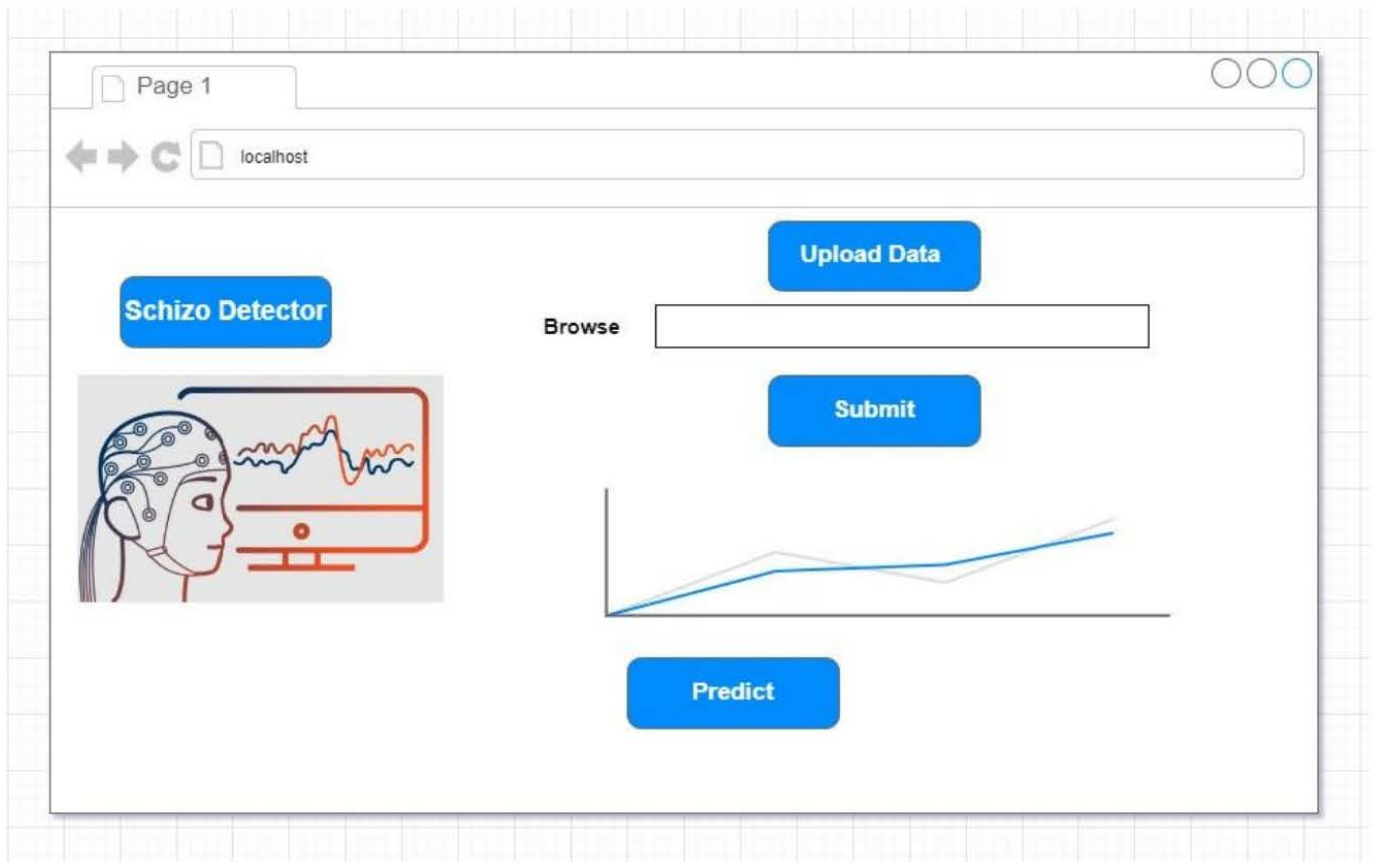


Figure 4 6 Wireframe for Uploading of Patient Data

CHAPTER 5 : SYSTEM IMPLEMENTATION AND TESTING

5.1 Introduction

System testing and implementation describe the series of steps taken to actualize and finish the system's design in accordance with system analysis and design. This particular research aimed to develop a Schizophrenia detection system by the use of EEG (electroencephalograph) data in a machine learning model. The entire design of the system as well as the required software and hardware are elucidated in the forthcoming sections. The procedure used to create, train, test, and verify the model in order to satisfy its functional requirements is also covered in length in this chapter.

5.2 System Development

5.2.1 Hardware and Software Environment

The model was created with Visual Studio text editor, which included the Jupyter Notebook and Python extensions installed. Jupyter Notebook is a great tool for running data science assignments this is mainly attributed to it allowing for running parts of code instead of an entire Python script, and because it uses markdown for documentation, it simplifies data science tasks.

The table below indicates the specifications of the computer device that was utilized in the system development.

Hardware	Specifications
Central Processing Unit	Intel ® Core™ i7-10750H Processor @ 2.60 GHz (12 CPUs)
Memory	32 GB RAM
Disk	512 GB SSD
Integrated Graphics	Intel UHD Graphics
Discrete Graphics	Nvidia 1650 Ti

In regard to software, the programming language that was used in the prediction model development was Python version 3.11.3. The Pandas and NumPy libraries were necessary due to the nature of the dataset and processing necessary. The table below shows the libraries that were to be used

Library	Version
Pandas	2.2.1
NumPy	1.24.3
Jupyter	5.7.1
Matplotlib	3.7.1
Keras	3.1.1
Scikit-learn	1.2.2
TensorFlow	2.16.1

5.3 System Implementation

5.3.1 Data Collection

ERP (Event Related Potential) data is retrieved from the Kaggle repository (Roach, n.d.) as a Comma Separated Value (CSV) file. The data consists of 12 parameter columns

- SUBJECT
- CONDITION
- 9 ELECTRODE VALUES 'Fz', 'FCz', 'Cz', 'FC3', 'FC4', 'C3', 'C4', 'CP3', 'CP4'
- TIME (MS)

5.3.2 Data preprocessing

Loading the data set on to the program and displaying the first five rows, this is done by importing the required libraries needed to manipulate and show the data.

```

import pandas as pd
import numpy as np

import pandas as pd

# Provide the path for your particular file
df = pd.read_csv('D:\Msc Thesis 2024\EEG data from basic sensory task in Schizophrenia\ERPdata Classified.csv')

# Display the first few rows of the DataFrame
df.head()

#df = pd.read_csv("ERPdata Classified.csv")

```

✓ 0.9s

subject	condition	Fz	FCz	Cz	FC3	FC4	C3	C4	CP3	CP4	Diagnosis	
0	1	1	5.533701	5.726507	5.469535	5.386723	4.588875	6.560092	4.542811	5.397492	5.103695	Healthy
1	1	1	5.651489	5.837326	5.773131	5.627975	4.822217	6.739976	4.811770	5.541357	5.379273	Healthy
2	1	1	5.717580	5.932924	5.948466	5.826460	4.979647	7.026199	5.053779	5.634972	5.600504	Healthy
3	1	1	5.703267	5.968103	5.851512	5.812192	4.992899	6.940671	5.106650	5.543577	5.589775	Healthy
4	1	1	5.571578	5.917541	5.812808	5.744715	4.963338	6.726491	5.158073	5.454069	5.614092	Healthy

Figure 5 1 Importing the data and viewing it

Checking the data for any missing values

```

# check the data for any missing values

df.isnull().sum()

```

```

subject      0
condition    0
Fz           0
FCz          0
Cz           0
FC3          0
FC4          0
C3           0
C4           0
CP3          0
CP4          0
time_ms     0
dtype: int64

```

Figure 5 2Checking data for missing values

There are no missing values in the data set

Checking for any zero or Naan values

```
▶ #checking the data for any zero or Naan values
df.isnull().sum().sum()
↳ 0
```

Figure 5.3 Checking for zero values

After loading the data, we check for any zero or Nan values in the data using the `df.isnull()` command, the return provided is zero meaning there are no values that are missing.

5.3.3 Training the model

i. KNN Classifier

The first classifier we are going to use for the model training is the KNN classifier (K Nearest Neighbor),

```
from sklearn.neighbors import KNeighborsClassifier
from sklearn.model_selection import train_test_split

# Create arrays for features and target variable
X = df[['Fz', 'FCz', 'Cz', 'FC3', 'FC4', 'C3', 'C4', 'CP3', 'CP4']].values
y = df['Diagnosis'].values

# Split the data into training and testing sets
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.25, random_state=42)

# Create a KNN classifier with 3 neighbors
knn = KNeighborsClassifier(n_neighbors=3)

# Train the classifier on the training data
knn.fit(X_train, y_train)

# Predict the labels for the test data
y_pred = knn.predict(X_test)

# Evaluate the accuracy of the classifier
from sklearn.metrics import accuracy_score
accuracy = accuracy_score(y_test, y_pred)
print(f"Accuracy: {accuracy}")

[4] ✓ 1m 2.0s Python
... Accuracy: 0.9337223508230452
```

Figure 5 4 Model with KNN Classifier

For the first classifier we chose KNN, with the value of K=3, after the data was uploaded and cleaned, arrays were created for the features and target variables. The data was then split into training and testing, the training data was 0.75 being 75% of the data and testing data was 0.25, being 25% of the data. The classifier was created and the model was trained and tested, the overall accuracy of the model was 93.37%

ii. Decision Tree Classifier

For the second model we chose the Decision Tree Classifier, the necessary libraries were loaded into the model, pandas, DecisionTreeClassifier, the training and test split and accuracy scores

```
import pandas as pd
from sklearn.tree import DecisionTreeClassifier
from sklearn.model_selection import train_test_split
from sklearn.metrics import accuracy_score

# Load the averages_without_condition_3.csv file into a pandas DataFrame
#df = pd.read_csv('averages_without_condition_3.csv')

# Create arrays for features and target variable
X = df[['Fz', 'FCz', 'Cz', 'FC3', 'FC4', 'C3', 'C4', 'CP3', 'CP4']].values
y = df['Diagnosis'].values

# Split the data into training and testing sets
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.25, random_state=42)

# Create a Decision Tree classifier
dt = DecisionTreeClassifier()

# Train the classifier on the training data
dt.fit(X_train, y_train)

# Predict the labels for the test data
y_pred = dt.predict(X_test)

# Evaluate the accuracy of the classifier
accuracy = accuracy_score(y_test, y_pred)

print(f"Accuracy: {accuracy}")
✓ 44.2s
Accuracy: 0.7739626200274349
```

Figure 5 5 Model with Decision Tree Classifier

With an equally train, test split as the KNN classifier of 0.75 and 0.25 respectively, the model gave an accuracy score of 77%

iii. Random Forest Classifier

The third model was trained with the Random Forest Classifier, with an estimator of n=100 trees, the results are shown below.

```
from sklearn.ensemble import RandomForestClassifier

import pandas as pd

#df = pd.read_csv('averages_without_condition_3.csv')

# Create arrays for features and target variable
X = df[['Fz', 'FCz', 'Cz', 'FC3', 'FC4', 'C3', 'C4', 'CP3', 'CP4']].values
y = df['Diagnosis'].values

# Split the data into training and testing sets
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.25, random_state=42)

# Create a Random Forest classifier with 100 trees
rf = RandomForestClassifier(n_estimators=100)

# Train the classifier on the training data
rf.fit(X_train, y_train)

# Predict the labels for the test data
y_pred = rf.predict(X_test)

# Evaluate the accuracy of the classifier
accuracy = accuracy_score(y_test, y_pred)
print(f"Accuracy: {accuracy}")

✓ 19m 2.7s
Accuracy: 0.8811567644032922
```

Figure 5.6 Model with Random Forest Classifier

The RF method of machine learning gave a prediction accuracy of 88%

Machine Learning Mode Accuracy

	Model Utilized	Model Accuracy
1	K Nearest Neighbor (k=3)	93.37%
2	K Nearest Neighbor (k=7)	90.98 %
3	Decision Trees	77.39%
4	Random Forest Classifier	88.11%

From the above comparison we can see the accuracies from the different models, KNN clearly gives the best overall accuracy of 93.37%, making it the best model for use in the Schizophrenia detection system. Decision trees accuracy was at 77.39% whilst the Random Forest Classifier was at 88.11%

5.3.4 Validating the Model

In terms of validating the model and having chosen KNN as our best algorithm for this scenario, we are going to evaluate its performance using metrics such as accuracy, precision and F1 score. KNN is a simple instance-based algorithm that does not have the concept of training and validation loss.

```
# Evaluate the accuracy of the classifier
accuracy = accuracy_score(y_test, y_pred)
print(f"Accuracy: {accuracy}")

# Calculate precision and F1 score
precision = precision_score(y_test, y_pred, average='macro')
f1 = f1_score(y_test, y_pred, average='macro')

# Print the results
print(f"Precision: {precision}")
print(f"F1 Score: {f1}")
```

```
↳ Accuracy: 0.9337223508230452
Precision: 0.9307741140604116
F1 Score: 0.9306233759521962
```

The model gave us an accuracy score of 93.37%, the precision was 93.07 and the F1 score as 93.06

CHAPTER 6 : DISCUSSIONS

6.1 Introduction

In this chapter we will discuss the results of the developed solution in regard to the objectives and research questions that are highlighted in the first chapter. The primary objective of this research was to come up with a tool that detects schizophrenia in patients through the use of electroencephalograph data and event related potential data

6.2 Model Validation

The process of model validation is a crucial stem in assessing the performance and reliability of machine learning algorithms. In the case for our KNN classifier we obtained good results indicating the effectiveness of the algorithm. With an accuracy score of 93.37 % our model demonstrates a high level of correctness in its prediction. Additionally, both the precision and F1 score at 93.07 and 93.03 respectively, highlight the model's capability to correctly classify instances of interest while minimizing false positives and negatives. These metrics validate the sturdiness of the KNN classifier. Making it capable as the preferred model for schizophrenia detection in patients.

6.3 Model Performance

The study continued to establish the performance of the KNN Model in comparison to other models. The KNN classifier performed better than the Decision Trees and Random Forest Models in terms of the accuracy of prediction, while the KNN model achieved an impressive score of 93.37% the Decision Tree model produced a lower accuracy of 77.39. Similarly, although the Random Forest model had a better performance compared to Decision Trees with an accuracy of 88.11%, it still was not up to par with the KNN classifier.

It is also important to consider that KNN model also showed higher precision and F1 scores, indicating a good balance between true positives false positives and false negatives.

However, it is essential to consider other factors such as how computationally efficient the model is and how easy it is to understand and explain how the model has performed.

KNN algorithm usually allows for easy interpretation of results but sometimes may lack the model sophistication to capture intricate patterns in the data compared to other complex models.

6.4 Contribution to Research

The model that has been developed provides a solution in early detection of Schizophrenia in patients. The model can be deployed in a production environment through a web interface. The platform enables a physician to upload a patient's ERP data (after the patient has been subjected to a similar stimuli). After the data has been uploaded it will be processed through the trained model in order to give a prediction on whether they have Schizophrenia or whether they are healthy.

CHAPTER 7 : CONCLUSION AND RECOMMENDATION

7.1 Conclusion

Early identification of the onset and advanced stages of schizophrenia is crucial for prompt and efficient therapy to stop or lessen the disease's progression. This particular study purposed to develop a tool to detect Schizophrenia in patients by use of a machine learning model, in order to accomplish this, it was paramount to understand the psychotic manifestations of the disease, we also explored the various implications the disease has on the brain's functional capacity and how it affects neural activities in response to certain stimulus. In addition to this we explored literature highlighting the analysis of EEG signals in relation to Schizophrenia. The second chapter also covered various forms of machine learning models seeking to differentiated patients with Schizophrenia and Healthy Controls.

The primary objective of the study was to develop Machine Learning model for Schizophrenia detection by use of Electroencephalograph EEG Data, EEG Data has proven to be very essential in identifying anomalies especially in mentally related disorders, this explains why it was advantageous to use this in order to enable the diagnosis of Schizophrenia in patients.

The model was tested on a variety of machine learning methods in order to obtain the best accuracy , out of the three methods that were used , K Nearest Neighbor managed to present the best accuracies for the project with an accuracy of 93% , this was KNN with 3 Neighbors, the other models that were used were Decision Trees , with an accuracy score of 77% , the use of Random Forest Classifier produced an accuracy of 88 % and the final model was Support Vector Machine.

7.2 Recommendation.

Considering the factors and possible drawbacks mentioned before, the following recommendations are made for creating a machine learning model to identify schizophrenia in patients:

- i. The model could be trained with more data to increase the accuracy levels in determining which patients have schizophrenia and which ones do not.
- ii. The model could be used to give diagnosis on a severity spectrum, especially for patients with Schizophrenia, this can be done in addition to having a diagnosis of two classes of either Schizophrenic or Healthy.
- iii. The use of neural networks as a classifier would greatly improve the model's performance and accuracy.

7.3 Future work

The use of machine learning as a tool for Computer Aided Diagnosis has great potential in the medical field, especially in the use of EEG data to diagnose mental disorders. As technology advances and more improvements are made, there are various areas which can be improved upon in the future.

- i. Data collection procedures can be improved upon, this can be done by the incorporation of multimodal data sources and the use of advanced machine learning methods.
- ii. Machine learning technologies can be integrated into healthcare systems and clinical workflows in order to improve diagnostic accuracy.

REFERENCES

- Anyanwu, M. N., & Shiva, S. G. (2009). Comparative Analysis of Serial Decision Tree Classification Algorithms. *International Journal of Computer Science and Security*, 3(3), 230–240.
<http://www.cscjournals.org/csc/manuscriptinfo.php?ManuscriptCode=72.73.66.82.82.44.55.49.99>
- Bae, Y. J., Shim, M., & Lee, W. H. (2021). Schizophrenia detection using machine learning approach from social media content. *Sensors*, 21(17), 1–18. <https://doi.org/10.3390/s21175924>
- Baess, P., Horváth, J., Jacobsen, T., & Schröger, E. (2011). Selective suppression of self-initiated sounds in an auditory stream: An ERP study. *Psychophysiology*, 48(9), 1276–1283.
<https://doi.org/10.1111/j.1469-8986.2011.01196.x>
- Barros, C., Silva, C. A., & Pinheiro, A. P. (2021). Advanced EEG-based learning approaches to predict schizophrenia: Promises and pitfalls. *Artificial Intelligence in Medicine*, 114(August 2020), 102039.
<https://doi.org/10.1016/j.artmed.2021.102039>
- Belgiu, M., & Drăgu, L. (2016). Random forest in remote sensing: A review of applications and future directions. *ISPRS Journal of Photogrammetry and Remote Sensing*, 114, 24–31.
<https://doi.org/10.1016/j.isprsjprs.2016.01.011>
- Beynon-Davies, P., Came, C., Mackay, H., & Tudhope, D. (1999). Rapid application development (Rad): An empirical review. *European Journal of Information Systems*, 8(3), 211–232.
<https://doi.org/10.1057/palgrave.ejis.3000325>
- Caprihan, A., Pearlson, G. D., & Calhoun, V. D. (2008). Application of principal component analysis to distinguish patients with schizophrenia from healthy controls based on fractional anisotropy measurements. *NeuroImage*, 42(2), 675–682. <https://doi.org/10.1016/j.neuroimage.2008.04.255>
- Charbuty, B., & Abdulazeez, A. (2021). Classification Based on Decision Tree Algorithm for Machine Learning. *Journal of Applied Science and Technology Trends*, 2(01), 20–28.
<https://doi.org/10.38094/jastt20165>
- Chatterjee, I., Agarwal, M., Rana, B., Lakhyani, N., & Kumar, N. (2018). Bi-objective approach for computer-aided diagnosis of schizophrenia patients using fMRI data. *Multimedia Tools and Applications*, 77(20), 26991–27015. <https://doi.org/10.1007/s11042-018-5901-0>
- Chauhan, T., Rawat, S., Malik, S., & Singh, P. (2021). Supervised and Unsupervised Machine Learning based Review on Diabetes Care. *2021 7th International Conference on Advanced Computing and Communication Systems, ICACCS 2021*, 581–585.

<https://doi.org/10.1109/ICACCS51430.2021.9442021>

Dulock, H. L. (1993). Research Design: Descriptive Research. *Journal of Pediatric Oncology Nursing*, 10(4), 154–157. <https://doi.org/10.1177/104345429301000406>

Feinberg, I., & Guazzelli, M. (1999). Schizophrenia - A disorder of the corollary discharge systems that integrate the motor systems of thought with the sensory systems of consciousness. *British Journal of Psychiatry*, 174(MAR.), 196–204. <https://doi.org/10.1192/bjp.174.3.196>

Ford, J. M., & Mathalon, D. H. (2005). Corollary discharge dysfunction in schizophrenia: Can it explain auditory hallucinations? *International Journal of Psychophysiology*, 58(2-3 SPEC. ISS.), 179–189. <https://doi.org/10.1016/j.ijpsycho.2005.01.014>

Ford, J. M., Palzes, V. A., Roach, B. J., & Mathalon, D. H. (2014). Did i do that? Abnormal predictive processes in schizophrenia when button pressing to deliver a tone. *Schizophrenia Bulletin*, 40(4), 804–812. <https://doi.org/10.1093/schbul/sbt072>

Gemino, A., & Parker, D. (2009). Use case diagrams in support of use case modeling: Deriving understanding from the picture. *Journal of Database Management*, 20(1), 1–24. <https://doi.org/10.4018/jdm.2009010101>

Golnaz Baghdadi, Farzad Towhidkhan, M. R. (2021). *Neurocognitive Mechanisms of Attention*.

Grabmeier, J. L., & Lambe, L. A. (2007). Decision trees for binary classification variables grow equally with the Gini impurity measure and Pearson's chi-square test. *International Journal of Business Intelligence and Data Mining*, 2(2), 213–226. <https://doi.org/10.1504/IJBIDM.2007.013938>

Herwig, U., Satrapi, P., & Schönfeldt-Lecuona, C. (2003). Using the International 10-20 EEG System for Positioning of Transcranial Magnetic Stimulation. *Brain Topography*, 16(2), 95–99. <https://doi.org/10.1023/B:BRAT.0000006333.93597.9d>

IBM. (n.d.). *What is the KNN Algorithm*. <https://www.ibm.com/topics/knn>

Insel, T. R. (2010). Rethinking schizophrenia. *Nature*, 468(7321), 187–193. <https://doi.org/10.1038/nature09552>

Jo, Y. T., Joo, S. W., Shon, S. H., Kim, H., Kim, Y., & Lee, J. (2020). Diagnosing schizophrenia with network analysis and a machine learning method. *International Journal of Methods in Psychiatric Research*, 29(1), 1–11. <https://doi.org/10.1002/mpr.1818>

Kramer, O. (2013). Dimensionality Reduction with Unsupervised Nearest Neighbors. *Intelligent Systems*

Reference Library, 51, 13–14. <https://doi.org/10.1007/978-3-642-38652-7>

- Li, F., Wang, J., Liao, Y., Yi, C., Jiang, Y., Si, Y., Peng, W., Yao, D., Zhang, Y., Dong, W., & Xu, P. (2019). Differentiation of Schizophrenia by Combining the Spatial EEG Brain Network Patterns of Rest and Task P300. *IEEE Transactions on Neural Systems and Rehabilitation Engineering*, 27(4), 594–602. <https://doi.org/10.1109/TNSRE.2019.2900725>
- Li, X., Liu, Z., & Jifeng, H. (2004). A formal semantics of UML sequence diagram. *Proceedings of the Australian Software Engineering Conference, ASWEC, 2004*, 168–177. <https://doi.org/10.1109/aswec.2004.1290469>
- Li, Z., Liu, F., Yang, W., Peng, S., & Zhou, J. (2022). A Survey of Convolutional Neural Networks: Analysis, Applications, and Prospects. *IEEE Transactions on Neural Networks and Learning Systems*, 33(12), 6999–7019. <https://doi.org/10.1109/TNNLS.2021.3084827>
- Light, G. A., Williams, L. E., Minow, F., Sprock, J., Rissling, A., Sharp, R., Swerdlow, N. R., & Braff, D. L. (2010). Electroencephalography (EEG) and event-related potentials (ERPs) with human participants. *Current Protocols in Neuroscience, SUPPL. 52*, 1–24. <https://doi.org/10.1002/0471142301.ns0625s52>
- Lindström, E., Wieselgren, I. -M, & von Knorring, L. (1994). Interrater reliability of the Structured Clinical Interview for the Positive and Negative Syndrome Scale for schizophrenia. *Acta Psychiatrica Scandinavica*, 89(3), 192–195. <https://doi.org/10.1111/j.1600-0447.1994.tb08091.x>
- M. I. Jordan, & T. M. Mitchell. (2015). Machine learning: Trends, perspectives, and prospects. *Science*, 349(6245), 255–260.
- Mahesh, B. (2020). Machine Learning Algorithms - A Review. *International Journal of Science and Research*, 9(1), 381–386. <https://doi.org/10.21275/ART20203995>
- Meltzer, H. Y. (2004). Cognitive factors in schizophrenia: Causes, impact, and treatment. *CNS Spectrums*, 9(10 SUPPL. 11), 15–24. <https://doi.org/10.1017/s1092852900025098>
- Mohri, M., Rostamizadeh, A., & T. (2018). *Foundations of machine learning*. Cambridge: MIT Press.
- Nayani, T. H., & David, A. S. (1996). The auditory hallucination: A phenomenological survey. *Psychological Medicine*, 26(1), 177–189. <https://doi.org/10.1017/s003329170003381x>
- Oh, S. L., Vicnesh, J., Ciaccio, E. J., Yuvaraj, R., & Acharya, U. R. (2019). Deep convolutional neural network model for automated diagnosis of Schizophrenia using EEG signals. *Applied Sciences*

(Switzerland), 9(14). <https://doi.org/10.3390/app9142870>

Petkovic, D., Altman, R., Wong, M., & Vigil, A. (2018). Improving the explainability of Random Forest classifier – User centered approach. *Pacific Symposium on Biocomputing*, 0(212669), 204–215. https://doi.org/10.1142/9789813235533_0019

Pynn, L. K., & DeSouza, J. F. X. (2013). The function of efference copy signals: Implications for symptoms of schizophrenia. *Vision Research*, 76, 124–133. <https://doi.org/10.1016/j.visres.2012.10.019>

Roach, B. J. (n.d.). *EEG Data for Basic Sensory Task in Schizophrenia*. <https://www.kaggle.com/datasets/broach/button-tone-sz/data>

Shim, M., Hwang, H. J., Kim, D. W., Lee, S. H., & Im, C. H. (2016). Machine-learning-based diagnosis of schizophrenia using combined sensor-level and source-level EEG features. *Schizophrenia Research*, 176(2–3), 314–319. <https://doi.org/10.1016/j.schres.2016.05.007>

Song, I.-Y., Evans, M., & Park, U. E. K. (1995). A Comparative Analysis of Entity-Relationship Diagrams 1. *Journal of Computer and Software Engineering*, 3(4), 427–459.

Sperry, R. W. (1950). Neural basis of the spontaneous optokinetic response produced by visual inversion. *Journal of Comparative and Physiological Psychology*, 43(6), 482–489. <https://doi.org/10.1037/h0055479>

Stam, C. J. (2004). Functional connectivity patterns of human magnetoencephalographic recordings: A “small-world” network? *Neuroscience Letters*, 355(1–2), 25–28. <https://doi.org/10.1016/j.neulet.2003.10.063>

Tibbets, P. E. (2013). Principles of Cognitive Neuroscience. *Quarterly Review of Biology Principles of Cognitive Neuroscience*, 88, 139–140.

Whiteford, H. A., Degenhardt, L., Rehm, J., Baxter, A. J., Ferrari, A. J., Erskine, H. E., Charlson, F. J., Norman, R. E., Flaxman, A. D., Johns, N., Burstein, R., Murray, C. J. L., & Vos, T. (2013). Global burden of disease attributable to mental and substance use disorders: Findings from the Global Burden of Disease Study 2010. *The Lancet*, 382(9904), 1575–1586. [https://doi.org/10.1016/S0140-6736\(13\)61611-6](https://doi.org/10.1016/S0140-6736(13)61611-6)