

**MACHINE LEARNING AND ITS APPLICATIONS**

**BY**

**OGUNODE OYINDAMOLA DEBORAH**

**PSC1809104**

**A PROJECT WORK SUBMITTED TO FACULTY OF  
PHYSICAL SCIENCES IN PARTIAL FULFILMENT OF THE  
AWARD OF BACHELOR OF SCIENCE (BSc.) MATHEMATICS,  
UNIVERSITY OF BENIN, BENIN CITY, EDO STATE, NIGERIA.**

**SEPTEMBER, 2023**

## DECLARATION

This project work was carried out by OGUNODE OYINDAMOLA DEBORAH. I have not copied the work of any author. All work has been duly cited and acknowledged.

---

Ogunode Oyindamola D.

---

Date

## CERTIFICATION

This is to certify that this project work was carried out and submitted by **OGUNODE OYINDAMOLA DEBORAH** with matriculation number **PSC1809104** under my supervision. It is adequate and satisfactory, both in scope and content, for the award of Bachelor of Science (B.Sc) in Mathematics of the University of Benin.

---

**Prof. D. Omorogbe**

**Project Supervisor**

---

**Date**

---

**Prof. Robert Okuoghae**

**HEAD OF DEPARTMENT**

---

**Date**

## **DEDICATION**

I dedicate this work to God Almighty for his love and strength to see this work through and also to my loving parents Mr Solomon Ogunode and Mrs Mercy Ogunode and my siblings.

## ACKNOWLEDGEMENT

I am genuinely grateful to the Almighty God, who has shown me great love and showering me with loving kindness all through my academic career in the University of Benin.

I also wish to express my sincere gratitude to my project supervisor, Professor .D. Omorogbe for his professional guidance and directions. God bless you immensely sir.

I will not fail to appreciate my other lecturers in the Department of Mathematics who have in one way or the other impacted me. God bless you all.

I am forever indebted to my parents, Mr Solomon Ogunode and Mrs Mercy Ogunode for your love, care, prayers and for financially supporting me all through my journey in school. May God continue to bless you. I'm grateful to my siblings, Mrs Agenyi Oluwafunmilayo, Mrs Adeola Oluwaseyi, Mrs Rasaan Ifeoluwa and Miss Ogunode Boluwatife and also to my mentor, Mr Animashawun Ebenezer for their constant love and care all through my journey in the university. God bless you all.

I will not fail to appreciate my friends, and also My Coursemates for being there for me. I love you all.

## ABSTRACT

Machine learning has emerged as a transformative technology with a profound impact on various industries.

In this course of study, this abstract provides an overview of ML, its significance, limitations, types of ML we have.

**Machine Learning:** This refers to the exploration of computer programs that utilize algorithm and statistical model to acquire knowledge by identifying patterns and making inferences all without explicit programming.

Some of the significance of Machine Learning includes its contribution to technological advancement in various industries, it equips individual with skills to automate processes and streamline operations and it also enables automation of repetitive tasks which enhance efficiency and productivity.

There are challenges or limitations associated with ML, some of which includes: Data dependency, Interpretability and Transparency, Overfitting and Generalization, Domain-specific expertise and Lack of casual understanding.

The types of Machine Learning are: Supervised Learning, Unsupervised Learning, Semi-supervised Learning and Reinforcement Learning.

## TABLE OF CONTENTS

Title page

Declaration

Certification

Dedications

Acknowledgement

Abstract

Table of contents

### CHAPTER ONE: GENERAL INTRODUCTION

|     |                          |       |     |
|-----|--------------------------|-------|-----|
| 1.1 | Background of the study  | ----- | 1-2 |
| 1.2 | Statement of the problem | ----- | 2-3 |

|     |                                  |       |     |
|-----|----------------------------------|-------|-----|
| 1.3 | Aims and Objectives of the study | ----- | 3   |
| 1.4 | Motivation for the study         | ----- | 4   |
| 1.5 | Significance of the study        | ----- | 4-5 |
| 1.6 | Scope of the study               | ----- | 5   |
| 1.7 | Limitations of the study         | ----- | 6   |
| 1.8 | Definitions of terms             | ----- | 6-7 |
| 1.9 | Conclusions                      | ----- | 7   |

## CHAPTER TWO: LITERATURE REVIEW

|     |  |       |       |
|-----|--|-------|-------|
| 2.1 | What is Machine Learning?                | ----- | 8-9   |
| 2.2 | Types of Machine learning                | ----- | 9-16  |
| 2.3 | Mathematical models for Machine Learning | ----- | 17-19 |
| 2.4 | Importance of Machine Learning           | ----- | 20    |
| 2.5 | Limitations of Machine Learning          | ----- | 21-22 |
| 2.6 | Conclusions                              | ----- | 22    |

## CHAPTER THREE: METHODOLOGY

|     |   |       |       |
|-----|---|-------|-------|
| 3.1 | Introduction                            | ----- | 23    |
| 3.2 | Mathematical models in Machine Learning | ----- | 23-65 |
| 3.3 | Conclusions                             | ----- | 66    |

## CHAPTER FOUR: MACHINE LEARNING AND ITS APPLICATIONS

|     |                                  |       |       |
|-----|----------------------------------|-------|-------|
| 4.1 | Introduction                     | ----- | 67    |
| 4.2 | Applications of Machine learning | ----- | 67-84 |
| 4.3 | Discussion of Findings           | ----- | 85    |
| 4.4 | Conclusions                      | ----- | 86    |

## CHAPTER FIVE: CONCLUSIONS AND RECOMMENDATIONS

|     |                |       |    |
|-----|----------------|-------|----|
| 5.1 | Recommendation | ----- | 87 |
| 5.2 | Conclusions    | ----- | 88 |

## References

# **CHAPTER ONE**

## **GENERAL INTRODUCTION**

### **1.1 BACKGROUND OF THE STUDY**

These days, ML (Machine Learning) techniques are being used in some amazing applications like self-driving cars, natural language processing, and facial recognition systems. The first paper on machine learning (McCulloch and Pitts, 1943), which examined and applied neurons to ML, was written in the year 1943. They built an electrical circuit model, and as a result, the neural network was born.

English mathematician Alan Turing was the one who created the renowned Turing test (1950). The goal of the test is to determine whether a machine can act intelligently in a manner that is comparable to that of a human. To find out if computers had true intelligence, we underwent this test. For the test to pass, it needs to convince the tester that it is a human and not a computer. The first computer program that could learn as it played the game of checkers and determine its chances of winning or losing was created (Arthur Samuel, 1952). The first neural network, dubbed the "Perceptron" (Frank Rosenblatt, 1957), was created. Due to the availability of vast quantities of data in the 1990s, there was a massive transition in machine learning (ML) from a knowledge-driven to a data-driven

technique. In the chess game, the first machine to defeat the world champion was IBM's Deep Blue, which was created in 1997.

Businesses are now aware that machine learning has the potential to increase the efficiency of complex calculations. These are some of the most recent projects: A deep neural network called Google Brain was created in 2012 with the goal of recognizing patterns in pictures and videos. Later, it was used to find items in YouTube videos. Facebook developed Deep Face in 2014, which can recognize individuals in a manner similar to that of humans. Alpha Go, a computer program developed by Deep Mind in 2014, defeated a skilled Go player in a game of Go. The complexity of the game is said to make it a difficult yet traditional AI game.

According to scientists Stephen Hawking and Stuart Russell, if AI develops the capacity to redesign itself at an increasingly rapid rate, an unstoppable "intelligence explosion" may cause the extinction of humanity. Musk refers to AI as the "biggest existential threat" to humanity. Elon Musk established the "Open AI" organization in 2015 with the goal of creating a friendly and secure AI that could be useful to humanity.

The purpose of this narrative is to examine machine learning and its implications for national development.

## **1.2 STATEMENT OF THE PROBLEM**

The issue with machine learning is the challenges involved in creating algorithms and systems that can analyze data, learn from it, and make predictions or decisions on their own. Here are a few of the major issues with ML:

- The scarcity of adequate data
- The presence of biases in data leading to unfair outcomes
- The lack of interpretability and explainability of complex models
- The challenge of achieving good generalization
- The computational demands of training models
- The vulnerability to adversarial attacks
- The ethical considerations surrounding the use of ML

However, the study is to examine machine learning and its applications for national development.

## **1.3 AIMS AND OBJECTIVES OF THE STUDY**

The aim of the study is to give an insight into machine learning and its applications for national development.

The objectives of the study:

- To do a brief literature search into the concept of machine learning.

- To gain better knowledge and insight into machine learning.
- To examine the applications of machine learning for national development.

#### **1.4 MOTIVATION FOR THE STUDY**

Studying machine learning has become popular as a way to take advantage of the rapid advancements in computing power, algorithm development, and data accessibility to develop intelligent systems and solve complex problems. It inspires and motivates people to study machine learning (ML) and explore new avenues, pushing the envelope of what is possible because it is recognized as a field with enormous untapped potential.

In general, the drive behind machine learning stems from an understanding of the importance of data in reaching well-informed decisions. Utilizing the potential for automation provided by ML algorithms, it is also possible to automate processes, lower the amount of manual labor required, and improve efficiency in a variety of domains.

However, the study is motivated to provide a better understanding and insight into machine learning and its applications for national development.

#### **1.5 SIGNIFICANCE OF THE STUDY**

Learning about ML advances technology and fosters innovation across a range of sectors, spurring the creation of novel algorithms, methods, and applications. The

significance is in using ML to extract important insights from data, enabling informed decision-making, process optimization, and improved overall performance.

By analyzing user data, machine learning (ML) enables the automation of repetitive tasks, increases efficiency, productivity, resource allocation, and creates personalized experiences. It also equips individuals with the skills to create tailored recommendations, content, and interfaces to increase user satisfaction and engagement.

The ability to build predictive models is a skill acquired through ML study, allowing for proactive decision-making, risk mitigation, and opportunity capitalization. Understanding the ethical ramifications and addressing fairness issues are also part of studying ML. It enables precise predictions based on historical data analysis. By taking into account variables like bias, interpretability, privacy, and transparency, it encourages the responsible development and deployment of impartial ML systems.

## **1.6 SCOPE OF THE STUDY**

The scope of this study extends from the definition, categories of ML and types of ML we have and how to put it in use for national development.

## 1.7 LIMITATIONS OF THE STUDY

It can be difficult to be a student while also writing a project and getting ready for exams. It was expensive, and there wasn't much time to read and do research for the exam.

Additionally, the lack of light can occasionally be exhausting, which results in the research being suspended temporarily due to dead equipment.

## 1.8 DEFINITION OF TERMS

- i. **Machine learning:** This area of AI focuses on creating models and algorithms that let computers learn from data without having to be explicitly programmed.
- ii. **Data:** Data in ML is any information or examples used to train a model. It may consist of various types of structured or unstructured information, such as numbers, text, images, or audio.
- iii. **Training data:** The subset of available data used to train an ML model is known as training data. The model can learn patterns and relationships by using input examples and their corresponding labels or outputs.

- iv. **Model:** A model is a mathematical or computational representation that can make predictions or decisions based on input data. It is created by training a ML algorithm using training data that can be used to predict on new data.
- v. **Algorithm:** An algorithm is a set of mathematical rules and procedures used to train a model and make predictions based on new data. It guides the step for processing data and adjusting the model's parameters to optimize performance.
- vi. **Prediction:** A trained model produces a prediction when presented with fresh data. It displays the model's estimation or decision made using training-derived patterns.

## CONCLUSIONS

In conclusion, machine learning is a sophisticated and promising field that has many potential applications. While there are limitations to its use, the pros far outweigh the challenges.

As this technology continues to develop, it will certainly have a significant influence on our lives for national development.

## **CHAPTER TWO**

### **LITERATURE REVIEW**

#### **2.1 WHAT IS MACHINE LEARNING?**

Machine learning is the study of computer programs that use algorithms and statistical models to learn by spotting patterns and drawing conclusions—all without any explicit programming. It is a feature of artificial intelligence that enables software programs to improve their capacity to predict outcomes without the need for explicit instructions. Machine learning algorithms use historical data as input to forecast future output values.

The ability of machines to learn on their own from data and past experiences is referred to as machine learning, and it falls under the umbrella of artificial intelligence (AI). This allows machines to identify patterns and predict outcomes with little to no human intervention. It is categorized as AI and focuses on giving machines human-like behavior and decision-making capabilities by allowing them to learn and evolve their own programs. As machines travel through the world, they learn and develop in a process that is automated and improved.

Iterative learning and pattern recognition are made possible by machine learning, which uses algorithms to draw insightful conclusions from large datasets. As a result, artificial intelligence (AI) has moved beyond its traditional function of carrying out predefined tasks. In the past, AI programs mainly automated fundamental tasks in businesses and enterprises, like intelligent automation and rule-based classification. As a result, the capabilities of AI algorithms were limited to the specific domains for which they were intended. However, the development of machine learning has allowed computers to go beyond their preprogrammed capabilities and advance through continuous iterations, constantly evolving their capabilities.

Inherent evolution and the use of a variety of programming techniques to quickly analyze massive amounts of data and derive insightful conclusions set machine learning apart from artificial intelligence.

## **2.2 TYPES OF MACHINE LEARNING**

- 1. Supervised Learning:** In this type of machine learning, known as supervised learning, machines are trained using labeled datasets to make predictions based on the provided training. The labeled dataset contains information where certain input and output parameters are already associated. Consequently, the machine is trained using the input data and the

corresponding output. In subsequent stages, the machine employs a test dataset to predict outcomes.

Consider a dataset of pictures of parrots and crows as an example. The machine is initially trained to recognize a variety of visual cues, including the color, shape, and size of both parrots and crows, as well as their eyes and feathers. After training, a parrot image is presented as input, and the computer is expected to recognize the bird and predict what will happen. To arrive at the final prediction, the trained machine looks at a variety of features in the input image, including the object's eyes, color, shape, etc. This method illustrates how object identification is accomplished using supervised machine learning.

The input variable and the output variable must be mapped as the technique's main goal.

### **Categories of Supervised Learning**

- **Classification:** This category includes categorical output variables and algorithms made to address classification issues. As a result, the output can be categorized into different groups, such as ‘yes or no’, “true or false”, “male or female”, and so on. The Random Forest algorithm, Decision Tree algorithm, Logistic Regression algorithm, Support Vector Machine algorithm, and others are examples of well-known classification algorithms.

- **Prediction:** Regression algorithms are used to solve prediction problems where there is a linear relationship between the input and output variables. These algorithms have the ability to forecast continuous output variables. Examples of these applications include market trend analysis, weather forecasting, and more. The Simple Linear Regression Algorithm, Multivariate Regression Algorithm, Decision Tree Algorithm, and Lasso Regression are all well-known regression algorithms.

### **Advantages of Supervised Learning**

- Supervised learning utilizes labeled datasets, providing precise information about the classes of objects.
- This type of algorithm is valuable for making predictions based on previous experiences.

### **Disadvantages**

- These algorithms may struggle with intricate tasks.
- If the test data differs significantly from the training data, they might produce incorrect predictions.
- Training these algorithms can be computationally time-consuming.

### **Applications of Supervised Learning**

- Image segmentation

- Medical diagnosis
- Fraud detection
- Spam detection
- Speech recognition

2. **Unsupervised Learning:** Unsupervised learning refers to a teaching method without direction or supervision. This technique trains a computer to make predictions on its own using data that has no labels or annotations. An unsupervised learning algorithm's goal is to cluster or group the disorganized dataset based on analogies, contrasts, and patterns in the inputs.

Take a group of pictures showing a container with various fruits in it as an example. The machine learning model doesn't know about these images. When we give the ML model this dataset, it must identify patterns among the objects, such as color, shape, or discrepancies seen in the input images, and classify them appropriately. When given a test dataset and categorization that has been completed, the machine can then make predictions.

### **Categories of Unsupervised Learning**

- **Clustering:** Using the clustering technique, objects are grouped according to characteristics like their similarities or differences. The classification of customers according to the goods they buy is an illustration of clustering. K-

Means clustering, Mean-Shift, DBSCAN, Principal Component Analysis (PCA), and Independent Component Analysis (ICA) are a few popular clustering algorithms.

- **Association:** Identifying common relationships between variables in sizable datasets is the goal of association learning. It entails mapping related variables and figuring out the dependencies between various data points. This method is frequently employed in applications like market data analysis and web usage mining. The Apriori, Eclat, and FP-Growth algorithms are well-known methods for identifying association rules.

### **Advantages Unsupervised Learning**

- Since unsupervised algorithms work with unlabeled datasets, they can handle more complicated tasks than their supervised counterparts.
- Unlabeled data is typically easier to obtain than labeled datasets, which is why unsupervised algorithms are preferred for a variety of tasks.

### **Disadvantages**

- As they operate on unlabeled datasets and are not trained with precise output information beforehand, unsupervised algorithms may produce results with lower accuracy.

- Unsupervised learning can be more difficult to work with because unlabeled datasets that don't directly correspond to particular output mappings must be handled.

## **Applications of Unsupervised Learning**

- Network analysis
- Recommendation system
- Anomaly Detection

3. **Semi-supervised learning:** Both supervised and unsupervised machine learning techniques are combined in semi-supervised learning. It trains its algorithms using both labeled and unlabeled datasets. Semi-supervised learning overcomes the drawbacks of the aforementioned options by utilizing both types of data.

Take a college student learning a concept as an example to help illustrate. Under the direct supervision and direction of a teacher in a classroom setting, a student learns the concept through supervised learning. Unsupervised learning, on the other hand, involves a student learning the same concept independently and without the guidance of a teacher at home. In the context of semi-supervised learning, the student reviews and reinforces the knowledge acquired after receiving

initial training from the college instructor. The student gains from both supervised and unsupervised learning experiences in this semi-supervised learning scenario.

### **Advantages of Semi-Supervised Learning**

- The algorithm is straightforward and easy to comprehend.
- It demonstrates high efficiency.
- It addresses the limitations of both supervised and unsupervised learning algorithms.

### **Disadvantages**

- The results obtained through iterations may lack stability.
- These algorithms are not applicable to network-level data.
- The accuracy achieved by these algorithms tends to be low.

4. **Reinforcement learning:** In a feedback-driven process known as reinforcement learning, an AI component dynamically evaluates its environment through trial and error. It acts, gains knowledge from its experiences, and works to enhance its performance. The goal of this learning strategy is to maximize the accumulated rewards by rewarding the component for positive actions and punishing it for bad ones.

Reinforcement learning, as opposed to supervised learning, relies solely on experiences for agent learning rather than labeled data. Take video games as an

example. The environment is created by the game, and each action the reinforcement agent takes determines its state. Rewards and penalties are used as feedback for the agent, which affects the final game score. Making the best decisions will help the agent ultimately reach a high score.

### **Categories of Reinforcement Learning**

- **Positive reinforcement learning:** Following a specific behavior displayed by the agent, a reinforcing stimulus is introduced to aid in positive reinforcement learning. This makes it more likely that the behavior will continue in the future. Giving a reward or incentive as a result of the behavior is one example.
- **Negative reinforcement learning:** The goal of negative reinforcement learning is to strengthen a particular behavior that aids in preventing a negative outcome. In this instance, emphasizing behaviors that prevent or eliminate undesirable outcomes is key.

### **Advantages of Reinforcement learning**

- It tackles complex real-world problems that are challenging to solve using conventional techniques.
- The learning model of reinforcement learning closely resembles human learning, resulting in highly accurate outcomes.

- It facilitates the attainment of long-term results

### **Disadvantages**

- RL algorithms are not well-suited for simple problems, as their complexity may outweigh their benefits.
- RL algorithms demand substantial data and computational resources for effective training.
- Excessive reinforcement learning can lead to an overwhelming number of states, potentially diminishing

### **Applications of Reinforcement Learning**

- Video Games
- Resource Management
- Robotics
- Text Mining

## **2.3 MATHEMATICAL MODELS FOR MACHINE LEARNING**

Mathematical models are at the core of machine learning algorithms. These models capture patterns and relationships within data to make predictions or decisions. Here are some common mathematical models and techniques used in machine learning:

1. **Linear Regression:** Linear regression models the relationship between a dependent variable and one or more independent variables by fitting a linear equation to the observed data. The model can be expressed as:

$Y = aX + b$  where  $Y$  is the dependent variable,  $X$  is the independent variable,  $a$  is the slope, and  $b$  is the intercept.

2. **Logistic Regression:** Logistic regression is used for binary classification problems. It models the probability of a binary outcome using the logistic function. The model can be expressed as:

$$P(Y = 1|X) = \frac{1}{1+e^{-(aX+b)}}$$

3. **Decision Trees:** Decision trees are hierarchical models used for both classification and regression tasks. They partition the data into subsets based on the values of input features and make decisions based on these partitions.

4. **Random Forests:** Random forests are an ensemble technique that combines multiple decision trees to improve accuracy and reduce overfitting.

5. **Support Vector Machines (SVM):** SVM is a supervised learning model that finds a hyperplane in a high-dimensional space to best separate classes. It is often used for classification.

6. **Neural Networks:** Neural networks are a class of models inspired by the structure of the human brain. They consist of layers of interconnected nodes (neurons) and are used for a wide range of machine learning tasks, including image recognition and natural language processing.
7. **Principal Component Analysis (PCA):** PCA is a dimensionality reduction technique that uses linear algebra to transform high-dimensional data into a lower-dimensional space while preserving as much variance as possible.
8. **Naive Bayes:** Naive Bayes is a probabilistic classifier based on Bayes' theorem. It assumes that features are conditionally independent, which simplifies the calculations.
9. **Gradient Boosting:** Gradient boosting is an ensemble method that combines weak learners (typically decision trees) to create a strong learner. Popular implementations include Gradient Boosting Machines (GBM), XGBoost, and LightGBM.
10. **Hidden Markov Models (HMMs):** HMMs are used for modeling sequential data, such as time series or speech recognition. They involve probabilistic state transitions and emissions.
11. **Recurrent Neural Networks (RNNs):** RNNs are a type of neural network designed for sequential data. They have loops to allow information to persist

through time steps, making them suitable for tasks like natural language processing and time series prediction.

**12. Long Short-Term Memory (LSTM) Networks:** LSTMs are a specific type of RNN that are capable of capturing long-range dependencies in sequential data.

**13. Convolutional Neural Networks (CNNs):** CNNs are specialized neural networks designed for processing grid-like data, such as images and audio spectrograms. They use convolutional layers to automatically learn hierarchical features.

**14. Reinforcement Learning Models:** These models, such as Q-learning and policy gradients, are used in reinforcement learning to make sequential decisions in environments with rewards and punishments.

## **2.4 IMPORTANCE OF MACHINE LEARNING**

One of the most well-known branches of artificial intelligence (AI) is machine learning. Nearly every industry, including healthcare, finance, infrastructure, marketing, self-driving cars, recommendation systems, chatbots, social media, gaming, cyber security, and many others, uses machine learning concepts.

Due to its capacity to analyze vast amounts of data, generate precise predictions, and automate processes, machine learning is extremely important. Through

recommendation systems, it enables data-driven decision-making, boosts productivity, and offers customized experiences.

Aside from its applications in fraud detection, cybersecurity, healthcare, and optimization, machine learning has had a significant impact on industries through advances in image and speech recognition. In general, machine learning enables organizations to use data for insights, automation, and innovation across a variety of domains.

## **2.5 LIMITATIONS OF MACHINE LEARNING**

- 1. Data dependency:** Inadequate or biased data can produce inaccurate results because machine learning heavily relies on data. The reliability, representativeness, and availability of data can present problems and limit how well ML models work.
- 2. Interpretability and transparency:** Deep learning neural networks are one example of an ML algorithm that can be complicated and challenging to understand. Particularly in areas where interpretability is essential, lack of transparency in model predictions and decision-making processes can impede understanding, trust, and acceptance.
- 3. Overfitting and generalization:** ML models can perform poorly on untrained data if they become overly focused on the training set. Overfitting can reduce the model's capacity to handle novel, unanticipated scenarios,

making it difficult to strike a balance between model complexity and generalization.

4. **Computational resources and scalability:** Deep learning algorithms in particular call for a lot of memory and processing power. Scaling up ML solutions can be computationally expensive, posing practical constraints in terms of cost, infrastructure, and scalability.
5. **Ethical and fairness considerations:** Due to biases in the training data, ML models are vulnerable to producing unfair or discriminatory results. In order to ensure fairness, accountability, and transparency and prevent the reinforcement of biases or the continuation of discrimination, it is essential to address ethical considerations.
6. **Domain-specific expertise:** Domain knowledge is frequently necessary for applying ML effectively. The complexity of the data, feature engineering, and choosing the right evaluation metrics can be difficult, necessitating cooperation between domain specialists and ML practitioners.
7. **Continuous learning and adaptability:** ML models may have trouble handling situations that are unrelated to their training data or adapting to new environments. It can be difficult to adjust models to changing data or concept drift; this necessitates ongoing monitoring, retraining, and updates.

8. **Lack of casual understanding:** Instead of establishing causality, ML focuses on finding correlations and patterns. The ability of ML models to understand cause-and-effect relationships may be limited, which limits the range of applications where they can be used.

## **CONCLUSIONS**

In conclusion, mathematical models play a significant role in machine learning's ability to analyze data, predict outcomes, and identify patterns. These models serve as the foundation for many machine learning (ML) algorithms, which facilitate the identification of patterns and relationships in datasets for the benefit of national development.

## **CHAPTER THREE**

### **METHODOLOGY**

#### **3.1 INTRODUCTION**

Machine learning is built on mathematical models because they provide a structured method for comprehending and resolving a wide range of issues in the industry. These models, which range from straightforward linear equations to intricate neural networks and probabilistic systems, represent patterns and relationships in data. They are essential for extending learned patterns to forecast future data that hasn't yet been observed. Models are trained using mathematical methods, with the aim of reducing error by optimizing the parameters. Mathematical metrics and criteria are used in the evaluation of these models to judge how well they perform. Machine learning models frequently use concepts from many different mathematical disciplines, including linear algebra, calculus, statistics, and optimization. As the field develops, new mathematical models and methods keep coming into existence to address challenging data-driven problems.

#### **3.2 MATHEMATICAL MODELS IN MACHINE LEARNING**

##### **1. Linear models**

- **Linear Regression:** models a linear relationship between a dependent variable (**Y**) and an independent variable (**X**).

Simple linear regression:  $Y = \mathbf{m}x + \mathbf{b}$

**Y:** this is the dependent variable or the target variable

**x:** this is the independent variable or the input feature

**m:** this is the slope of the line

**b:** this is the y-intercept.

In order to best fit the line to the provided data points, the coefficients (**m** and **b**) in machine learning are typically estimated using data. This is accomplished by reducing the sum of squared differences between the target values (**y**) and the predicted values (**mx + b**).

- **Logistic Regression:** Uses the logistic function to model the likelihood of binary classification outcomes based on one or more independent variables, ensuring that the predicted values are between 0 and 1.

Logistic regression (Binary Classification):

$$P\left(Y = \frac{1}{X}\right) = \frac{1}{1 + e^{-(b_0 + b_1x_1 + b_2x_2 + \dots + b_nx_n)}}$$

- **P(Y=1|X):** This is the probability that the binary outcome variable **Y** is equal to **1** given the input features **X**.

- **X**: This is a vector of input features ( $x_1, x_2, \dots, x_n$ ).
- **b<sub>0</sub>**: This is the intercept term (also called the bias term) of the logistic regression model.
- **b<sub>1</sub>, b<sub>2</sub>, ..., b<sub>n</sub>**: These are the coefficients (weights) associated with each input feature. They represent how much each feature contributes to the probability of the positive class ( $Y=1$ ).
- **e**: This is the base of the natural logarithm, approximately equal to 2.71828.

The logistic function  $\frac{1}{1+e^{-(b_0+b_1x_1+b_2x_2+\dots+b_nx_n)}}$  maps the linear combination of input features and coefficients into a probability value between 0 and 1. Using techniques like maximum likelihood estimation, logistic regression frequently estimates the values of the coefficients **b<sub>1</sub>, b<sub>2</sub>, ..., b<sub>n</sub>** based on the training data.

The predicted probability of making a binary decision is applied during classification using a threshold, which is typically 0.5. The class is designated as 1 if the predicted probability exceeds or equals the threshold, and 0 otherwise.

## 2. Tree-Based Models

- **Decision Tree**: This non-linear model is employed for both regression and classification tasks. To create homogeneous subgroups with respect to the target variable, they are built by recursively partitioning the data into subsets

based on feature values. Unlike logistic or linear regression, the mathematical expression for decision trees is not frequently expressed as a single equation. Instead, a set of rules or conditions are used to create decision trees.

How decision tree works:

- **Splitting Nodes:** At each internal node of the tree, a decision is made based on one of the input features. This decision is represented as a rule or condition. For example, "if feature  $X_i$  is less than a threshold value, go left; otherwise, go right."
- **Leaf Nodes:** The leaves of the decision tree represent the final predicted values for regression tasks or the class labels for classification tasks. Each leaf node corresponds to a specific outcome.
- **Training:** The decision tree is trained using a dataset by recursively splitting the data at each node based on the chosen conditions. The goal is to maximize homogeneity within each leaf node.
- **Prediction:** To make predictions on new data, you traverse the tree from the root node to a leaf node, following the conditions at each node until you reach a leaf. The value or class associated with that leaf is the prediction.

Here's a simplified representation of a decision tree for a binary classification task:

if  $X_i < \text{threshold}$ :

|

|--> Class A

else:

|

|--> Class B

The above representation shows a single decision node and two leaf nodes, one for each class (A and B). Depending on the dataset and the number of features, decision trees can be much more complex with multiple levels of nodes.

- **Random Forests:** This is an ensemble of decision trees that improve predictive accuracy and reduce overfitting.

Simplified representation of random forests:

1. **Ensemble of Decision Trees:** A Random Forest consists of a collection of individual decision trees.
2. **Bootstrapped Samples:** Each decision tree in the Random Forest is trained on a bootstrapped (randomly sampled with replacement) subset of the training data. This introduces variability into the training process.
3. **Feature Randomization:** At each node of each decision tree, a random subset of features is considered for splitting. This introduces diversity among the individual trees.

4. **Voting or Averaging:** For classification tasks, the final prediction is often determined by a majority vote among the decision trees (i.e., the class that appears most frequently). For regression tasks, the final prediction is typically the average of the predictions from individual trees.

The key idea behind Random Forest is that by combining multiple decision trees, each trained on different subsets of data and features, you can reduce overfitting and improve generalization performance. The randomness in the bootstrapped samples and feature selection helps to decorrelate the individual trees.

So, while there isn't a single mathematical expression that represents Random Forest as a whole, it can be thought of as an ensemble learning technique that leverages decision trees with bootstrapping and feature randomization to make predictions more robust and accurate.

### 3. Kernel Methods

- **Support Machine Vectors (SVM):** Uses mathematical kernels (e.g. linear, polynomial, Gaussian) to find optimal hyperplanes for classification and fits the data well in the case of regression.

Support Vector Machine (Linear):  $f(x) = w \cdot x + b$

- **f(x)**: This is the decision function that represents the output of the SVM. It's the linear combination of the input features **x**, weighted by the vector **w**, and with an added bias term **b**.
- **w**: This is the weight vector, which is perpendicular to the hyperplane and determines the orientation of the hyperplane.
- **x**: This is the input feature vector.
- **b**: This is the bias term or intercept, which shifts the hyperplane away from the origin.

The decision boundary of the SVM is the hyperplane defined by the equation  $\mathbf{f}(\mathbf{x})=0$ . In a binary classification problem, the SVM aims to find the hyperplane that maximizes the margin between the two classes, subject to the constraint that all data points are correctly classified. The support vectors are the data points that are closest to the decision boundary, and they are the ones that have a non-zero margin.

In practice, SVMs often involve solving optimization problems to find the optimal values of **w** and **b** while satisfying the margin and classification constraints. The most common formulation for the optimization problem in a linear SVM is the quadratic programming problem, which aims to maximize the margin while minimizing the classification error.

The expression becomes more involved when you consider non-linear SVMs, where a kernel function is used to map the data into a higher-dimensional space, making it linearly separable. However, the basic principles of finding a hyperplane to separate data with support vectors remain the same.

- **Kernel PCA:** Applies kernel tricks to Principal Component Analysis for non-linear dimensionality reduction by using kernel functions.

The mathematical expression for Kernel PCA can be a bit complex due to the use of kernel functions, but here's a high-level overview of how it works:

1. **Kernel Function:** Kernel PCA uses a kernel function, denoted as  $\mathbf{K}(\mathbf{x}, \mathbf{x}')$ , which computes the similarity between two data points  $\mathbf{x}$  and  $\mathbf{x}'$  in the original feature space. Common kernel functions include the Gaussian (RBF) kernel, polynomial kernel, and sigmoid kernel.
2. **Gram Matrix:** A Gram matrix (also called a kernel matrix), denoted as  $\mathbf{K}$ , is constructed from the data points. Each entry  $\mathbf{K}(\mathbf{i}, \mathbf{j})$  of the Gram matrix represents the similarity between data points  $\mathbf{x}_i$  and  $\mathbf{x}_j$  using the kernel function:  $\mathbf{K}(\mathbf{i}, \mathbf{j}) = \mathbf{K}(\mathbf{x}_i, \mathbf{x}_j)$
3. **Centering:** Center the Gram matrix  $\mathbf{K}$  by subtracting the mean of each row, column, and the entire matrix to obtain a centered Gram matrix  $\mathbf{K}_c$ .

4. **Eigenvalue Decomposition:** Perform eigenvalue decomposition (eigendecomposition) on the centered Gram matrix  $\mathbf{K}_c$  to obtain its eigenvectors (principal components) and eigenvalues.
5. **Projection:** The top  $k$  eigenvectors corresponding to the largest eigenvalues are selected, and data points are projected into the lower-dimensional space spanned by these eigenvectors.

Simplified mathematical expression for the projection step in Kernel PCA:

**Projection into Lower-Dimensional Space:**

$$\mathbf{z}_i = \begin{bmatrix} \phi(\mathbf{x}_i) \cdot \mathbf{V}_1 \\ \phi(\mathbf{x}_i) \cdot \mathbf{V}_2 \\ \vdots \\ \phi(\mathbf{x}_i) \cdot \mathbf{V}_k \end{bmatrix}$$

- $\mathbf{z}_i$ : The projected representation of data point  $\mathbf{x}_i$  in the lower-dimensional space.
- $\phi(\mathbf{x}_i)$ : The mapping of data point  $\mathbf{x}_i$  into the higher-dimensional feature space using the kernel function.
- $\mathbf{V}_1, \mathbf{V}_2, \dots, \mathbf{V}_k$ : The top  $k$  eigenvectors (principal components) obtained from the eigendecomposition of the centered Gram matrix.

In practice, the computational complexity of Kernel PCA can be high, especially when dealing with large datasets, as it involves the computation of the Gram matrix and eigendecomposition. Common kernels include the Radial Basis Function (RBF) kernel and polynomial kernel.

#### 4. Neural Networks

➤ **Artificial Neural Networks (ANN):** these are composed of layers of interconnected neurons or nodes that perform linear combination and nonlinear transformations. Here's a mathematical expression for a simplified feedforward neural network with one hidden layer:

- **Input Layer:** The input layer consists of neurons that represent the input features. If you have 'n' input features, you'll have 'n' neurons in this layer. The output of the input layer is the input itself.

$$\text{Input: } \mathbf{x} = [x_1, x_2, \dots, x_n]$$

- **Hidden Layer:** The hidden layer(s) perform transformations on the input data using weights and biases. In this example, we have one hidden layer with 'm' neurons.

For the  $i_{th}$  neuron in the hidden layer, the activation  $\mathbf{a}_i^{(1)}$  is computed as:

$$\mathbf{a}_i^{(1)} = \sigma(\mathbf{w}_i^{(1)} \cdot \mathbf{x} + \mathbf{b}_i^{(1)})$$

Here:

- $\mathbf{w}_i^{(1)}$  represents the weights for the connections from the input layer to the hidden layer for the  $i_{th}$  neuron.
- $\mathbf{b}_i^{(1)}$  represents the bias term for the  $i_{th}$  neuron.
- $\sigma$  is an activation function (commonly used functions include the sigmoid, ReLU, or tan  $h$  functions).
- **Output Layer:** The output layer takes the activations from the hidden layer and produces the final output. In this case, we assume a binary classification problem, so we have a single neuron in the output layer.

For the output neuron, the activation  $\mathbf{a}^2$  is computed as:

$$\mathbf{a}^2 = \sigma(\mathbf{w}^2 \cdot \mathbf{a}^1 + \mathbf{b}^2)$$

Here:

- $\mathbf{w}^2$  represents the weights for the connections from the hidden layer to the output layer.
- $\mathbf{b}^2$  represents the bias term for the output neuron.

- $\sigma$  is again the activation function, which can be different from the one used in the hidden layer.

**Output:** The final output of the neural network is the activation of the output neuron: Output:  $a^2$

Training a neural network involves adjusting the weights ( $\mathbf{w}$ ) and biases ( $\mathbf{b}$ ) through a process known as backpropagation and minimizing a loss function (e.g., mean squared error or cross-entropy) to make the network's predictions more accurate.

- **Convolutional Neural Networks (CNN):** Designed for spatial or grid-like data like images and video using convolutional layers. Mathematical expression for a basic convolutional layer in a CNN:

Given an input feature map or image, the convolutional layer performs a convolution operation with a set of learnable filters (also called kernels) to produce feature maps that capture local patterns or features.

- **Input Feature Map:** Let's denote the input feature map as  $\mathbf{I}$ .

$\mathbf{I}$  represents a 3D tensor with dimensions:

Width:  $W_{in}$

Height:  $H_{in}$

Number of Channels (e.g., RGB channels):  $C_{in}$

- **Filters (Kernels):** A convolutional layer consists of multiple filters. Let's denote the set of filters as  $\mathbf{K}$ .

$\mathbf{K}$  is a collection of 3D tensors with dimensions:

Filter Width:  $F_w$

Filter Height:  $F_h$

Number of Channels (must match  $C_{in}$ ):  $C_{in}$

- **Convolution Operation:** For each filter  $k$  in  $\mathbf{K}$ , the convolution operation is performed, resulting in a feature map or activation map  $A_k$ .

$A_k$  is computed by sliding the filter  $k$  over the input feature map  $\mathbf{I}$ , performing element-wise multiplications and summing the results for each position.

$$A_k(i, j) = \sum_{m=0}^{F_h-1} \sum_{n=0}^{F_w-1} \sum_{c=0}^{C_{in}-1} K(i-m, j-n, c) \cdot I(m, n, c)$$

$(i, j)$  represents the spatial position in the output feature map.

$(m, n)$  represents the spatial position in the filter.

$c$  indexes the channels.

- **Activation Function:** Typically, an activation function (e.g., ReLU) is applied element-wise to each value in the feature map  $A_k$ .

$A'_k$  is the activated feature map:

$$A'_k(i, j) = \text{ReLU}(A_k(i, j))$$

- **Output Feature Maps:** The convolutional layer generates multiple output feature maps, one for each filter in  $\mathbf{K}$ .

The set of output feature maps is represented as  $\mathbf{A}$ :

$$\mathbf{A} = A'_1, A'_2, \dots, A'_n$$

- **Recurrent Neural Networks (CNN):** These are a class of neural networks designed for processing sequences of data, making them particularly useful for tasks involving sequential or time-series data. RNNs have recurrent connections that allow them to maintain and update a hidden state as they process each element of a sequence. Here's a mathematical expression for a simplified recurrent layer in an RNN:

- **Recurrent Layer:** Given a sequence of input vectors  $\mathbf{X} = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_T\}$  an RNN processes each input element sequentially and maintains a hidden state  $\mathbf{h}_t$  at each time step  $t$ .

**Input Sequence:** Let's denote the input sequence as  $\mathbf{X}$ .

$$\mathbf{X} = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_T\}$$

Each  $\mathbf{x}_t$  represents the input vector at time step  $t$ .

**Hidden State:** At each time step  $t$ , the RNN maintains a hidden state vector  $\mathbf{h}_t$ .

$\mathbf{h}_t$  is computed based on the current input  $\mathbf{x}_t$  and the previous hidden state  $\mathbf{h}_{\{t-1\}}$  using a recurrent function:

$$\mathbf{h}_t = f(W \cdot \mathbf{x}_t + U \cdot \mathbf{h}_{t-1})$$

Here:

$f()$  is an activation function, often the hyperbolic tangent (tanh) or rectified linear unit (**ReLU**).

$W$  and  $U$  are weight matrices that are learned during training.

$\cdot$  represents matrix-vector multiplication.

**Output:** At each time step, the RNN may produce an output vector  $\mathbf{y}_t$  based on the hidden state  $\mathbf{h}_t$ .

$\mathbf{y}_t$  can be computed using another set of weights and an activation function:

$$\mathbf{y}_t = g(V \cdot \mathbf{h}_t)$$

Here:

$g()$  is an activation function that is often used depending on the problem (e.g., softmax for classification).

The RNN processes the entire sequence  $\mathbf{X}$  by sequentially updating the hidden state at each time step and optionally producing outputs. The final hidden state can  $\mathbf{h}_T$  be used for various purposes, such as making predictions or capturing information from the entire sequence.

➤ **Long Short-Term Memory (LSTM) and Gated Recurrent Unit (GRU):**

These are advanced variants of recurrent neural networks (RNNs) designed to address the vanishing gradient problem and better capture long-term dependencies in sequential data. Here are mathematical expressions for LSTM and GRU:

**Long Short-Term Memory (LSTM):** LSTM introduces the concept of cell state, which allows it to selectively add or remove information from the cell state using three main gates: the forget gate, input gate, and output gate.

**Forget Gate:** The forget gate determines what information from the previous cell state  $\mathbf{C}_{\{t-1\}}$  should be thrown away or kept. It produces a forget gate activation  $\mathbf{f}_t$  based on the current input  $\mathbf{x}_t$  and the previous hidden state  $\mathbf{h}_{\{t-1\}}$ .

$$\mathbf{f}_t = \sigma(\mathbf{W}_f \cdot [\mathbf{h}_{t-1}, \mathbf{x}_t] + \mathbf{b}_f)$$

**Input Gate:** The input gate decides what new information should be added to the cell state. It produces an input gate activation ( $i_t$ ) and a candidate cell state  $C_{\{t-1\}}$  based on the current input and the previous hidden state.

$$i_t = \sigma(W_i \cdot [h_{t-1}, x_t] + b_i)$$

$$C_{tilde}_t = \tanh(W_c \cdot [h_{t-1}, x_t] + b_c)$$

**Update Cell State:** The cell state ( $C_t$ ) is updated by adding the information that the input gate decided to keep and removing what the forget gate decided to forget.

$$C_t = f_t \cdot C_{t-1} + i_t \cdot C_{tilde}_t$$

**Output Gate:** The output gate determines the hidden state output based on the current input and the updated cell state. It produces an output gate activation ( $o_t$ ) and the hidden state ( $h_t$ ).

$$o_t = \sigma(W_o \cdot [h_{t-1}, x_t] + b_o)$$

$$h_t = o_t \cdot \tanh(C_t)$$

**Gated Recurrent Unit (GRU):** GRU simplifies the architecture compared to LSTM by combining the cell state and hidden state into a single hidden state, using two main gates: the reset gate and the update gate.

**Update Gate:** The update gate ( $z_t$ ) determines how much of the previous hidden state should be retained and how much of the new candidate state should be used to update the current hidden state.

$$z_t = \sigma(W_z \cdot [h_{t-1}, x_t] + b_z)$$

**Reset Gate:** The reset gate ( $r_t$ ) determines how much of the previous hidden state should be "forgotten" and not used in the current hidden state calculation.

$$r_t = \sigma(W_r \cdot [h_{t-1}, x_t] + b_r)$$

**Candidate Hidden State:** The candidate hidden state ( $h_{t-1}$ ) is computed based on the current input and the reset gate, allowing the network to decide which information from the previous hidden state should be kept.

$$h_{tilde}_t = \tanh(W_h \cdot [r_t \cdot h_{t-1}, x_t] + b_h)$$

**Update Hidden State:** The hidden state  $h_t$  is updated by interpolating between the candidate hidden state and the previous hidden state based on the update gate.

$$h_t = (1 - z_t) \cdot h_{t-1} + z_t \cdot h_{tilde}_t$$

Both LSTM and GRU architectures allow for better gradient flow during training and have shown improved performance on tasks involving sequential data compared to basic RNNs.

## 5. Clustering Models

➤ **K-Means:** This is an unsupervised machine learning algorithm used for clustering data points into a specified number of clusters. The algorithm aims to minimize the sum of squared distances between data points and their assigned cluster centroids. Here's a mathematical expression for the K-Means algorithm:

Given a dataset with  $n$  data points  $\{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n\}$  and a predetermined number of clusters  $K$ , the goal of K-Means is to find  $K$  cluster centroids  $\{\mathbf{c}_1, \mathbf{c}_2, \dots, \mathbf{c}_k\}$  and assign each data point to one of these clusters.

**Initialization:** Initialize the cluster centroids randomly or using a specific initialization method.

**Assignment Step (Expectation):** For each data point  $\mathbf{x}_i$ , calculate its distance to each of the  $K$  cluster centroids and assign it to the cluster associated with the nearest centroid. The assignment is typically based on the Euclidean distance.

$$\text{minimize } \{S_1, S_2, \dots, S_k\} \sum_{k=1}^k \sum_{x \in S_k} \|x - C_k\|^2$$

$S_k$  represents the set of data points assigned to cluster  $k$ .

$C_k$  is the centroid of cluster  $k$ .

The goal is to minimize the sum of squared distances within each cluster.

**Update Step (Maximization):** Recalculate the cluster centroids by taking the mean of all data points assigned to each cluster:

$$C_k = \frac{1}{S_k} \sum_{x \in S_k} x$$

**Convergence Check:** Repeat the assignment and update steps iteratively until convergence. Convergence is typically determined by a predefined criterion, such as a maximum number of iterations or a small change in cluster assignments.

**Output:** Once the algorithm converges, the final cluster centroids and cluster assignments are obtained.

The above steps describe the basic K-Means algorithm, which is commonly used for clustering data. K-Means is sensitive to the initial placement of centroids, and different initializations can lead to different solutions. To mitigate this, practitioners often use techniques like K-Means++ initialization, which selects initial centroids in a way that improves the quality of the resulting clusters.

- **Hierarchical Clustering:** This is a family of clustering algorithms that build a hierarchy of clusters by successively merging or splitting groups of data points. Two common approaches to hierarchical clustering are agglomerative (bottom-up) and divisive (top-down). Here's an expression for agglomerative hierarchical clustering:

**Agglomerative Hierarchical Clustering:** Given a dataset with  $n$  data points  $\{x_1, x_2, \dots, x_n\}$ , the goal of agglomerative hierarchical clustering is to create a hierarchy of clusters by iteratively merging individual data points or existing clusters until a single cluster containing all data points is formed. The process involves the following steps:

**Initialization:** Start with each data point as its own cluster, resulting in  $n$  initial clusters, each containing a single data point.

**Merge Step (Agglomeration):** At each iteration, merge the two closest clusters into a single cluster. The closeness between clusters is typically defined using a linkage criterion, such as single linkage, complete linkage, or average linkage. The linkage criterion determines the distance between clusters and guides the merging process.

**Single Linkage:** The distance between two clusters is defined as the shortest distance between any pair of data points, one from each cluster.

**Complete Linkage:** The distance between two clusters is defined as the longest distance between any pair of data points, one from each cluster.

**Average Linkage:** The distance between two clusters is defined as the average distance between all pairs of data points, one from each cluster. The merging continues until there is only one cluster left.

**Hierarchical Structure:** As clusters are merged, a hierarchical structure (dendrogram) is constructed, where the root of the tree represents a single cluster containing all data points, and the leaves represent the individual data points.

The hierarchical structure can be visualized as a dendrogram, which provides a graphical representation of the clustering process, showing how clusters are merged step by step. The height at which clusters are merged in the dendrogram reflects the distance or dissimilarity between the merged clusters.

Agglomerative hierarchical clustering does not require specifying the number of clusters beforehand, which can be a benefit. However, you can decide to cut the dendrogram at a certain height to obtain a specific number of clusters, if desired.

The mathematical expressions for the distances and linkage criteria used in agglomerative hierarchical clustering can vary depending on the specific algorithm and implementation, but the overall process involves iteratively merging clusters based on distance or similarity measures.

## **6. Dimensionality Reduction**

- **Principal Component Analysis:** This is a dimensionality reduction technique used for finding the orthogonal axes (principal components) along which the data varies the most. Here's a mathematical expression for PCA:

Given a dataset with  $n$  data points  $\{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n\}$ , where each  $\mathbf{x}_i$  is a  $p$ -dimensional vector, the goal of PCA is to find a set of  $k$  principal components (orthogonal unit vectors  $\{\mathbf{V}_1, \mathbf{V}_2, \dots, \mathbf{V}_k\}$ , such that the variance of the data along these components is maximized.

**Mean-Centering:** Calculate the mean vector  $\boldsymbol{\mu}$  of the data:

$$\boldsymbol{\mu} = \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i$$

Subtract the mean vector from each data point:

$$\mathbf{x}_i = \mathbf{x}_i - \boldsymbol{\mu}$$

**Covariance Matrix:**

Compute the  $p \times p$  covariance matrix  $\boldsymbol{\Sigma}$ :

$$\boldsymbol{\Sigma} = \frac{1}{n-1} \sum_{i=1}^n (\mathbf{x}_i)(\mathbf{x}_i)^T$$

$\boldsymbol{\Sigma}$  measures how the features of the data are correlated with each other.

**Eigenvalue Decomposition:**

Perform eigenvalue decomposition on the covariance matrix  $\boldsymbol{\Sigma}$  to obtain the eigenvalues  $\lambda_1, \lambda_2, \dots, \lambda_p$ , and corresponding eigenvectors  $\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_p$ :

$$\sum V_i = \lambda_i V_i$$

$\lambda_i$  represents the eigenvalues, which indicate the amount of variance explained by each principal component.

$v_i$  are the eigenvectors, which are the principal components.

**Select Principal Components:** Sort the eigenvalues in descending order and select the top  $k$  eigenvectors (principal components) corresponding to the largest eigenvalues. These  $k$  eigenvectors form a new  $p \times k$  matrix  $V_k$ .

**Project Data onto Principal Components:** Project the mean-centered data points onto the selected principal components to obtain the reduced-dimensional representation:

$$Z_i = X_i^T V_k$$

Each  $z_i$  is a  $k$ -dimensional vector representing the data point in the reduced space.

PCA reduces the dimensionality of the data while retaining as much of the original variance as possible. The selected principal components represent the directions in the data space along which the data varies the most. The amount of variance retained in the reduced space is determined by the eigenvalues, with larger eigenvalues indicating greater variance along the corresponding principal component.

- **t-Distributed Stochastic Neighbor Embedding (t-SNE):** t-Distributed Stochastic Neighbor Embedding (t-SNE) is a nonlinear dimensionality reduction technique used for visualizing high-dimensional data in lower-dimensional spaces, typically 2D or 3D. t-SNE aims to preserve pairwise similarities between data points by transforming them into a probability distribution over pairs of points in the lower-dimensional space. Here's an overview of the mathematical expression for t-SNE:

Given a dataset with  $n$  data points  $\{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n\}$  in a high-dimensional space and a desired lower-dimensional space with  $d$  dimensions (usually  $d = 2$  or  $3$ ), the goal of t-SNE is to find a set of lower-dimensional data points  $\{\mathbf{y}_1, \mathbf{y}_2, \dots, \mathbf{y}_n\}$  that best represents the pairwise similarities between the data points in the high-dimensional space.

### **Pairwise Similarities in the High-Dimensional Space:**

Calculate pairwise similarities between data points in the high-dimensional space.

The similarity between data points  $\mathbf{x}_i$  and  $\mathbf{x}_j$  is typically computed as a conditional probability:

$$P_{j/i} = \frac{\exp(-\|\mathbf{x}_i - \mathbf{x}_j\|^2 / 2\sigma_i^2)}{\sum_{k=j} \exp(-\|\mathbf{x}_i - \mathbf{x}_k\|^2 / 2\sigma_i^2)}$$

The similarity  $P_{j/i}$  represents the conditional probability of choosing data point  $x_j$  as a neighbor of  $x_i$  given a Gaussian distribution centered at  $x_i$  with variance  $\sigma_i^2$ . The value of  $\sigma_i^2$  is determined based on the data and helps control the neighborhood size for each data point.

### Pairwise Similarities in the Lower-Dimensional Space:

Calculate pairwise similarities between data points in the lower-dimensional space. These similarities are computed using a similar conditional probability, but with a different variance:

$$q_{j/i} = \frac{(1 + \|y_i - y_j\|^2)^{-1}}{\sum_{k=j} (1 + \|y_i - y_k\|^2)^{-1}}$$

The similarity  $q_{j/i}$  represents the conditional probability of choosing data point  $y_j$  as a neighbor of  $y_i$  given a Student's t-distribution with one degree of freedom (Cauchy distribution) centered at  $y_i$ .

**Minimize the Kullback-Leibler (KL) Divergence:** The goal of t-SNE is to minimize the KL divergence between the pairwise similarities in the high-dimensional space and the pairwise similarities in the lower-dimensional space:

$$KL(P\|Q) = \sum_i \sum_j P_{j/i} \log \left( \frac{P_{j/i}}{q_{j/i}} \right)$$

**P** represents the conditional probabilities in the high-dimensional space

$$(P = \{p_{j/i}\}).$$

**Q** represents the conditional probabilities in the lower-dimensional space

$$(Q = \{q_{j/i}\}).$$

**Optimization:** Use optimization techniques (e.g., gradient descent) to find the lower-dimensional representations  $\{y_1, y_2, \dots, y_n\}$  that minimize the KL divergence.

t-SNE is effective for visualizing high-dimensional data in a way that preserves local similarities while emphasizing global structure. The optimization process helps distribute dissimilar data points far apart in the lower-dimensional space and keeps similar data points close together.

## 7. Bayesian Models

- **Naïve Bayes:** This is a classification algorithm based on Bayes' theorem and the assumption of independence among features. It's often used for text classification and is particularly effective for tasks like spam detection and document categorization. Mathematical expression for the Naive Bayes classification:

Given a dataset with  $n$  data points  $\{x_1, x_2, \dots, x_n\}$  and their corresponding class labels  $\{y_1, y_2, \dots, y_n\}$ , where each  $x_i$  is a feature vector, and each  $y_i$  is the class

label associated with  $x_i$ , the goal of Naive Bayes is to predict the most likely class label for a new input feature vector  $x_{new}$ .

The expression for Naive Bayes involves calculating the probability of a class label given the input features using Bayes' theorem:

$$P\left(\frac{y_j}{x_{new}}\right) = \frac{P y_j \cdot P(x_{new}|y_j)}{P(x_{new})}$$

$P\left(\frac{y_j}{x_{new}}\right)$ : The probability of class  $y_j$  given the new input vector  $x_{new}$ .

$P(y_j)$ : The prior probability of class  $y_j$ . It represents the likelihood of encountering class  $y_j$  in the dataset.

$P\left(\frac{y_j}{x_{new}}\right)$ : The likelihood probability of observing the input vector  $x_{new}$  given class  $y_j$ . This probability is typically estimated using the assumption of feature independence (hence the "Naive" in Naive Bayes):

$$P(x_{new}|y_j) = P(x_{new} \cdot 1|y_j) \cdot P(x_{new} \cdot 2|y_j) \cdot \dots \cdot P(x_{new} \cdot p|y_j)$$

$p$  is the number of features in the input vector.

$P(x_{new}|y_j)$  represents the probability of observing the **i-th** feature given class  $y_j$ .

These probabilities are estimated from the training data using techniques like Maximum Likelihood Estimation (MLE).

$P(\mathbf{x}_{new})$  : The marginal likelihood or evidence, which is the probability of observing the input vector  $\mathbf{x}_{new}$  across all possible classes. It can be computed as a normalizing factor:

$$P(\mathbf{x}_{new}) = \sum_j P(\mathbf{y}_j) \cdot P(\mathbf{x}_{new}|\mathbf{y}_j)$$

Finally, Naive Bayes assigns the class label  $\mathbf{y}_{new}$  to the new input vector  $\mathbf{x}_{new}$  by selecting the class label that maximizes the posterior probability:

$$\mathbf{y}_{new} = \mathit{argmax}_j P(\mathbf{y}_j|\mathbf{x}_{new})$$

In practice, Naive Bayes models are often used with different assumptions about the distribution of the likelihood probabilities (e.g., Gaussian Naive Bayes for continuous features or Multinomial Naive Bayes for discrete count data) depending on the nature of the data.

- **Bayesian Networks:** Bayesian Networks, also known as Bayes Networks or Bayesian Belief Networks, are graphical models that represent probabilistic relationships among a set of random variables. They are based on Bayesian probability theory and conditional independence assumptions. The mathematical expression for Bayesian Networks involves the joint probability distribution of the variables and conditional probabilities. Here's an overview:

Consider a set of random variables  $\{X_1, X_2, \dots, X_n\}$ , where each variable can take on different values. A Bayesian Network represents the joint probability distribution of these variables as a product of conditional probabilities.

**Directed Acyclic Graph (DAG):** A Bayesian Network is represented as a directed acyclic graph (DAG), where each node in the graph corresponds to a random variable, and directed edges between nodes represent conditional dependencies.

**Conditional Probability Distribution for Each Node:** For each node (random variable)  $X_i$  in the network, specify its conditional probability distribution given its parent nodes (nodes with edges pointing to  $X_i$ ). This conditional probability distribution is represented as:

$$P(X_i | \text{Parents}(X_i))$$

$\text{Parents}(X_i)$  represents the set of parent nodes of  $X_i$ , i.e., the nodes with edges pointing to  $X_i$ .

The specific form of this conditional distribution depends on the nature of the variable (e.g., discrete or continuous) and the type of distribution used (e.g., Gaussian, multinomial).

**Joint Probability Distribution:** The joint probability distribution of all variables in the Bayesian Network can be represented as a product of the conditional probability distributions, following the chain rule of probability:

$$P(X_1, X_2, \dots, X_n) = \prod_{i=1}^n P(X_i | \text{Parents}(X_i))$$

This equation captures how the joint probability of all variables is factorized into conditional probabilities based on the DAG structure.

**Inference:** Bayesian Networks are useful for performing various types of probabilistic inference, such as:

**Marginalization:** Calculating the marginal probability of a specific variable by summing or integrating over all other variables.

**Conditional Probability:** Computing the conditional probability of one variable given evidence about other variables.

**Most Likely Explanation (MLE):** Finding the most likely state of a set of variables given evidence.

**Belief Propagation:** Propagating information through the network to update probabilities.

The specific expressions for conditional probability distributions and joint probability distributions can vary depending on the structure of the Bayesian Network and the types of variables involved. Bayesian Networks are widely used in various applications, including medical diagnosis, natural language processing,

and image recognition, to model and reason about uncertain or probabilistic relationships between variables.

- **Gaussian Mixture Models (GMM):** These are probabilistic models used for modeling and representing complex data distributions by combining multiple Gaussian distributions. GMMs are commonly used for tasks like density estimation, clustering, and data generation. Here's a mathematical expression for a Gaussian Mixture Model:

Consider a dataset consisting of  $\mathbf{n}$  data points  $\{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n\}$  where each  $\mathbf{x}_i$  is a  $\mathbf{d}$ -dimensional vector. A GMM models the probability density function (PDF) of this data as a weighted sum of  $\mathbf{K}$  Gaussian components:

$$p(\mathbf{x}) = \sum_{k=1}^{\mathbf{K}} \pi_k \cdot \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$$

$\mathbf{P}(\mathbf{x})$ : The probability density of data point  $\mathbf{x}$ .

$\mathbf{K}$ : The number of Gaussian components in the mixture.

$\pi_k$ : The mixing coefficient of the  $\mathbf{k}$ -th Gaussian component, representing the probability that a data point belongs to that component. These coefficients satisfy  $\sum \pi_k = 1$ , and  $0 \leq \pi_k \leq 1$ .

$\boldsymbol{\mu}_k$ : The mean vector of the  $\mathbf{k}$ -th Gaussian component, which determines the center of the component.

$\Sigma_k$ : The covariance matrix of the **k-th** Gaussian component, which defines the shape and orientation of the distribution.

$\Sigma_k$  can be a diagonal covariance matrix (diagonal elements represent variances along each dimension) or a full covariance matrix (allowing for correlations between dimensions).

The Gaussian component is represented as:

$$\mathfrak{N}(x|\mu_k, \Sigma_k) = \frac{1}{(2\pi)^{d/2} |\Sigma_k|^{1/2}} \exp \left( -\frac{1}{2} (x - \mu_k)^T \Sigma_k^{-1} (x - \mu_k) \right)$$

**d**: The dimensionality of the data.

$|\Sigma_k|$ : The determinant of the covariance matrix  $\Sigma_k$ .

The goal in GMMs is to estimate the model parameters:  $\{\pi_k, \mu_k, \Sigma_k\}$  for each Gaussian component, as well as the number of components **K**, from the observed data. This is typically done using techniques like Expectation-Maximization (EM) or Maximum Likelihood Estimation (MLE).

Once the GMM is trained, it can be used for various tasks, including density estimation, clustering (assigning data points to components), and generating new data points by sampling from the learned mixture.

## 8. Reinforcement Learning

➤ **Markov Decision Processes (MDP):** This is a mathematical framework used in reinforcement learning and decision-making under uncertainty. It models a decision-making problem as a tuple  $(S, A, P, R)$ , where:

**S (State Space):**  $S$  represents the set of all possible states in the environment. It is a finite or countably infinite set.

**A (Action Space):**  $A$  represents the set of all possible actions an agent can take. It is a finite or countably infinite set.

**P (Transition Probabilities):**  $P$  defines the transition probabilities between states when actions are taken. It is often represented as a function:

$$P(s'|s, a)$$

This function represents the probability of transitioning to state  $s'$  from state  $s$  when action  $a$  is taken.

**R (Reward Function):**  $R$  defines the rewards or penalties associated with state-action pairs. It is often represented as a function:

$$R(s, a, s')$$

This function represents the immediate reward (or penalty) received when transitioning from state  $s$  to state  $s'$  by taking action  $a$ .

The dynamics of an MDP are Markovian, meaning that the future state and reward only depend on the current state and action and are independent of the past history.

The objective in solving an MDP is to find a policy  $\pi$ , which is a mapping from states to actions, that maximizes the expected cumulative reward over time. This is typically formulated as the expected return, also known as the cumulative reward, which is the sum of rewards received over time, often discounted by a factor:

$$G_t = \sum_{k=0}^{\infty} \gamma^k R_{t+k+1}$$

$G_t$ : Expected return at time  $t$ .

$\gamma$  (gamma): Discount factor that trades off immediate rewards against future rewards.

$R_{t+k+1}$ : Reward received at time  $t+k+1$ .

The goal is to find the policy  $\pi$  that maximizes the expected return:

$$\pi^* = \underset{\pi}{\operatorname{argmax}} \mathbf{E} \left[ \sum_{k=0}^{\infty} \gamma^k R_{t+k+1} \mid \pi \right]$$

Solving an MDP involves finding the optimal policy  $\pi^*$  that achieves the maximum expected return. This can be done using various algorithms such as Dynamic Programming, Monte Carlo methods, Temporal Difference Learning, and deep reinforcement learning techniques for complex problems.

- **Q-Learning and Deep Q-Networks (DQN):** Q-Learning and Deep Q-Networks (DQNs) are key components of reinforcement learning algorithms used for solving Markov Decision Processes (MDPs) and related problems. Here, I'll provide mathematical expressions for both Q-Learning and DQNs:

**Q-Learning:** Q-Learning is a model-free reinforcement learning algorithm that estimates the Q-values (action-value function) for each state-action pair. The Q-value represents the expected cumulative reward an agent can achieve by taking a specific action in a given state and following a particular policy. The update rule for Q-Learning is as follows:

- Initialize the Q-values arbitrarily for all state-action pairs:  $Q(s, a)$ .
- Repeat until convergence or a fixed number of iterations:
- Choose an action  $a$  from the current state  $s$  using an exploration strategy (e.g.,  $\epsilon$ -greedy).
- Execute action  $a$  and observe the next state  $s'$  and the immediate reward  $r$ .
- Update the Q-value for the current state-action pair using the Q-learning update rule:

$$Q(s, a) \leftarrow Q(s, a) + \alpha(r + \gamma \max_{a'} Q(s', a') - Q(s, a))$$

$\alpha$  (**alpha**) is the learning rate, controlling the step size for updating Q-values.

$\gamma$  (**gamma**) is the discount factor, representing the weight given to future rewards.

$\max_a Q(s', a')$  is the maximum Q-value for the next state  $s'$ .

- Continue learning and improving the policy until convergence.

Q-Learning converges to the optimal Q-values and allows the agent to choose the best action in each state to maximize the expected cumulative reward.

**Deep Q-Networks (DQNs):** Deep Q-Networks are a variant of Q-Learning that use deep neural networks to approximate the Q-values. The Q-value function is parameterized by a neural network with weights  $\theta$ , and the goal is to find  $\theta$  such that  $Q(s, a; \theta)$  approximates the true Q-values as accurately as possible.

The DQN loss function is defined as:

$$L(\theta) = \mathbf{E}[(Q(s, a; \theta) - (r + \gamma \max_{a'} Q(s', a'; \theta^-)))^2]$$

$\theta$ : The weights of the neural network.

$\theta^-$ : The target network's weights (used for stability).

$Q(s, a; \theta)$ : The estimated Q-value for state  $s$  and action  $a$  with network weights  $\theta$ .

$r$ : The immediate reward received after taking action  $a$  in state  $s$ .

$\gamma$ : The discount factor.

$\max_{a'} Q(s', a'; \theta^-)$ : The maximum Q-value over actions in the next state  $s'$ .

The goal of training DQNs is to minimize this loss function by updating the network weights  $\theta$  using gradient descent or related optimization techniques.

DQNs have been particularly successful in solving complex reinforcement learning problems, including playing Atari games and controlling robotic systems. They leverage the power of deep learning to approximate Q-values for high-dimensional state spaces.

## 9. Ensemble Models

- **Bagging (Bootstrap Aggregating):** Bagging (Bootstrap Aggregating) is an ensemble machine learning technique that improves the stability and accuracy of a predictive model by combining the results of multiple base models, often decision trees. Bagging achieves this by training each base model on a different subset of the training data, obtained through a resampling technique known as bootstrapping. Here's a mathematical expression for the bagging process:

Given a dataset with  $n$  data points  $\{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n\}$  and corresponding labels  $\{\mathbf{y}_1, \mathbf{y}_2, \dots, \mathbf{y}_n\}$ , where each  $\mathbf{x}_i$  is a feature vector and  $\mathbf{y}_i$  is the true label for data point  $\mathbf{x}_i$ :

**Bootstrap Sampling:** For each base model in the ensemble (typically a decision tree), create a dataset by randomly selecting  $n$  data points from the original dataset

with replacement. This means that some data points will be duplicated in the new dataset, and some may be omitted.

Mathematically, for each base model ( $\mathbf{m}$ ) in the ensemble:

- Create a bootstrap sample dataset  $\mathbf{D}_m$  of size  $\mathbf{n}$ .
- Each data point in  $\mathbf{D}_m$  is randomly selected with replacement from the original dataset.

**Base Model Training:** Train a base model (e.g., a decision tree) on each of the bootstrap sample datasets  $\mathbf{D}_m$ . Each base model learns to map input feature vectors to labels based on the bootstrap sample it was trained on.

Mathematically, for each base model ( $\mathbf{m}$ ) in the ensemble:

- Train a model, denoted as  $\mathbf{H}_m$ , on the bootstrap sample dataset  $\mathbf{D}_m$ .

**Aggregation (Voting for Classification):** For classification tasks, combine the predictions of all base models by taking a majority vote. The final prediction for a given data point is the class label that receives the most votes among the base models.

Mathematically, for a classification task:

- Let  $\mathbf{H}_m(\mathbf{x}_i)$  represent the prediction of base model  $\mathbf{m}$  for data point  $\mathbf{x}_i$ .

- The ensemble prediction  $H(x_i)$  for data point  $x_i$  is determined by majority voting:

$$H(x_i) = \mathit{mode}\{H_1(x_i), H_2(x_i), \dots, H_M(x_i)\}$$

where  $M$  is the number of base models.

**Aggregation (Averaging for Regression):** For regression tasks, combine the predictions of all base models by taking their average. The final prediction for a given data point is the average of predictions from all base models.

Mathematically, for a regression task:

- Let  $H_m(x_i)$  represent the prediction of base model  $m$  for data point  $x_i$ .
- The ensemble prediction  $H(x_i)$  for data point  $x_i$  is the average of all base model predictions:

$$H(x_i) = \frac{1}{M} \sum_{m=1}^M H_m(x_i)$$

Bagging helps reduce overfitting and improves the robustness of the model by introducing diversity among the base models. It is a widely used technique in ensemble learning, and it is especially effective when combined with high-variance models like decision trees.

➤ **Boosting:** Boosting is an ensemble learning technique that combines multiple weak learners (typically simple models like decision trees) to create a strong learner. Boosting focuses on iteratively giving more weight to misclassified data points, allowing the ensemble to improve its performance with each iteration. The mathematical model for boosting can be expressed as follows:

Given a dataset with  $n$  data points  $\{x_1, x_2, \dots, x_n\}$  and corresponding labels  $\{y_1, y_2, \dots, y_n\}$ , where each  $x_i$  is a feature vector and  $y_i$  is the true label for data point  $x_i$ :

Initialize a set of weights for the data points:  $w_1 = 1/n, w_2 = 1/n, \dots, w_n = 1/n$ .

These weights indicate the importance of each data point.

For  $T$  iterations (where  $T$  is the number of boosting rounds):

- Train a weak learner (e.g., decision tree) on the dataset with weights  $\{w_1, w_2, \dots, w_n\}$ . This weak learner is referred to as  $h_t(x)$ , where  $t$  represents the current iteration.
- Calculate the weighted error (weighted misclassification rate) of the weak learner:

$$\epsilon_t = \sum_{i=1}^n w_i \cdot \mathbf{I}(h_t(x_i) \neq y_i)$$

$\epsilon_t$ : Weighted error of the weak learner in the **t-th** iteration.

**I**: Indicator function (1 if the condition is true, 0 otherwise).

- Calculate the importance (weight) of the weak learner in the final ensemble:

$$\alpha_t = \frac{1}{2} \ln \left( \frac{1 - \epsilon_t}{\epsilon_t} \right)$$

$\alpha_t$ : Importance weight for the weak learner in the **t-th** iteration.

- Update the weights for the data points:

$$w_{t+1,i} = w_{t,i} \cdot \exp(-\alpha_t \cdot I(h_t(x_i) \neq y_i))$$

$w_{t+1,i}$ : Updated weight for data point **i** in the **(t+1)-th** iteration.

- Normalize the weights so that they sum to 1:

$$w_{t+1,i} = \frac{w_{t+1,i}}{\sum_{i=1}^n w_{t+1,i}}$$

- Combine the weak learners into a strong ensemble model:

$$H(x) = \sum_{t=1}^T \alpha_t h_t(x)$$

$H(x)$ : The final ensemble model.

$\alpha_t$ : Importance weight for weak learner **t**.

$h_t(x)$ : Prediction of weak learner **t**.

Boosting focuses on giving more weight to misclassified data points in each iteration, allowing subsequent weak learners to focus on the examples that previous learners found difficult. The final ensemble combines the predictions of all weak learners with different importance weights to create a strong learner that performs well on the dataset. Popular boosting algorithms include AdaBoost, Gradient Boosting, and XGBoost.

## 10. Matrix Factorization

➤ **Singular Value Decomposition (SVD):** This is a fundamental matrix factorization technique used in linear algebra, data compression, and dimensionality reduction. Given a matrix  $A$ , SVD decomposes it into three matrices:  $U$ ,  $\Sigma$  (Sigma), and  $V^T$  (the transpose of  $V$ ). Here's the mathematical expression for Singular Value Decomposition:

Given a real or complex matrix  $A$  of size  $\mathbf{m} \times \mathbf{n}$ :

$$A = U\Sigma V^T$$

Where:

$U$  ( $m \times m$ ): The left singular vectors or left singular matrix. The columns of  $U$  are orthogonal unit vectors.

$\Sigma$  ( $m \times n$ ): The diagonal matrix of singular values. The singular values are non-negative and appear in decreasing order on the diagonal.

$V^T$  (n x n): The right singular vectors or right singular matrix. The rows of  $V^T$  are orthogonal unit vectors.

The elements of  $\Sigma$  (the singular values) are typically denoted as  $\sigma_1, \sigma_2, \dots, \sigma_r$ , where  $r$  is the rank of the matrix  $A$ . The singular values are non-negative and represent the importance of the corresponding singular vectors in describing the original matrix  $A$ . They are ordered in decreasing order, so  $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_r > 0$ .

Using SVD, one can perform various operations, such as matrix approximation, dimensionality reduction, and solving linear equations. SVD is a powerful tool in linear algebra and has applications in various fields, including machine learning, signal processing, and image compression.

➤ **Non-Negative Matrix Factorization:** This is a dimensionality reduction technique used to factorize a non-negative matrix into two lower-dimensional non-negative matrices. NMF is commonly applied in data analysis, feature extraction, and topic modeling. Mathematical expression for Non-Negative Matrix Factorization:

Given a non-negative matrix  $X$  of size  $m \times n$ , NMF decomposes it into two non-negative matrices  $W$  ( $m \times k$ ) and  $H$  ( $k \times n$ ):

$$X \approx WH$$

Where:

**X** (m x n): The original non-negative matrix that you want to factorize.

**W** (m x k): The non-negative matrix representing the basis vectors (also called dictionary or feature matrix). Each column of **W** is a non-negative basis vector.

**H** (k x n): The non-negative matrix representing the coefficients (also called encoding matrix). Each column of **H** contains the coefficients corresponding to the basis vectors for reconstructing the columns of **X**.

The goal of NMF is to find the matrices **W** and **H** that best approximate **X**. This is typically achieved by minimizing the Frobenius norm (a measure of the difference between matrices) between **X** and the approximate reconstruction **WH**:

$$\mathit{min}_{\mathbf{W}, \mathbf{H}} \|\mathbf{X} - \mathbf{WH}\|_F$$

subject to the constraints that **W** and **H** are non-negative (all elements are greater than or equal to zero).

NMF is often used for tasks like topic modeling, image processing, and data compression, where the non-negativity constraints help ensure meaningful and interpretable factorizations. The choice of the factorization rank **k** determines the dimensionality of the approximation and the level of detail captured in the factorization.

## **CONCLUSION**

Mathematical models serve as the cornerstone of machine learning, offering a well-organized and rigorous framework for both constructing and enhancing predictive models within this rapidly evolving domain for national development.

## **CHAPTER FOUR**

### **MACHINE LEARNING AND ITS APPLICATIONS**

#### **4.1 INTRODUCTION**

Machine learning offers a solution to the problem of repetitive coding for new applications by narrowing down potential issues, with examples like speech recognition, self-driving cars, and search recommendation engines already in widespread use today.

Building computer programs to carry out specific tasks, which, when fed with data, can automatically learn from that data on their own (experience), and thus improve their performance (performance), is the core idea behind machine learning.

Experience makes this performance better. The procedure is iterative. These ML algorithms are used to solve a variety of issues, and the solutions they produce are solid and trustworthy enough to be used as a prototype and relied upon.

## 4.2 APPLICATIONS OF MACHINE LEARNING

1. **Healthcare:** ML plays a pivotal role in healthcare, assisting in disease diagnosis, drug discovery, and treatment optimization. Algorithms can analyze medical images (e.g., X-rays, MRIs) to detect abnormalities, predict disease risk factors, and personalize treatment plans based on patient data.

- **Prediction of medical outcomes: Linear Regression:** In linear regression, you can use a mathematical expression like this to predict a medical outcome (e.g., patient's blood pressure) based on one or more input features (e.g., age, weight):

$$\hat{y} = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_n x_n$$

Where:

$\hat{y}$  is the predicted outcome.

$\beta_0, \beta_1, \beta_2, \dots, \beta_n$  are the model coefficients.

$x_1, x_2, \dots, x_n$  are the input features.

- **Logistic Regression for Disease Classification:** In logistic regression, you can use a mathematical expression to model the probability of a binary event (e.g., disease presence or absence):

$$P(Y = 1) = \frac{1}{1 + e^{-(\beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_n x_n)}}$$

Where:

$P(Y = 1)$  is the probability of the event occurring.

$\beta_0, \beta_1, \beta_2, \dots, \beta_n$  are the model coefficients.

$x_1, x_2, \dots, x_n$  are the input features.

- **Neural Networks for Medical Image Analysis:** There is a complex mathematical expression that describes the layers, activations, and weights of convolutional neural networks (CNNs), which are used in deep learning, particularly for the analysis of medical images. Although it may be quite complicated, this is typically represented as a series of matrix multiplications, convolutions, and non-linear activation functions. For example, a simple mathematical expression for a single-layer perceptron might look like:

$$Z = W \cdot X + b$$

Where:

**Z** is the output.

**W** is the weight matrix.

**X** is the input data.

**b** is the bias term.

- **Support Vector Machine (SVM) for Disease Classification:** The mathematical expression for SVMs involves finding the hyperplane that best separates data points into different classes. For a linear SVM, it might be represented as:

$$f(x) = \left( \sum_{i=1}^n \alpha_i y_i K(x_i, x) + b \right)$$

Where:

$f(x)$  is the classifier's decision function.

$\alpha_i$  are the Lagrange multipliers.

$y_i$  are the class labels.

$K(x_i, x)$  is the kernel function.

$b$  is the bias term.

- **Cox Proportional Hazards Model for Survival Analysis:** In survival analysis, you can use the Cox proportional hazards model to analyze time-to-event data:

$$\lambda(t) = \lambda_0(t) \cdot e^{\beta_1 x_1 + \beta_2 x_2 + \dots + \beta_n x_n}$$

Where:

$\lambda(t)$  is the hazard rate at time  $t$ .

$\lambda_0(t)$  is the baseline hazard rate.

$\beta_1, \beta_2, \dots, \beta_n$  are the coefficients associated with covariates  $x_1, x_2, \dots, x_n$

2. **Finance:** It can analyze vast amounts of financial data to identify unusual patterns and make real-time investment decisions.

- **Linear Regression for Asset Price Prediction:** In linear regression, you can use a mathematical expression to predict the price of a financial asset based on historical data and relevant features:

$$\hat{P} = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_n x_n$$

Where:

$\hat{P}$  is the predicted price.

$\beta_0, \beta_1, \beta_2, \dots, \beta_n$  are the model coefficients.

$x_1, x_2, \dots, x_n$  are the input features (e.g., historical prices, trading volume).

- **Black-Scholes Equation for Option Pricing:** The Black-Scholes equation is a fundamental expression in options pricing:

$$C(S, t) = S \cdot N(d_1) - X e^{-r(T-t)} \cdot N(d_2)$$

Where:

$C(S, t)$  is the option price.

$S$  is the current stock price.

$X$  is the strike price.

$r$  is the risk-free interest rate.

$T$  is the time to expiration.

$t$  is the current time.

$N(d_1)$  and  $N(d_2)$  are cumulative distribution functions of the standard normal distribution.

$d_1$  and  $d_2$  are parameters derived from the above variables.

**Portfolio Optimization:** Portfolio optimization aims to find the optimal allocation of assets in a portfolio. The mathematical expression for portfolio returns and risk is often used in this context:

$$R_p = \sum_{i=1}^n w_i \cdot R_i$$

$$\sigma_p^2 = \sum_{i=1}^n \sum_{j=1}^n w_i \cdot w_j \cdot \sigma_i \cdot \sigma_j \cdot \rho_{ij}$$

Where:

$R_p$  is the portfolio return.

$\sigma_p$  is the portfolio risk.

$w_i$  is the weight of asset  $i$  in the portfolio.

$R_i$  is the return of asset  $i$ .

$\sigma_i$  is the standard deviation of asset  $i$ 's returns.

$\rho_{ij}$  is the correlation between the returns of assets  $i$  and  $j$ .

3. **Natural Language Processing (NLP):** NLP applications include sentiment analysis, language translation, chatbots, and text summarization. These

technologies enable machines to understand and generate human language, improving customer service and automating various textual tasks.

- **Term Frequency-Inverse Document Frequency (TF-IDF):** TF-IDF is used to represent the importance of words in a document relative to a collection of documents. The TF-IDF weight of a term  $\mathbf{t}$  in a document  $\mathbf{d}$  is calculated as:

$$TD - IDF(\mathbf{t}, \mathbf{d}) = TF(\mathbf{t}, \mathbf{d}) \cdot IDF(\mathbf{t})$$

Where:

$TF(\mathbf{t}, \mathbf{d})$  is the term frequency of term  $\mathbf{t}$  in document  $\mathbf{d}$ .

$IDF(\mathbf{t})$  is the inverse document frequency of term  $\mathbf{t}$  across the document collection

- **Word Embeddings (e.g., Word2Vec, GloVe):** Word embeddings represent words as dense vectors in a continuous vector space. Word2Vec, for example, is trained to predict the context of words in a large text corpus. The cosine similarity between word vectors is often used to measure word similarity:

$$similarity(\mathbf{w}_1, \mathbf{w}_2) = \frac{\mathbf{w}_1 \cdot \mathbf{w}_2}{\|\mathbf{w}_1\| \cdot \|\mathbf{w}_2\|}$$

Where:

$w_1$  and  $w_2$  are the word vectors.

• represents the dot product.

$\|\cdot\|$  denotes vector norm.

4. **E-commerce and recommendation systems:** Machine learning powers recommendation engines on platforms like Netflix and Amazon, suggesting products or content based on user behavior and preferences. These systems enhance user experience and increase sales.

- Collaborative Filtering for Recommendation Systems: Collaborative filtering models, such as matrix factorization, use mathematical expressions to predict user preferences based on historical user-item interactions.

Expression for the matrix factorization equation:

$$R \approx P \cdot Q^T$$

Where:

$R$  is the user-item interaction matrix.

$P$  is the user matrix.

$Q$  is the item matrix.

• represents matrix multiplication.

- **Content-Based Filtering for Recommendation Systems:** Content-based filtering models use mathematical expressions to recommend items to users based on their attributes and preferences. One common expression is the cosine similarity between user profiles and item attributes.

$$\mathit{similarity}(\mathbf{u}, \mathbf{i}) = \cos(\theta) = \frac{\mathbf{u} \cdot \mathbf{i}}{\|\mathbf{u}\| \cdot \|\mathbf{i}\|}$$

Where:

$\mathbf{u}$  is the user's profile vector.

$\mathbf{i}$  is the item's attribute vector.

- **Logistic Regression for Fraud Detection:** Logistic regression models can be used to detect fraudulent transactions based on features like transaction amount, location, and user behavior. The logistic regression equation is used to model the probability of a transaction being fraudulent.

$$P(\mathit{fraudulent}) = \frac{1}{1 + e^{-(\beta_0 + \beta_1 x_1 + \dots + \beta_n x_n)}}$$

Where:

$P(\mathit{fraudulent})$  is the probability of fraud.

$\beta_0, \beta_1, \dots, \beta_n$  are the model coefficients.

$x_1, \dots, x_n$  are the input features.

5. **Autonomous Vehicles:** Self-driving cars rely heavily on machine learning algorithms to perceive their environment, make decisions, and navigate safely. Computer vision and sensor data analysis are key components of autonomous vehicle systems.

- **Sensor Fusion and State Estimation:** Autonomous vehicles typically use multiple sensors, such as LiDAR, radar, cameras, and GPS, to perceive their surroundings. Sensor fusion techniques, such as Kalman filters or particle filters, are used to estimate the vehicle's state, including its position, velocity, and orientation. The state estimation equations can be represented as:

$$\mathbf{x}_t = f(\mathbf{x}_{t-1}, \mathbf{u}_t, \mathbf{z}_t)$$

Where:

$\mathbf{x}_t$  is the vehicle's state at time  $t$ .

$\mathbf{u}_t$  is the control input.

$\mathbf{z}_t$  is the sensor measurement.

$f$  represents the state transition function.

- **Path Planning:** Path planning algorithms help autonomous vehicles navigate from their current location to a desired destination while avoiding

obstacles. A common mathematical representation is the use of cost functions and optimization techniques to find an optimal path:

$$\mathit{argmin} \sum_{i=1}^n C_i$$

Where:

$C_i$  represents the cost associated with each potential path segment.

- **Simultaneous Localization and Mapping (SLAM):** SLAM algorithms help autonomous vehicles build and update maps of their environment while simultaneously localizing themselves within these maps. SLAM can be represented as a mathematical optimization problem:

$$\mathit{argmin} \sum_{i=1}^N \mathit{error}(\mathbf{m}_i, \mathbf{x})$$

Where:

$N$  is the number of map features.

$\mathbf{m}_i$  represents the  $i$ th map feature.

$\mathbf{x}$  is the vehicle's pose.

6. **Manufacturing and Predictive Maintenance:** Machine learning helps optimize manufacturing processes, predict equipment failures, and reduce

downtime. Sensors collect data on machinery, which is then used to predict maintenance needs and prevent breakdowns.

- **Predictive Maintenance with Failure Prediction:** Predictive maintenance models aim to predict when a machine or equipment is likely to fail. One common model used for this purpose is the Weibull distribution, often employed in reliability engineering. The Weibull distribution's probability density function is given by:

$$f(t; \lambda, k) = \frac{k}{\lambda} \left(\frac{t}{\lambda}\right)^{k-1} e^{-(t/\lambda)^k}$$

Where:

$f(t; \lambda, k)$  is the probability density function at time  $t$ .

$\lambda$  is the scale parameter (related to the mean time between failures).

$k$  is the shape parameter (related to the failure distribution).

Other models, such as survival analysis and machine learning techniques like Random Forest or LSTM, can also be used for predictive maintenance.

- **Regression Models for Process Optimization:** Regression models are used to optimize manufacturing processes. For example, in chemical manufacturing, response surface models are used to relate process variables

to product quality. The mathematical expression for a response surface model might be a polynomial equation.

$$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_n X_n + \epsilon$$

Where:

$Y$  is the response variable (e.g., product quality).

$X_1, X_2, \dots, X_n$  are process variables.

$\beta_0, \beta_1, \beta_2, \dots, \beta_n$  are model coefficients.

$\epsilon$  represents the error term.

- **Time Series Analysis for Sensor Data:** In manufacturing, sensor data often involves time series data. Time series analysis techniques, such as autoregressive models (AR), moving average (MA) models, and autoregressive integrated moving average (ARIMA) models, are used to model and forecast sensor data.

For example, an ARIMA(p, d, q) model might be represented as:

$$X_t = C + \phi_1 X_{t-1} + \phi_2 X_{t-2} + \dots + \phi_p X_{t-p} + \theta_1 \epsilon_{t-1} + \theta_2 \epsilon_{t-2} + \dots + \theta_q \epsilon_{t-q} + \epsilon_t$$

Where:

$X_t$  is the observed value at time  $t$ .

$\phi_1, \phi_2, \dots, \phi_p$  are autoregressive parameters.

$\theta_1, \theta_2, \dots, \theta_q$  are moving average parameters

$C$  is a constant term.

$\epsilon_t$  is white noise at time  $t$ .

7. **Marketing and Customer Segmentation:** Marketers use machine learning to analyze customer data and segment audiences for targeted advertising. It can also optimize ad placements and budget allocation for better ROI.

- **K-Means Clustering:** K-means clustering is a popular unsupervised learning algorithm used for customer segmentation. Given  $K$  clusters and data points  $X$ , it can be represented as:

$$\sum_{i=1}^k \sum_{j=1}^{N_i} \|x_j - \mu_i\|^2$$

Where:

$N_i$  is the number of data points in cluster  $i$ .

$\mu_i$  is the centroid of cluster  $i$ .

- **Logistic Regression for Targeted Marketing:** Logistic regression models can be used to predict customer response to marketing campaigns. The mathematical expression for logistic regression is the sigmoid function:

$$P(\text{response}) = \frac{1}{1 + e^{-(\beta_0 + \beta_1 x_1 + \dots + \beta_n x_n)}}$$

Where:

$P(\text{response})$  is the probability of a customer responding.

$\beta_0, \beta_1, \dots, \beta_n$  are the model coefficients.

$x_1, \dots, x_n$  are the input features.

8. **Image and Video Analysis:** Beyond healthcare, computer vision is applied in fields like security (facial recognition), agriculture (crop monitoring), and entertainment (special effects). It can identify objects, track motion, and recognize patterns in visual data.

- **Convolutional Neural Networks (CNNs):** CNNs are a fundamental model for image analysis. The mathematical expressions in CNNs involve convolution operations, activation functions, and pooling layers. The forward pass through a convolutional layer can be represented as:

$$Z = W \cdot X + b$$

Where:

$Z$  is the output feature map.

$W$  is the filter weights.

$X$  is the input image.

$b$  is the bias term.

- **Activation Functions:** Various activation functions, such as the Rectified Linear Unit (ReLU), are used to introduce non-linearity into neural networks.

The ReLU activation function can be expressed as:

$$f(x) = \max(0, x)$$

9. **Gaming:** In the gaming industry, ML is used to create intelligent non-player characters (NPCs), adapt game difficulty based on player performance, and generate realistic game environments.

- **Physics Simulation:** Physics engines in games often use mathematical expressions to simulate physical interactions, such as Newton's second law for motion:

$$F = ma$$

Where:

$F$  is the force applied.

$m$  is the mass of the object.

$a$  is the acceleration.

- **Game Mechanics:** Game mechanics often involve mathematical equations to calculate player scores, health, experience points, and other in-game attributes. For example, experience points (XP) may be calculated using an equation like:

$$XP = baseXP \times multiplier$$

- **Reinforcement Learning (RL):** RL is used to train game agents to make decisions in dynamic environments. The Bellman equation is a fundamental equation in RL:

$$Q(s, a) = R(s, a) + \gamma \sum_{s'} P(s' | s, a) \max_a Q(s', a)$$

Where:

$Q(s, a)$  is the action-value function.

$R(s, a)$  is the immediate reward of taking action  $a$  in state  $s$ .

$\gamma$  is the discount factor.

$P(s' | s, a)$  is the transition probability to state  $s'$  from state  $s$  after taking action  $a$ .

**10. Supply Chain and Inventory Management:** Machine learning can optimize inventory levels, predict demand, and streamline supply chain operations, reducing costs and improving delivery times.

- **Regression Models:** Regression models, such as linear regression, are used to forecast demand by considering various features like seasonality, marketing promotions, and economic indicators. The mathematical expression for linear regression is:

$$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_n X_n + \epsilon$$

Where:

$Y$  is the predicted demand.

$X_1, X_2, \dots, X_n$  are the input features.

$\beta_0, \beta_1, \beta_2, \dots, \beta_n$  are the model coefficients.

$\epsilon$  represents the error term.

- **Economic Order Quantity (EOQ):** The EOQ model calculates the optimal order quantity to minimize inventory holding costs and ordering costs. The mathematical expression for EOQ is:

$$EOQ = \sqrt{\frac{2DS}{H}}$$

Where:

$D$  is the annual demand.

$S$  is the ordering cost per order.

$H$  is the holding cost per unit per year.

- **Inventory Turnover:** Inventory turnover models aim to find the balance between maintaining sufficient stock levels and minimizing excess inventory.

The mathematical expression for inventory turnover is:

$$\text{Inventory Turnover} = \frac{\text{Cost of Goods Sold (COGS)}}{\text{Average Inventory Value}}$$

## CONCLUSIONS

In summary, mathematical models and expressions are fundamental components of machine learning, serving as the foundation for developing algorithms across various domains. These models are versatile and applicable in fields like healthcare, finance, marketing, and more, enabling solutions to complex problems. Machine

learning encompasses diverse model types, including regression, classification, clustering, and neural networks, each with its mathematical principles tailored to specific tasks for national development.

## **CHAPTER FIVE**

### **RECOMMENDATION AND CONCLUSION**

#### **5.1 RECOMMENDATION**

1. **Continuous Learning:** Stay updated with the latest developments in machine learning. The field is rapidly evolving, and ongoing learning is

crucial to remain competitive and make informed decisions regarding its applications.

2. **Ethical Considerations:** Prioritize ethical practices in machine learning applications. Be vigilant about data biases, transparency and fairness. Ensure that your ML systems adhere to ethical guidelines and regulations.
3. **Interdisciplinary collaboration:** Collaborate with experts from diverse domains. Combining ML expertise with domain can lead to innovative and effective solutions in various industries.
4. **Data Quality:** Invest in data quality and data management practices. High-quality, clean data is the foundation of successful ML applications.
5. **Interpretability:** Focus on making ML models more interpretable, especially in critical applications where transparency is essential.

## 5.2 CONCLUSION

Machine Learning has emerged as a transformative force with vast applications in numerous sectors, offering efficiency, accuracy, and automation. While it presents exciting opportunities, it also poses challenges related to data quality, ethical

considerations and interpretability. The ability to choose appropriate ML algorithms and mathematical models is essential for achieving success.

To harness the full potential of ML, practitioners must continuously update their knowledge, collaborate across disciplines, and prioritize ethical AI development.

By adhering to these recommendations and staying vigilant, the world of machine learning can continue to advance, bringing innovative solutions to complex problems while maintaining ethical and responsible practices for national development.

## **REFERENCES**

*Artificial Intelligence: A Modern Approach* - by Stuart Russell, Peter Norvig

*Neural Networks in Computer Intelligence* - by Limin Fu

Artificial Intelligence: *A Guide to Intelligent Systems* - by Michael Negnevitsky

Data Mining: *Practical Machine Learning Tools and Techniques* - by Ian H. Witten, Eibe Frank

Davenport, T. H., Ronanki, R., *Artificial Intelligence for the Real World*, 2018.  
Accessed July 21, 2020.

*An Executive Primer on Artificial General Intelligence*, McKinsey & Company,  
April 2020. Accessed July 21, 2020.

Sutton, R. S., Barto, A. G., *Reinforcement Learning: An Introduction*, MIT Press.  
Accessed September 22, 2020.

