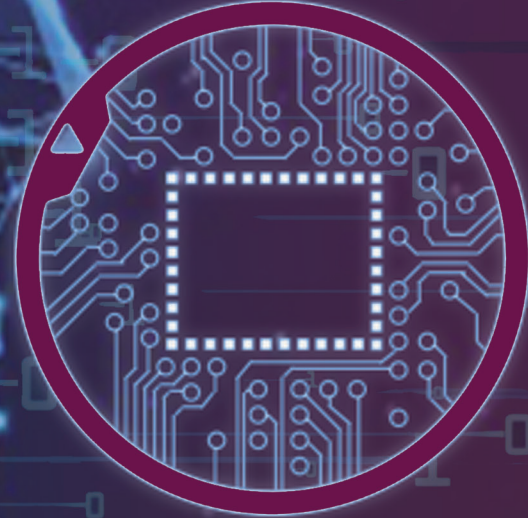




VOLUME - 4  
ISSUE - 2  
2024

e-ISSN: 2791-8335

# JOURNAL OF ARTIFICIAL INTELLIGENCE AND DATA SCIENCE



**İZMİR KÂTİP ÇELEBİ UNIVERSITY**

Artificial Intelligence and Data Science  
Research and Application Center



# JAIDA

[HTTPS://DERGIPARK.ORG.TR/PUB/JAIDA](https://dergipark.org.tr/pub/jaida)

**Privilege Owner**

Prof. Dr. Saffet Köse, Rector (İzmir Kâtip Çelebi University)

**Editor in Chief**

Prof. Dr. Ayşegül Alaybeyoğlu (İzmir Kâtip Çelebi University)

**Associate Editors**

Assoc. Prof. Dr. Levent Aydın (İzmir Kâtip Çelebi University)

Dr. Ümit Sarp (İzmir Kâtip Çelebi University)

**Managing Editor**

Dr. Ümit Sarp (İzmir Kâtip Çelebi University)

**Grammar Editor**

Feyyaz Demirer (İzmir Demokrasi University)

**International Advisory Board**

Prof. Dr. Abd Samad Hasan Basari (Universiti Tun Hussein Onn Malaysia)

Prof. Dr. Filiz Güneş (Yıldız Teknik University)

Prof. Dr. Nejat Yumuşak (Sakarya University)

Prof. Dr. Chirag Paunwala (Sarvajanic College of Engineering and Tech.)

Prof. Dr. Narendra C. Chauhan (A D Patel Institute of Technology)

Prof. Dr. Saurabh Shah (GSFC University)

Prof. Dr. H. Seçil Artem (İzmir Institute of Technology)

Prof. Dr. Doğan Aydın (İzmir Kâtip Çelebi University)

Assoc. Prof. Dr. Amit Thakkar (Charusat University)

Assoc. Prof. Dr. Cheng Jin (Beijing Institute of Technology)

Assoc. Prof. Dr. Mustafa Emiroğlu (Tepecik Training and Research Hospital)

Assoc. Prof. Dr. Ali Turgut (Tepecik Training and Research Hospital)

Assoc. Prof. Dr. Mohd Sanusi Azmi (Universiti Teknikal Malaysia Melaka)

Assoc. Prof. Dr. Peyman Mahouti (İstanbul University- Cerrahpaşa)

Assoc. Prof. Dr. Mehmet Ali Belen (İskenderun Technical University)

Assoc. Prof. Dr. Ferhan Elmalı (İzmir Kâtip Çelebi University)

Assoc. Prof. Dr. Sharnil Pandya (Symbiosis International University)

Assist. Prof. Dr. Kadriye Filiz Balbal (Dokuz Eylül University)

Assist. Prof. Dr. Mansur Alp Toçoğlu (İzmir Kâtip Çelebi University)

Assist. Prof. Dr. Emre Şatır (İzmir Kâtip Çelebi University)

Dr. Çağdaş Eşiyok (İzmir Kâtip Çelebi University)

Dr. Zihao Chen (Harbin Institute of Technology)

Dr. Aysu Belen (İskenderun Technical University)

**Editorial Board****Engineering and Architecture**

Prof. Dr. Merih Palandöken (İzmir Kâtip Çelebi University)

Prof. Dr. Ayтуğ Onan (İzmir Kâtip Çelebi University)

Assoc. Prof. Dr. Mehmet Ali Belen (İskenderun Technical University)

Assist. Prof. Dr. Esra Aycan Beyazıt (İzmir Kâtip Çelebi University)

Assist. Prof. Dr. Mehmet Erdal Özbek (İzmir Kâtip Çelebi University)

Assist. Prof. Dr. Nesibe Yalçın (Erciyes University)

Assist. Prof. Dr. Olgun Aydın (Gdansk University of Technology)

Assist. Prof. Dr. Onan Güren (İzmir Kâtip Çelebi University)

Assist. Prof. Dr. Osman Gökalg (İzmir Institute of Technology)

Assist. Prof. Dr. Serpil Yılmaz (İzmir Kâtip Çelebi University)

Assist. Prof. Dr. Semih Çakır (Zonguldak Bülent Ecevit University)

Dr. M. Mustafa Bahşı (Manisa Celal Bayar University)

Dr. Sümeyye Sınır (İzmir Kâtip Çelebi University)

**Natural Science**

Assoc. Prof. Dr. Seda Oğuz Ünal (Sivas Cumhuriyet University)

Assoc. Prof. Dr. Şule Ayar Özbal (Yaşar University)

Assist. Prof. Dr. Ahmet Emin (Karabük University)

Assist. Prof. Dr. Ezgi Kaya (İğdır University)

Dr. Göknur Giner (The University of Melbourne)

Dr. Ümit Sarp (İzmir Kâtip Çelebi University)

**Social Science**

Prof. Dr. Murat Kayacan (İzmir Kâtip Çelebi University)

Assoc. Prof. Dr. Bekir Emiroğlu (İzmir Kâtip Çelebi University)

Assoc. Prof. Dr. Gizay Daver (Zonguldak Bülent Ecevit University)

Assoc. Prof. Dr. Ersin Kanat (Zonguldak Bülent Ecevit University)

Assist. Prof. Dr. Hilal Kahraman (İzmir Kâtip Çelebi University)

Assist. Prof. Dr. Ümit Aydoğan (İzmir Kâtip Çelebi University)

**Health Science**

Assoc. Prof. Dr. Mustafa Ağâh Tekindal (İzmir Kâtip Çelebi University)

Assist. Prof. Dr. Ünzile Yaman (İzmir Kâtip Çelebi University)

**Education Science**

Assoc. Prof. Dr. Mustafa Ergun (Ondokuz Mayıs University)

Assist. Prof. Dr. Kadriye Filiz Balbal (Dokuz Eylül University)

**Arts and Design**

Assoc. Prof. Dr. Cem Çırak (İzmir Kâtip Çelebi University)

Assoc. Prof. Dr. Mucahit Yalçın Öztüfekci (İzmir Kâtip Çelebi University)

Assoc. Prof. Dr. Sehran Dilmaç (İzmir Kâtip Çelebi University)

**Aim & Scope**

The Journal of Artificial Intelligence and Data Science (JAIDA) is an international, scientific, peer-reviewed, and open-access e-journal. It is published twice a year and accepts only manuscripts written in English. The aim of JAIDA is to bring together interdisciplinary research in the fields of artificial intelligence and data science. Both fundamental and applied research are welcome. Besides regular papers, this journal also accepts research field review articles. Paper submission/processing is free of charge.

**Contact**

Web site: <https://dergipark.org.tr/pub/jaida> E-mail: [ikcujaida@gmail.com](mailto:ikcujaida@gmail.com)

Phone: +90 (232) 329 35 35/3731/ 1072

Fax: +90 (232) 325 33 60

Mailing address: İzmir Kâtip Çelebi Üniversitesi, Yapay Zekâ ve Veri Bilimi Uygulama ve Araştırma Merkezi, Balatçık Kampüsü, Çiğli Ana Yerleşkesi, 35620, İzmir, TÜRKİYE

## ÖNSÖZ

Yapay Zeka ve Veri Bilimi alanındaki teknolojik ve bilimsel gelişmeler; Yapay Zekanın endüstri, sağlık, otomotiv, ekonomi, eğitim gibi bir çok farklı alanda uygulanmasına imkan sağlamıştır. Ülkemiz Ulusal Yapay Zeka Stratejisinde; yeni bir çağın eşiğine gelindiği, yapay zekayla üretim süreçleri, meslekler, gündelik yaşam ve kurumsal yapıların yeni bir dönüşüm sürecine girdiği vurgulanarak, Yapay Zekanın öneminden bahsedilmiştir.

Sayın Cumhurbaşkanımızın da belirttiği gibi ülkemiz adına insan odaklı yeni bir atılım yapmanın zamanının geldiğine inanıyoruz. Yapay zeka çağına geçiş noktasında Türkiye'nin lider ülkelerden biri olması motivasyonu ile üniversitemizde yapay zeka teknolojilerinin kullanıldığı projeler gerçekleştirmekte, kongreler ve bilimsel etkinlikler düzenlemekteyiz.

Günümüz dünyasına rengini veren dijital teknolojilerin odağındaki ana unsurun yapay zeka teknolojilerinin olduğu düşüncesi ile yola çıkarak hazırlamış olduğumuz Yapay Zekâ ve Veri Bilimi Dergisinin, Ülkemiz Ulusal Yapay Zeka Stratejisinde belirtilen "Dijital Türkiye" vizyonu ve "Milli Teknoloji Hamlesi" kalkınma hedefleri doğrultusunda katkı sağlayacağı inancındayız.

Dergimizin hazırlanmasında emeği geçen üniversitemiz Yapay Zekâ ve Veri Bilimi Uygulama ve Araştırma Merkez Müdürü, Baş Editör Prof. Dr. Ayşegül ALAYBEYOĞLU'na, Editör ve Danışma kurulu üyelerine, akademik çalışmalarını ile sağladıkları destek için tüm yazarlara, hakem olarak görev alan değerli bilim insanlarına teşekkür eder, dergimizin yeni sayısının ülkemize hayırlı olmasını dilerim.

**Prof. Dr. Saffet KÖSE, Rektör**

**Dergi Sahibi**

## **PREFACE**

**Technological and scientific developments in Artificial Intelligence and Data Science enabled the application of Artificial Intelligence in many different fields such as industry, health, automotive, economy and education. In our country's National Artificial Intelligence Strategy; the importance of Artificial Intelligence was mentioned by emphasizing the transformation process of production processes, occupations, daily life and corporate structures with artificial intelligence.**

**As stated by our President, we believe that the time has come to make a new human-oriented breakthrough on behalf of our country. With the motivation of Turkey being one of the leading countries at the point of transition to the age of artificial intelligence, we realize projects in which artificial intelligence technologies are used, and organize congresses and scientific events at our university.**

**We have prepared the Journal of Artificial Intelligence and Data Science with the idea that the main element in the focus of digital technologies that color today's world is artificial intelligence technologies, and we believe that our journal will contribute to the development goals of the "Digital Turkey" vision and "National Technology Move" stated in the National Artificial Intelligence Strategy of our country.**

**I would like to thank Prof. Dr. Ayşegül ALAYBEYOĞLU, the Director of Artificial Intelligence and Data Science Application and Research Center of our university. I would also like to thank to Editor and Advisory Board members, to all authors for their supports with their academic studies and to reviewers for their contributions to the preparation of our journal. I wish the new issue of our journal to be beneficial for our country.**

**Prof. Dr. Saffet KÖSE, Rector**

**Privilege Owner**

## **BAŞ EDITÖR'DEN**

**Değerli Araştırmacılar ve Dergi Okuyucuları;**

**İzmir Kâtip Çelebi Üniversitesi Yapay Zekâ ve Veri Bilimi Uygulama ve Araştırma Merkezi olarak Rektörümüz Prof. Dr. Saffet Köse sahipliğinde Yapay Zekâ ve Veri Bilimi Dergisinin 4. cilt 2. sayısını sizlerle buluşturmanın gururunu yaşamaktayız.**

**İzmir Kâtip Çelebi Üniversitesi Yapay Zekâ ve Veri Bilimi Uygulama ve Araştırma Merkezi olarak hedefimiz; Cumhurbaşkanlığı Dijital Dönüşüm Ofisi Başkanlığı ve Sanayi ve Teknoloji Bakanlığı tarafından hazırlanan “Ulusal Yapay Zekâ Stratejisi” hedefleri doğrultusunda dergi, kongre, eğitim, bilimsel etkinlikler ve proje faaliyetleri gerçekleştirilerek ülkemizin yapay zekâ alanındaki gelişim sürecine katkı sağlamaktır.**

**Farklı üniversitelerden, bilimsel disiplinlerden ve alanlardan değerli araştırmacıların İngilizce dilinde hazırlamış oldukları 6 adet araştırma ve 1 adet derleme makalesi bu sayı kapsamında sunulmaktadır. Siz değerli araştırmacılarımızın destekleri ile kaliteyi daha da arttırarak en kısa sürede ulusal ve uluslararası indekslerde daha çok taranan bir dergi olmayı hedeflemekteyiz.**

**Dergimizin yayın hayatına başlaması ve tüm merkez faaliyetlerinde büyük desteklerini gördüğümüz başta Rektörümüz Prof. Dr. Saffet KÖSE olmak üzere; dergimize olan destekleri için tüm yazarlara, dergimizin yayına hazırlanmasında heyecanla çalışan ve çok büyük emek harcayan Baş Editör Yardımcılarına, Editör ve Danışma kurulu üyelerimize, hakem olarak görev alan tüm değerli bilim insanlarına en derin şükranlarımı sunarım.**

**Saygılarımla,**

**Prof. Dr. Ayşegül ALAYBEYOĞLU**

**Baş Editör**

## **LETTER FROM THE EDITOR-IN-CHIEF**

**Dear Researchers and Readers of the Journal,**

**As İzmir Katip Çelebi University Artificial Intelligence and Data Science Application and Research Center, we are proud to present you the volume 4 issue 2 of the Journal of Artificial Intelligence and Data Science (JAIDA), hosted by our Rector Prof. Dr. Saffet Köse.**

**As İzmir Katip Çelebi University Artificial Intelligence and Data Science Application and Research Center, our goal is; to contribute to the development process of our country in the field of artificial intelligence by carrying out journals, congresses, education, scientific events and project activities in line with the objectives of the "National Artificial Intelligence Strategy" prepared by the Digital Transformation Office of the Presidency of Türkiye and the Ministry of Industry and Technology.**

**6 research articles and 1 review article prepared by valuable researchers from different universities, scientific disciplines and fields are presented within the scope of this issue. With the support of esteemed researchers, we aim to increase the quality even more and become a journal that is scanned in national and international indexes more as soon as possible.**

**I would like to express my deepest gratitude to Our Rector, Prof. Dr. Saffet KÖSE, who supported the publication of our journal and the research center's activities; to all the authors for their support to our journal; to our Associate Editors, who worked enthusiastically and put great efforts into the preparation of our journal; to our Editorial and Advisory Board members, and all esteemed scientists who served as reviewer.**

**Best Regards,**

**Prof. Dr. Ayşegül ALAYBEYOĞLU**

**Editor-in-Chief**

# CONTENTS

Earthquake Probability Prediction with Decision Tree Algorithm: The Example of Izmir, Türkiye (Research Article) <b>İsmahan ERMIŞ, İsa CÜREBAL</b> .....	59
Optimization of Thermal Management for Cooling System of Power Electronics Modules Consisting Insulated-Gate Bipolar Transistor Using Neuro-Regression Analysis and Non-Traditional Algorithms (Research Article) <b>Melih SAVRAN, Ece Nur YÜNCÜ, Levent AYDIN</b> .....	68
Time Series Prediction of Heart Rate Using Deep Learning Models (Research Article) <b>Emir EVCİL</b> .....	79
Artificial Intelligence Applications in Drug Discovery and Research (Review Article) <b>Seyma MINTAS, Canan SEVİMLİ-GUR</b> .....	87
Artificial Intelligence Based Customer Risk Classification for Receivables Management of Businesses (Research Article) <b>Şaban Can TİRYAKİ, Adnan KAVAK</b> .....	97
Production of a Cost Effective Microstrip Antenna operating at 2.4 GHz and 5 GHz (Research Article) <b>Burak DOKMETAS</b> .....	104
Machine Learning Approaches for Prediction of Alzheimer’s Disease (Research Article) <b>Kadriye Filiz BALBAL</b> .....	110

# Earthquake Probability Prediction with Decision Tree Algorithm: The Example of Izmir, Türkiye

İsmahan ERMİŞ<sup>1\*</sup>, İsa CÜREBAL<sup>2</sup>

## Abstract

This study investigates earthquake records in the Izmir province of western Türkiye, focusing on seismic activity prediction through the application of decision tree models. Utilizing earthquake data from 1900 to 2024, including magnitude, depth, latitude, and longitude variables, the aim is to estimate future seismic events in a region known for its significant earthquake risks. The decision tree model, a machine learning approach, was trained with 80% of the dataset and tested on the remaining 20%. Performance was assessed using metrics such as precision, recall, F1 score, and overall accuracy, with the model achieving an accuracy rate of 92%. However, its ability to predict larger earthquakes was hindered due to the limited availability of data for higher-magnitude events. A chi-square test demonstrated a statistically significant relationship between earthquake depth and magnitude. Additionally, a risk analysis map was created using Geographic Information Systems (GIS), highlighting fault lines and areas prone to frequent seismic activity. The study concludes that while the decision tree model is effective for predicting smaller earthquakes, the accuracy for larger events could be improved with more comprehensive data. These findings underscore the importance of targeted earthquake preparedness in Izmir, particularly in coastal areas susceptible to both seismic events and secondary hazards like tsunamis.

**Keywords:** *Artificial Intelligence; Decision Trees; Earthquake; Izmir.*

## 1. Introduction

This study aims to analyze the earthquake records of Izmir, one of the most important provinces in western Türkiye, and its immediate surroundings (between 37° 45' and 39° 15' north latitude, 26° 15' and 28° 20' east longitude) in the paleo-seismology based on magnitude, depth, latitude and longitude variables. This earthquake data analysis aimed to estimate the magnitudes and locations of earthquakes that may occur in the future. The study's outputs are expected to be beneficial for Izmir, which is in a position to have serious risks in terms of seismic activity in Türkiye.

Various methods have been used to predict and analyze earthquakes. Other machine learning algorithms such as Decision Trees (DTs), Support Vector Machines (SVM), Random Forests (RF) and Neural Networks (NN) have also been used to improve earthquake prediction models [31]. For example, Support Vector Machines (SVM) [18], Random Forests Algorithm (RFA) [17], K-Nearest Neighbors (KNN) [22], Long Short-Term Memory (LSTM) [7] and Decision Trees (DTs) [2, 3, 6, 8, 12, 24, 29, 33, 36, 37, 39]. In this study, the decision trees method, which is assumed to provide accurate and reliable earthquake predictions, was preferred [30].

The decision trees method, a popular machine learning algorithm, is widely used in various studies to predict earthquake magnitudes and assess the impact of seismic events [16, 26]. The decision tree algorithm, especially the C4.5 variant, shows promise in characterizing factors that predict earthquakes and developing decision rules based on environmental variables and seismic data [4, 19]. Decision trees can identify important predictors by analyzing historical earthquake data and rank attributes using information theory criteria such as entropy and frequency of occurrence [4].

Located in the Coastal Aegean section of the Aegean Region in western Türkiye, Izmir is the third largest city in the country in terms of both population and economic importance. Izmir and the surrounding region is seismically active and has experienced damaging earthquakes in the past important human and economic effects. The tectonic structure of Western Anatolia, characterized by complex fault systems, contributes to the high seismic hazard in the region [13, 15, 27, 34]. The study focuses on the seismic activity in Izmir and examines whether decision trees can effectively handle complex seismic patterns in the region.

\*Corresponding author

İsmahan ERMİŞ; Balıkesir University, Institute of Social Sciences, Division of Geography, Balıkesir, Türkiye; e-mail: [ermisismahan@hotmail.com](mailto:ermisismahan@hotmail.com);

 0009-0007-7899-645X

İsa CÜREBAL; Balıkesir University, Faculty of Art and Science, Department of Geography, Balıkesir, Türkiye; e-mail: [curebal@balikesir.edu.tr](mailto:curebal@balikesir.edu.tr);

 0000-0002-3449-1595



The research questions of this study are as follows;

**MAIN QUESTION:** When the past earthquake data are analyzed, can the characteristics and distribution of earthquakes that may occur in the future in Izmir be predicted?

**SUBQUESTION 1:** How accurately can future earthquakes be predicted with the decision tree model?

**SUBQUESTION 2:** What are the advantages and limitations of the decision tree model?

**SUBQUESTION 3:** Is there a statistical relationship or similarity between the magnitudes and depths of earthquakes in Izmir province?

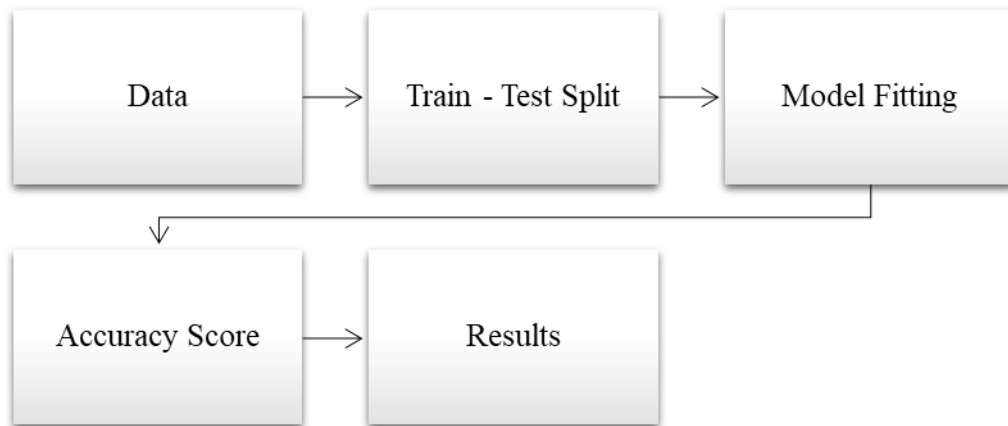
**SUBQUESTION 4:** What is the distribution of earthquakes that occurred in Izmir province in the last century and in which areas are these earthquakes concentrated?

## 2. Methodology

In this study, a dataset of earthquake records provided by Boğaziçi University Kandilli Observatory and Earthquake Research Institute (KRDAE) was used. This dataset, which covers the earthquakes that occurred between 1900 and 2024 (01.10.2024), includes the year, depth, latitude, longitude and magnitude information of the earthquakes.

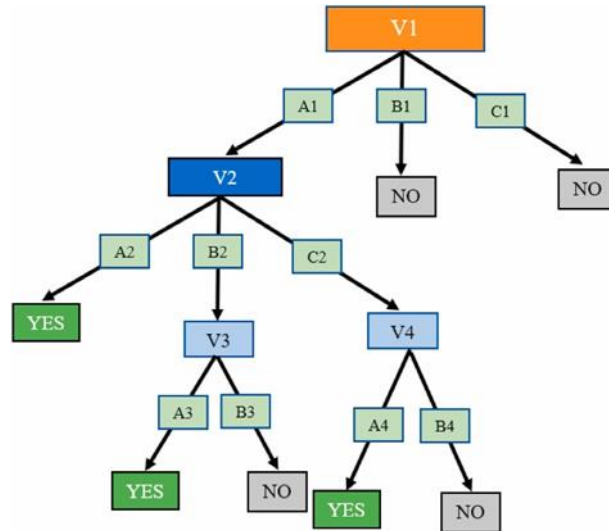
The study focuses on earthquakes with magnitudes ( $x_M$ ) between 3.0 and 8.0, and the data set is organized and visualized before applying the decision tree algorithm.

The decision tree model was implemented in Colab environment and 80% of the data was allocated for training and 20% for testing (Figure 1). Model fitting was applied on the separated test - training set and the results of this model were obtained. When latitude, longitude and depth information was given to the model as the outputs of the results, the model was trained to give as a result output the magnitude of potential future earthquakes based on the given information.



**Figure 1.** Flow Diagram of Methodology

The decision tree algorithm is first applied to the seismic dataset to generate binary predictions. The process starts by selecting a root node from the original dataset. In each iteration, the algorithm calculates the information gain for all available attributes. The attribute with the highest information gain is selected and the dataset is divided into subsets based on this attribute. The algorithm then processes each subset iteratively, considering only the remaining attributes that were not used in the previous splits. This iterative process continues until the decision tree is fully formed (Figure 2). The path from the root to the leaf node represents the values of the input variables, while each leaf node corresponds to a predicted value of the target variable. Decision trees are highly effective as classifiers that can capture complex variations in data [5]. This algorithm was used to estimate the magnitude of earthquakes based on their latitude, longitude and depth.



**Figure 2.** Decision Trees (DTs) Flowchart [21, 25]

We started by first analyzing the prediction performance of the method used for the prediction results of the model with a self-consistency test on the training set. Then, in order to evaluate the overall performance of the method more comprehensively, accuracy, recall and F1 score were calculated among the metrics used for performance analysis.

$$\text{Precision } (P) = \frac{TP}{TP + FP} \quad (1)$$

$$\text{Recall } (R) = \frac{TP}{TP + FN} \quad (2)$$

$$\text{F1 - Score} = \frac{2PR}{P + R} \quad (3)$$

$$\text{Accuracy} = \frac{TP + TN}{TP + FP + TN + FN} \quad (4)$$

In the above equations TP, FN, TN and FP are true positives, false negatives, true negatives and false positives respectively. Precision is the proportion of predicted hotspots that are true hotspots. Recall is the proportion of true positive hotspots that are predicted hotspots. F1 score is a measure to balance recall and precision. Accuracy is calculated as the ratio of correctly predicted areas in the model to the total dataset [9, 32].

After the implementation phase of the model was completely completed, the Confusion Matrix was used to visualize the distributions of the data and how many of these distributions were correctly predicted. Confusion Matrix is a tool that shows the morphs that are most likely to be confused with each other and improves performance prediction in ensemble detection tasks. It is also a technique used to measure the success rate of detection approaches [1, 38]. This method was applied to visually demonstrate the success rate. As a result of all the calculations of these processes, it is aimed to see what the advantages and limitations of the model are.

Chi-Square test was applied to reach a conclusion on whether there is a statistical relationship or similarity between the magnitudes and depths of earthquakes in Izmir region. Chi-square test is a statistical analysis method used to examine the relationship between categorical variables. Based on the research hypothesis, it determines whether there is a difference in the proportions of risk factors between groups by evaluating the differences between rows and columns [14, 28, 35]. The relationship between depth and magnitude was analyzed and visualized with the applied test.

Geographic Information Systems (GIS) is software that models geospatial environments, enhances analysis functions, and provides human-centered geographic information for better understanding and communication [23]. GIS supports a wide range of spatial queries and plays an important role in future location model development and application in fields such as Geography, Civil Engineering, and Computer Science [10]. The

data obtained as a result of queries can be presented in visualization methods such as points, lines or areas. The data is effectively integrated, managed and analyzed with geographic information from maps, images and text. These features provide a powerful tool for solving geospatial problems [11 ,20]. Since the data obtained in the study were in the form of point data, a general risk analysis map was created using these data. An earthquake map of the relevant study area was prepared, and thanks to this map, the distribution of earthquakes in and around Izmir province and where the magnitude ranges are more intense were visualized.

The answers to the questions asked in this study were obtained with Excel 2016, Google Colab and ArcMap 10.8 software. The Accuracy, Precision and F1-Score of the model, the results of the model were evaluated by examining the training and test datasets with the aim of successfully evaluating the prediction performance by processing the data within the framework of certain rules. GIS and data visualization techniques have made significant contributions to the analysis process of risky areas. The statistical relationship between the magnitudes, depths and occurrence regions of earthquakes in the Izmir region was visualized in tables.

**3. Findings**

Under this heading, the research questions given in the introduction of the study have been comprehensively answered. In order to make the findings obtained throughout the research process more understandable, they were visualized with different graphical methods and detailed using various tables. Thus, it is ensured that the results of the research are presented more clearly and clearly both quantitatively and qualitatively.

**MAIN QUESTION:** When the past earthquake data are analyzed, can the characteristics and distribution of earthquakes that may occur in the future in Izmir be predicted?

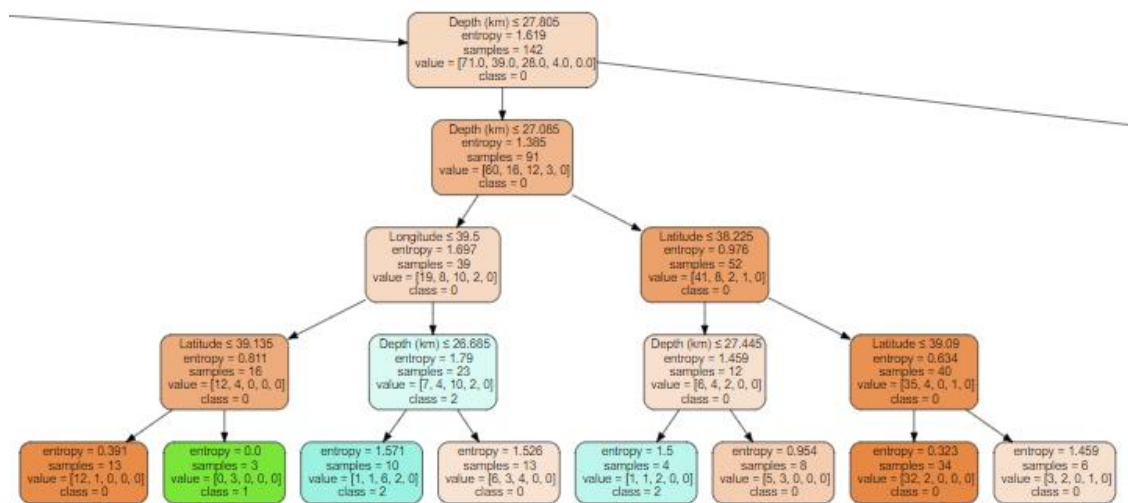
Yes, the data is predictable. As a result of the statistical analysis of the data, the classification success report of the model (Table 1) was created. Accordingly, the probabilities of predicting the magnitudes of earthquakes are objectified with Precision, Recall and F1-Score values.

**SUBQUESTION 1:** How accurately can future earthquakes be predicted with the decision tree model?

Table 1 shows that while high precision, recall and f1-score values were obtained especially for small earthquakes, these metrics decreased significantly for large earthquakes. The imbalance of the data, especially the limited number of high magnitude earthquakes, was found to negatively affect the success of the model in predicting these earthquakes.

**Table 1. Classification Report Table**

Class	Magnitud (xM)	Precision	Recall	F1-Score	Support
0	3.0 – 3.9	0.93	1.00	0.96	3014
1	4.0 – 4.9	0.24	0.03	0.05	189
2	5.0 – 5.9	0.14	0.03	0.05	35
3	6.0 and upper	0.00	0.00	0.00	3



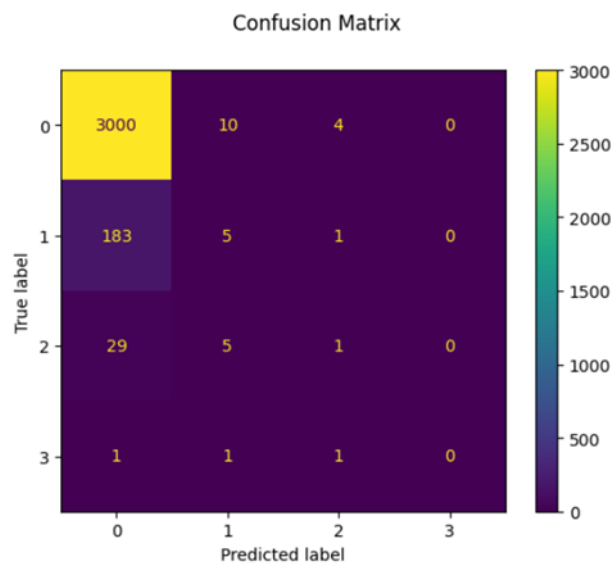
**Figure 3. Decision Trees' (DTs) One Tree Example**

The decision tree model in Figure 3 has one branch. In the output of this model, the answers given to the this question are given in Table 2 and Figure 3. In the table, it is determined at which level earthquake magnitudes can

be trained and at which level they can be estimated. In Figure 3, the data at the root node is divided into two based on the 'Depth' feature according to the limit of 27.805 km. The entropy value at this node is calculated as 1.619, which indicates that the dataset is quite diverse and complex. The value of “Samples” is 142 and there are a total of 142 data samples in this node. “Value” is [71, 39, 28, 4, 0], which means that 71 of the 142 samples belong to class 0, 39 to class 2 and 28 to class 3. The “Class” value is assigned as 0, indicating that class 0 has the most instances. Each “Value” value describes in detail to which classes the instances in the node belong. In this context, the decision tree performs the final classification by decomposing the data step by step. At leaf nodes, the entropy value approaches zero or becomes zero, indicating that the datasets have become completely homogeneous and there are only examples belonging to a single class.

**Table 2.** Model Performance Metrics

1. Training Performance For Regression Tree	1. Test Performance For Regression Tree
0.93	0.92
<b>Accuracy</b>	<b>Weighted Avg.</b>
0.92	0.90



**Figure 4.** Confusion Matrix

**SUBQUESTION 2:** What are the advantages and limitations of the decision tree model?

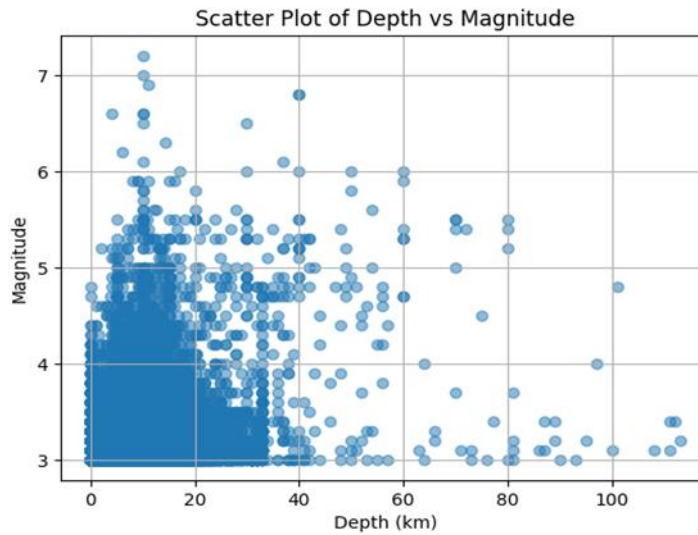
The results obtained from the model (Figure 4) show that the model is inadequate in predicting earthquakes with larger magnitudes. In fact, there are only 35 cases of earthquakes with magnitude 5.0 and only 3 cases of earthquakes with magnitude 6.0 and above. The small data set limited the model's ability to predict earthquakes of this magnitude. In such cases, the model was unable to accurately predict rare events (such as large earthquakes). However, when the overall model performance is analyzed in Table 3, the model provided a good prediction result of 92%.

**Table 3.** Model Performance Table

Accuracy	Precision	Recall	F1-Score
0.92	0.88	0.92	0.89

**SUBQUESTION 3:** Is there a statistical relationship or similarity between the magnitudes and depths of earthquakes in Izmir province?

In order to visually see the relationship between depth and magnitude, the table in Figure 5 was created. In the relevant table, a high density is observed especially in the 0-20 km depth range. It was revealed that many earthquakes occurred at shallow depths (close to the surface). In terms of magnitude, it is generally observed that values between 3-5 are more frequently observed, but in earthquakes deeper than 40 km, the magnitude shows a wider distribution and can reach values of 5 and above.



**Figure 5.** Scatter Plot of Depth vs Magnitude ( $xM$ )

As a result of the chi-square test performed to reveal the statistical relationship between depth and magnitude, the  $\chi^2$  (chi-square) value was 3924.06 and the p-value was 0.0. This result (Table 4) shows that there is a statistically significant relationship between depth and magnitude categories. The p-value is well below 0.05 (in this case very close to 0), indicating that there is a significant non-random relationship between depth and magnitude.

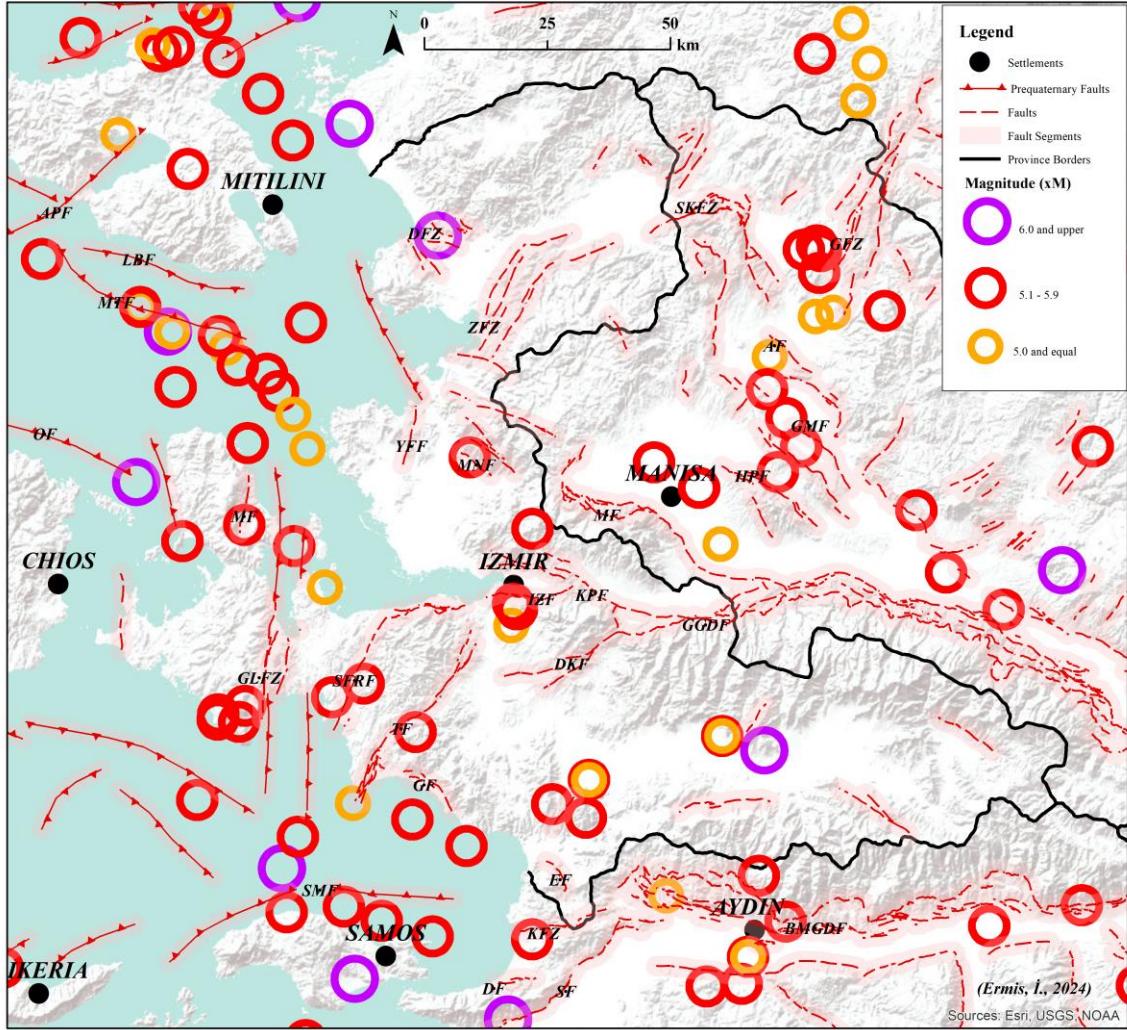
**Table 4.** Chi-Square Result

Chi-Square Test	P-Score
3924.06	0.0

**SUBQUESTION 4:** What is the distribution of earthquakes that occurred in Izmir province in the last century and in which areas are these earthquakes concentrated?

The distribution and concentration areas of earthquakes of 5.0 and above occurring in Izmir province were visualized using ArcGIS software. As a result of this visualization, an intensity map (Figure 6) showing the distribution of earthquakes and their hectic areas on specific fault segments was obtained.

Izmir province is a region with a high earthquake risk due to active fault lines on both land and seafloor (Figure 6). Frequent earthquakes, especially along active fault lines on the seafloor (north of Lesvos and south of Lesvos and north of Samos), pose a great danger to Izmir and the surrounding settlements. On the map, these earthquakes are particularly concentrated in coastal areas and the faults on the seafloor, which are very close to Izmir, significantly increase the seismic risk of the city. Izmir's proximity to the seashore also brings up secondary hazards such as earthquake-induced sea flooding. Therefore, coastal areas stand out as the riskiest areas due to their proximity to faults in the sea. Looking at the earthquakes that have occurred in Izmir, it is observed that earthquakes with a magnitude of 5.1 and above occur more frequently on the MTF (Mitilini Fault), IZF (Izmir Fault), GLFZ (Gülbaççe Fault Zone) and SMF (Samos Fault Zone) compared to other faults.



**Figure 6.** Distribution Map of Past Earthquakes (< 5.0) and Main Fault Zones: Information on fault lines was taken from the website of the General Directorate of MTA (<https://www.mta.gov.tr/>) and the faults on the seabed were taken from the EMODNet website (<https://emodnet.ec.europa.eu/>). While preparing the map, only earthquakes with magnitude 5.0 and above were selected because too much data causes image pollution.

The faults in Figure 6 and their descriptions are as follows: APF (Aghia Paraskevi Fault); LBF (Lesvos Basin Fault); MTF (Mitilini Fault); OF (Oinousses Fault); MF (Manisa Fault); GLFZ (Gülbahçe Fault Zone); SFRF (Seferihisar Fault); TF (Tuzla Fault); GF (Gümüldür Fault); SMF (Samos Fault); DF (Davutlar Fault); SF (Söke Fault); KFZ (Kuşadası Fault Zone); EF (Ephesus Fault); BMGDF (Büyük Menderes Graben Detachment Fault); DKF (Dagkızılca Fault); IZF (Izmir Fault); KPF (Kemalpaşa Fault); GGDF (Gediz Graben Detachment Fault); MNF (Menemen Fault Zone); YFF (Yenifoça Fault); ZFZ (Zeytindag Fault Zone); DFZ (Dikili Fault Zone); SKFZ (Soma-Kırkağaç Fault Zone); GFZ (Gelenbe Fault Zone); AF (Akhisar Fault); GMF (Gölmarmara Fault); HPF (Halitpaşa Fault).

#### 4. Conclusions and Suggestions

The findings of the study clearly show that Izmir is a region with high seismic risk. The analyses revealed that although we can successfully predict small earthquakes, the prediction accuracy for large earthquakes is limited. This is due to the lack of data on large earthquakes.

Forecasts using the decision tree model showed that the overall accuracy of the model was as high as 92%. However, it was observed that the model failed especially for large earthquakes. The limited data set for large earthquakes negatively affected the performance of the model. Especially for earthquakes of 5.0 and above, the precision, recall and f1-score values were almost zero. This shows that the model is inadequate in predicting such large events. According to Somodevilla et al., 2012 in the related literature review, it was mentioned that it is indeed possible to determine the seismic depth based on its magnitude.

According to the results of this study, the relationship between depth and magnitude was found to be statistically significant in the decision tree model. According to Asim et al., 2017, the Random Forest algorithm, which is included in the decision tree algorithm, gives 77% accuracy, while the decision trees algorithm, which is a detailed study of RFA in this study, gives 92% accuracy rate. However, according to Mignan, A., & Broccardo, M. (2020), the argument that it may not be fine-tuned when there is a basic machine learning classifier (Support Vector Machine, Decision Trees, Naïve Bayes, etc.) in this study due to the uneven distribution of the data confirmed this view. However, according to Mignan, A., & Broccardo, M. (2020), when it is a basic machine learning classifier (support vector machine, decision trees, naive Bayes, etc.), it may not be fine-tuned. According to Ridzwan and Yusoff, 2023, decision tree is the most accurate predictor algorithm compared to other machine learning models, whereas in this study, although decision tree gave very high rates, it was an algorithm that suffered from overfitting and difficulty in classifying small numbers of data.

When the relationship between depth and magnitude is analyzed, it is revealed with a statistically significance result that as the depth increases, the magnitude of the earthquakes also increases. The results of the chi-square test show that there is a strong non-random relationship between depth and magnitude. In particular, earthquakes occurring at shallow depths (0-20 km) are generally smaller, while deeper earthquakes have larger magnitudes.

When the earthquake risks in the region are analyzed, it is seen that İzmir province and its surroundings have a high risk especially due to the active fault lines located on the coast. According to Polat et al., 2009, in addition to the study that the earthquakes experienced in Izmir are generally around the Gulf of Izmir and that these areas are in danger, the study conducted by considering the faults located in the sea has contributed to the study on this subject. The study conducted by Tepe et al., 2021 provided valuable information on the intensity distribution of past earthquakes. In the study, it was emphasized that the earthquakes occurring on the Izmir Fault were destructive. With the revision of the study, it has been determined that the Gülbahçe Fault Zone (GLFZ), Mitilini Fault (MTF) and Samos Fault (SF) also have destructive effects. In addition to this contribution, the distribution of earthquakes between 1900 and 2024 has been revised and analyzed. In addition, the study maintains its importance since there is no study that has been conducted by applying a decision tree model in magnitude estimation based on the distribution in the context of Izmir province. Faults in the seafloor (north of Lesbos, south of Lesbos, north of Samos) pose a significant threat to the settlements in and around Izmir. They pose not only earthquake risk but also secondary hazards such as tsunamis for the population living in coastal areas.

In conclusion, the model needs to be strengthened by using more data to predict the magnitude of future earthquakes in Izmir. The lack of data, especially for large earthquakes, limits the performance of the prediction models. The findings of the study show that Izmir province is a region that should be carefully monitored in terms of seismic risk and emphasize that earthquake measures should focus on coastal areas.

## References

- [1] D. Abdullah and E. D. Putra, "Comparasi edge detection roberts dan morfologi pada deteksi plat nomor kendaraan roda dua," *J. Sci. Appl. Informatics*, vol. 1, no. 3, pp. 66–69, 2018.
- [2] S. Ahamed and E. G. Daub, "Machine learning approach to earthquake rupture dynamics," arXiv preprint arXiv:1906.06250, 2019.
- [3] F. Ahmed et al., "Earthquake magnitude prediction using machine learning techniques," in *2024 IEEE Int. Conf. Interdiscip. Approaches Technol. Manag. Soc. Innov. (IATMSI)*, 2024, vol. 2, pp. 1–5.
- [4] A. Ardakani and V. Kohestani, "Evaluation of liquefaction potential based on CPT results using C4.5 decision tree," *J. AI Data Mining*, vol. 3, pp. 85–92, 2015.
- [5] K. M. Asim, A. Idris, F. Martínez-Álvarez, and T. Iqbal, "Short term earthquake prediction in Hindukush region using tree based ensemble learning," in *2016 Int. Conf. Front. Inf. Technol. (FIT)*, 2016, pp. 365–370.
- [6] K. M. Asim, F. Martínez-Álvarez, A. Basit, and T. Iqbal, "Earthquake magnitude prediction in Hindukush region using machine learning techniques," *Nat. Hazards*, vol. 85, pp. 471–486, 2017.
- [7] Y. Cai, M. L. Shyu, Y. X. Tu, Y. T. Teng, and X. X. Hu, "Anomaly detection of earthquake precursor data using long short-term memory networks," *Appl. Geophys.*, vol. 16, pp. 257–266, 2019.
- [8] M. Cassel, "Machine learning and the construction of a seismic attribute-seismic facies analysis data base," M.S. thesis, Univ. Oklahoma, 2018. [Online]. Available: <https://shareok.org/>
- [9] R. Chen et al., "Rigorous assessment and integration of the sequence and structure based features to predict hot spots," *BMC Bioinformatics*, vol. 12, pp. 1–14, 2011.
- [10] R. L. Church, "Geographical information systems and location science," *Comput. Oper. Res.*, vol. 29, no. 6, pp. 541–562, 2002.
- [11] İ. Cürebal and E. Özşahin, *Harita Bilgisi (Bilgisayar Uygulamalı Tasarım ve Analiz)*. Bursa, Türkiye: Ekin Basın Yayın Dağıtım, 2022.
- [12] K. Demertzis, K. Kostinakis, K. Morfidis, and L. Iliadis, "A comparative evaluation of machine learning algorithms for the prediction of R/C buildings' seismic damage," arXiv preprint arXiv:2203.13449, 2022.
- [13] A. Doğru, E. Görgün, H. Ozener, and B. Aktuğ, "Geodetic and seismological investigation of crustal deformation near Izmir (Western Anatolia)," *J. Asian Earth Sci.*, vol. 82, pp. 21–31, 2014.

- [14] T. M. Franke, T. Ho, and C. A. Christie, "The chi-square test: Often used and more often misinterpreted," *Am. J. Eval.*, vol. 33, no. 3, pp. 448–458, 2012.
- [15] E. Gok and O. Polat, "An assessment of the microseismic activity and focal mechanisms of the Izmir (Smyrna) area from a new local network (IzmirNET)," *Tectonophysics*, vol. 635, pp. 154–164, 2014.
- [16] S. Goswami, S. Chakraborty, S. Ghosh, A. Chakrabarti, and B. Chakraborty, "A review on application of data mining techniques to combat natural disasters," *Ain Shams Eng. J.*, vol. 9, no. 3, pp. 365–378, 2018.
- [17] J. Han, J. Kim, S. Park, S. Son, and M. Ryu, "Seismic vulnerability assessment and mapping of Gyeongju, South Korea using frequency ratio, decision tree, and random forest," *Sustainability*, vol. 12, no. 18, p. 7787, 2020.
- [18] C. Jiang, X. Wei, X. Cui, and D. You, "Application of support vector machine to synthetic earthquake prediction," *Earthq. Sci.*, vol. 22, pp. 315–320, 2009.
- [19] A. Karbassi, B. Mohebi, S. Rezaee, and P. Lestuzzi, "Damage prediction for regular reinforced concrete buildings using the decision tree algorithm," *Comput. Struct.*, vol. 130, pp. 46–56, 2014.
- [20] G. Q. King, "Geography and GIS technology," *J. Geogr.*, vol. 90, no. 2, pp. 66–72, 1991.
- [21] S. B. Kotsiantis, "Decision trees: A recent overview," *Artif. Intell. Rev.*, vol. 39, pp. 261–283, 2013.
- [22] A. Li and L. Kang, "Knn-based modeling and its application in aftershock prediction," in *2009 Int. Asia Symp. Intell. Interact. Affect. Comput.*, 2009, pp. 83–86.
- [23] G. Lü et al., "Reflections and speculations on the progress in geographic information systems (GIS): A geographic perspective," *Int. J. Geogr. Inf. Sci.*, vol. 33, no. 2, pp. 346–367, 2019.
- [24] A. Mignan and M. Broccardo, "Neural network applications in earthquake prediction (1994–2019): Meta-analytic and statistical insights on their limitations," *Seismol. Res. Lett.*, vol. 91, no. 4, pp. 2330–2342, 2020.
- [25] O. Onat and H. Tanyıldızı, "Machine learning-based estimation of the out-of-plane displacement of brick infill exposed to earthquake shaking," *Eng. Appl. Artif. Intell.*, vol. 136, p. 109007, 2024.
- [26] G. V. Otari and R. V. Kulkarni, "A review of application of data mining in earthquake prediction," *Int. J. Comput. Sci. Inf. Technol.*, vol. 3, no. 2, pp. 3570–3574, 2012.
- [27] O. Polat et al., "IzmirNet: A strong-motion network in metropolitan Izmir, Western Anatolia, Türkiye" *Seismol. Res. Lett.*, vol. 80, pp. 831–838, 2009.
- [28] R. Rana and R. Singhal, "Chi-square test and its application in hypothesis testing," *J. Primary Care Specialties*, vol. 1, no. 1, pp. 69–71, 2015.
- [29] N. S. M. Ridzwan and S. H. M. Yusoff, "Machine learning for earthquake prediction: A review (2017–2021)," *Earth Sci. Inform.*, vol. 16, no. 2, pp. 1133–1149, 2023.
- [30] M. Senkaya, A. Silahtar, E. F. Erkan, and H. Karaaslan, "Prediction of local site influence on seismic vulnerability using machine learning: A study of the 6 February 2023 Türkiye earthquakes," *Eng. Geol.*, p. 107605, 2024.
- [31] I. Sikder and T. Munakata, "Application of rough set and decision tree for characterization of premonitory factors of low seismic activity," *Expert Syst. Appl.*, vol. 36, pp. 102–110, 2009.
- [32] M. Sokolova, N. Japkowicz, and S. Szpakowicz, "Beyond accuracy, F-score and ROC: A family of discriminant measures for performance evaluation," in *Australas. Joint Conf. Artif. Intell.*, pp. 1015–1021, 2006.
- [33] M. J. Somodevilla, A. B. Priego, E. Castillo, I. H. Pineda, D. Vilarinho, and A. Nava, "Decision support system for seismic risks," *J. Comput. Sci. Technol.*, vol. 12, 2012.
- [34] Ç. Tepe et al., "Updated historical earthquake catalog of İzmir region (western Anatolia) and its importance for the determination of seismogenic source," *Turk. J. Earth Sci.*, vol. 30, no. 8, pp. 779–805, 2021.
- [35] N. S. Turhan, "Karl Pearson's Chi-Square Tests," *Educ. Res. Rev.*, vol. 16, no. 9, pp. 575–580, 2020.
- [36] P. Xiong et al., "Towards advancing the earthquake forecasting by machine learning of satellite data," *Sci. Total Environ.*, vol. 771, p. 145256, 2021.
- [37] C. E. Yavas, L. Chen, C. Kadlec, and Y. Ji, "Predictive modeling of earthquakes in Los Angeles with machine learning and neural networks," *IEEE Access*, 2024.
- [38] E. ZeeAbrahamsen and J. Haberman, "Correcting 'confusability regions' in face morphs," *Behav. Res. Methods*, vol. 50, pp. 1686–1693, 2018.
- [39] J. Zhang, "Dive into decision trees and forests: A theoretical demonstration," *arXiv preprint arXiv:2101.08656*, 2021.



# Optimization of Thermal Management for Cooling System of Power Electronics Modules Consisting Insulated-Gate Bipolar Transistor Using Neuro-Regression Analysis and Non-Traditional Algorithms

Melih SAVRAN<sup>1,\*</sup>, Ece Nur YÜNCÜ<sup>2</sup>, Levent AYDIN<sup>3</sup>

## Abstract

Thermal management and extreme temperatures critically influence the performance of power electronics systems, especially those utilizing Insulated-Gate Bipolar Transistors (IGBTs) and diode components. Various parameters govern the cooling efficiency of these systems. In this study, the IGBT temperature was selected as the objective function. To achieve temperature minimization, optimum values of design variables: coolant flow rate (L/min), distance from the vortex generator (mm), height ( $\mu\text{mm}$ ), and width of the first pin-fin ( $\mu\text{mm}$ ), and distance of the vortex generator from the surface ( $\mu\text{mm}$ ) were determined. The mathematical modeling process employed Neuro-Regression analysis. The prediction performance of proposed 14 different regression models was evaluated using  $R^2$  Training,  $R^2$  Testing,  $R^2$  Validation indexes and boundedness check criteria. Differential Evolution, Nelder Mead, Simulated Annealing, and Random Search algorithms were applied to minimize IGBT temperature. The First Order Logarithmic Nonlinear (FOLN) model emerged as the most successful, achieving a minimum temperature lower than the experimental dataset given in literature. The results indicate a 12 % reduction in the minimum IGBT temperature.

**Keywords:** Optimization, neuro-regression analysis, thermal management, IGBT, cooling system

## 1. Introduction

Effective thermal management systems are essential for the practical use of lithium-ion battery packs. Air cooling alone is insufficient to maintain battery pack temperatures within a safe operating range under high-stress conditions without substantial fan power consumption [1, 2]. Insulated gate bipolar transistor (IGBT) modules have recently become prevalent in various industries, notably in high-power converters for wind turbines, trains, and HVDC systems [3]. Thermal management, encompassing battery temperature regulation and air conditioning cabinet, poses a significant challenge for electric vehicles (EVs), where traditional engines and oil tanks are replaced by electric motors and battery assemblies [4]. Optimizing thermal management is critical for the performance of IGBT-based power modules in hybrid electric vehicles [5]. Jun He et al. have studied the thermal design and assessment of IGBT power modules under both transient and steady-state conditions, suggesting that optimizing wire bond configurations and bonding pad positions can significantly reduce temperature gradients and peak temperatures on the IGBT surface [6]. Thermal resistance ( $R_{th}$ ), defined as the ratio of the temperature difference between the heat output and input ends to the power, is a crucial parameter for IGBT modules and an important measure of their heat dissipation efficiency [7].

Efficient thermal management not only enhances performance but also enables the miniaturization of power electronics equipment [8]. In the application of IGBTs, particularly in high-voltage heater systems, the challenge lies in managing the additional heat generated by the heating elements applied via plasma deposition technology. This makes the thermal management requirements even more stringent. In the application of an IGBT, it is crucial to analyze the heat generation and transfer behavior to minimize chip temperature. Within a high voltage heater system, the IGBT is secured to the heat exchanger using bolts, while the heating element is directly applied to the heat exchanger using plasma deposition technology. As a result, during IGBT operation, in addition to the heat

\*Corresponding author

Melih SAVRAN\*; İzmir Kâtip Çelebi University, Department of Mechanical Engineering, İzmir, Türkiye; e-mail: [mlhsvrn@gmail.com](mailto:mlhsvrn@gmail.com)



0000-0001-8343-1073

Ece Nur YÜNCÜ; İzmir Kâtip Çelebi University, Department of Electrical and Electronics Engineering, İzmir, Türkiye; e-mail: [ecenuryuncu79@gmail.com](mailto:ecenuryuncu79@gmail.com)



0009-0002-3195-124X

Levent AYDIN; İzmir Kâtip Çelebi University, Department of Mechanical Engineering, İzmir, Türkiye; e-mail: [leventaydinn@gmail.com](mailto:leventaydinn@gmail.com)



0000-0003-0483-0071

produced by the chip itself, heat is also transferred from the heating elements, leading to more stringent thermal management requirements for the IGBT in the high voltage heater system. Current cooling methods for IGBTs are inadequate for maintaining a safe temperature range under the high-power heating conditions of the high voltage heater system, posing a significant risk to system reliability. So far, many studies on IGBT cooling have concentrated on designing cooling structures to solve the problem of effective heat dissipation for IGBTs [9].

Rao et al. optimized a plate-fin heat exchanger by minimizing the total number of entropy generation units for a specific heat duty requirement within given space constraints, reducing the total volume, and lowering the total annual cost [10]. Lee et al. utilized a multi-objective genetic algorithm combined with surrogate modeling techniques to maximize heat transfer and minimize pressure drop in a heat exchanger [11]. Mishra et al. used GA for optimal design of plate-fin heat exchangers [12, 13]. Some authors used particle swarm optimization for rolling fin-tube heat exchanger optimization [14].

The main goal of this study is to optimize the cooling of a power electronic system that consists of an IGBT and a diode, along with the associated connections and joints, by leveraging the data from Pourfattah Farzad et al. [15]. This study introduces a novel approach to address the shortcomings in the design, modeling, and optimization of the thermal management for cooling systems of power electronics modules. The proposed method employs multiple nonlinear neuro-regression analyses, integrating artificial neural networks (ANN), regression analysis, and stochastic optimization techniques to achieve suitable designs that meet desired specifications. This approach allows for diverse alternative mathematical models, transcending traditional limitations to specific polynomial forms or activation functions such as sigmoid, unit step, and hyperbolic tangent.

Furthermore, model assessment incorporates both the  $R^2$  value and a boundedness check criterion, which provides a more holistic evaluation of model reliability. The boundedness check is vital for developing dependable mathematical models, as all engineering parameters must be finite. Realistic modeling in engineering systems necessitates that models are bounded within specified parameter intervals; thus, verifying this boundedness prior to optimization is essential. In contrast to modeling techniques reliant on artificial neural networks, this method circumvents the need for fine-tuning parameters such as the number of neurons and hidden layers, which are often adjusted to enhance ANN-based models. This modeling approach significantly enhances the thermal management for cooling systems of power electronics module in the existing literature. Algorithms; Differential Evolution, Nelder Mead, Simulated Annealing, and Random Search are employed to identify the optimal design parameters and IGBT temperature for efficient thermal management.

## 2. Materials and Methods

### 2.1 Mathematical modelling

In the modeling stage, a combined method of regression analysis and artificial neural networks are utilized to enhance the accuracy of predictions. The dataset is divided into three parts: 80% for training, 15% for testing, and 5% for validation. During training, various regression models outlined in Table 1 were employed to minimize the disparity between experimental and predicted values. In the testing and validation phase, the objective was to generate prediction outcomes while mitigating inconsistencies among regression models. Evaluating the boundedness of the models was crucial for assessing their realism. Following the selection of suitable models based on  $R^2$  index for training, testing, and validation, the maximum and minimum values for each design parameter were computed. In q. (1),  $R^2$  is coefficient of determination that indicates how well the data fit a regression model.  $R^2$  value range from 0 to 1. As  $R^2$  value is closer to 1 it indicates that there is a good fit with that model. SSE stands for ‘Sum of Squared Errors’ and measures the total deviation of the observed values from the predicted values produced by the regression model. SST stands for ‘sum of squares total’ and measures the total deviation of the observed values from their mean.

$$R^2 = 1 - \frac{SSE}{SST} \quad (1)$$

**Table 1.** Multiple regression model types including linear, quadratic, trigonometric, logarithmic, and their rational forms [16]

Model Name	Nomenclature	Formula
MultipleLinear	L	$(a_0 + a_1 x_1 + a_2 x_2 + a_3 x_3 + a_4 x_4 + a_5 x_5)$
Multiple Linear Rational	LR	$(a_0 + a_1 x_1 + a_2 x_2 + a_3 x_3 + a_4 x_4 + a_5 x_5) / (b_0 + b_1 x_1 + b_2 x_2 + b_3 x_3 + b_4 x_4 + b_5 x_5)$
Second-Order Multiple Nonlinear	SON	$(a_0 + a_1 x_1 + a_2 x_2 + a_3 x_3 + a_4 x_4 + a_5 x_5 + a_6 x_1 x_1 + a_7 x_2 x_2 + a_8 x_3 x_3 + a_9 x_4 x_4 + a_{10} x_5 x_5 + a_{11} x_1 x_2 + a_{12} x_1 x_3 + a_{13} x_1 x_4 + a_{14} x_1 x_5 + a_{15} x_2 x_3 + a_{16} x_2 x_4 + a_{17} x_2 x_5 + a_{18} x_3 x_4 + a_{19} x_3 x_5 + a_{20} x_4 x_5)$
Second-Order Multiple Nonlinear Rational	SONR	$(a_0 + a_1 x_1 + a_2 x_2 + a_3 x_3 + a_4 x_4 + a_5 x_5 + a_6 x_1 x_1 + a_7 x_2 x_2 + a_8 x_3 x_3 + a_9 x_4 x_4 + a_{10} x_5 x_5 + a_{11} x_1 x_2 + a_{12} x_1 x_3 + a_{13} x_1 x_4 + a_{14} x_1 x_5 + a_{15} x_2 x_3 + a_{16} x_2 x_4 + a_{17} x_2 x_5 + a_{18} x_3 x_4 + a_{19} x_3 x_5 + a_{20} x_4 x_5) / (b_0 + b_1 x_1 + b_2 x_2 + b_3 x_3 + b_4 x_4 + b_5 x_5 + b_6 x_1 x_1 + b_7 x_2 x_2 + b_8 x_3 x_3 + b_9 x_4 x_4 + b_{10} x_5 x_5 + b_{11} x_1 x_2 + b_{12} x_1 x_3 + b_{13} x_1 x_4 + b_{14} x_1 x_5 + b_{15} x_2 x_3 + b_{16} x_2 x_4 + b_{17} x_2 x_5 + b_{18} x_3 x_4 + b_{19} x_3 x_5 + b_{20} x_4 x_5)$
First-Order Trigonometric Multiple Nonlinear	FOTN	$(a_0 + a_1 \sin[x_1] + a_2 \sin[x_2] + a_3 \sin[x_3] + a_4 \sin[x_4] + a_5 \sin[x_5] + a_6 \cos[x_1] + a_7 \cos[x_2] + a_8 \cos[x_3] + a_9 \cos[x_4] + a_{10} \cos[x_5])$
First-Order Trigonometric Multiple Nonlinear Rational	FOTNR	$(a_0 + a_1 \sin[x_1] + a_2 \sin[x_2] + a_3 \sin[x_3] + a_4 \sin[x_4] + a_5 \sin[x_5] + a_6 \cos[x_1] + a_7 \cos[x_2] + a_8 \cos[x_3] + a_9 \cos[x_4] + a_{10} \cos[x_5]) / (b_0 + b_1 \sin[x_1] + b_2 \sin[x_2] + b_3 \sin[x_3] + b_4 \sin[x_4] + b_5 \sin[x_5] + b_6 \cos[x_1] + b_7 \cos[x_2] + b_8 \cos[x_3] + b_9 \cos[x_4] + b_{10} \cos[x_5])$
Second-Order Trigonometric Multiple Nonlinear	SOTN	$(a_0 + a_1 \sin[x_1] + a_2 \sin[x_2] + a_3 \sin[x_3] + a_4 \sin[x_4] + a_5 \sin[x_5] + a_6 \cos[x_1] + a_7 \cos[x_2] + a_8 \cos[x_3] + a_9 \cos[x_4] + a_{10} \cos[x_5] + a_{11} \sin[x_1] \sin[x_1] + a_{12} \sin[x_2] \sin[x_2] + a_{13} \sin[x_3] \sin[x_3] + a_{14} \sin[x_4] \sin[x_4] + a_{15} \sin[x_5] \sin[x_5] + a_{16} \cos[x_1] \cos[x_1] + a_{17} \cos[x_2] \cos[x_2] + a_{18} \cos[x_3] \cos[x_3] + a_{19} \cos[x_4] \cos[x_4] + a_{20} \cos[x_5] \cos[x_5])$
Second-Order Trigonometric Multiple Nonlinear Rational	SOTNR	$(a_0 + a_1 \sin[x_1] + a_2 \sin[x_2] + a_3 \sin[x_3] + a_4 \sin[x_4] + a_5 \sin[x_5] + a_6 \cos[x_1] + a_7 \cos[x_2] + a_8 \cos[x_3] + a_9 \cos[x_4] + a_{10} \cos[x_5] + a_{11} \sin[x_1] \sin[x_1] + a_{12} \sin[x_2] \sin[x_2] + a_{13} \sin[x_3] \sin[x_3] + a_{14} \sin[x_4] \sin[x_4] + a_{15} \sin[x_5] \sin[x_5] + a_{16} \cos[x_1] \cos[x_1] + a_{17} \cos[x_2] \cos[x_2] + a_{18} \cos[x_3] \cos[x_3] + a_{19} \cos[x_4] \cos[x_4] + a_{20} \cos[x_5] \cos[x_5]) / (b_0 + b_1 \sin[x_1] + b_2 \sin[x_2] + b_3 \sin[x_3] + b_4 \sin[x_4] + b_5 \sin[x_5] + b_6 \cos[x_1] + b_7 \cos[x_2] + b_8 \cos[x_3] + b_9 \cos[x_4] + b_{10} \cos[x_5] + b_{11} \sin[x_1] \sin[x_1] + b_{12} \sin[x_2] \sin[x_2] + b_{13} \sin[x_3] \sin[x_3] + b_{14} \sin[x_4] \sin[x_4] + b_{15} \sin[x_5] \sin[x_5] + b_{16} \cos[x_1] \cos[x_1] + b_{17} \cos[x_2] \cos[x_2] + b_{18} \cos[x_3] \cos[x_3] + b_{19} \cos[x_4] \cos[x_4] + b_{20} \cos[x_5] \cos[x_5])$
First-Order Logarithmic Multiple Nonlinear	FOLN	$(a_0 + a_1 \log[x_1] + a_2 \log[x_2] + a_3 \log[x_3] + a_4 \log[x_4] + a_5 \log[x_5])$
First-Order Logarithmic Multiple Nonlinear Rational	FOLNR	$(a_0 + a_1 \log[x_1] + a_2 \log[x_2] + a_3 \log[x_3] + a_4 \log[x_4] + a_5 \log[x_5]) / (b_0 + b_1 \log[x_1] + b_2 \log[x_2] + b_3 \log[x_3] + b_4 \log[x_4] + b_5 \log[x_5])$
Second-Order Logarithmic Multiple Nonlinear	SOLN	$(a_0 + a_1 \log[x_1] + a_2 \log[x_2] + a_3 \log[x_3] + a_4 \log[x_4] + a_5 \log[x_5] + a_6 \log[x_1] \log[x_1] + a_7 \log[x_2] \log[x_2] + a_8 \log[x_3] \log[x_3] + a_9 \log[x_4] \log[x_4] + a_{10} \log[x_5] \log[x_5] + a_{11} \log[x_1] \log[x_2] + a_{12} \log[x_1] \log[x_3] + a_{13} \log[x_1] \log[x_4] + a_{14} \log[x_1] \log[x_5] + a_{15} \log[x_2] \log[x_3] + a_{16} \log[x_2] \log[x_4] + a_{17} \log[x_2] \log[x_5] + a_{18} \log[x_3] \log[x_4] + a_{19} \log[x_3] \log[x_5] + a_{20} \log[x_4] \log[x_5])$
Second-Order Logarithmic Multiple Nonlinear	SOLNR	$(a_0 + a_1 \log[x_1] + a_2 \log[x_2] + a_3 \log[x_3] + a_4 \log[x_4] + a_5 \log[x_5] + a_6 \log[x_1] \log[x_1] + a_7 \log[x_2] \log[x_2] + a_8 \log[x_3] \log[x_3] + a_9 \log[x_4] \log[x_4] + a_{10} \log[x_5] \log[x_5] + a_{11} \log[x_1] \log[x_2] + a_{12} \log[x_1] \log[x_3] + a_{13} \log[x_1] \log[x_4] + a_{14} \log[x_1] \log[x_5] + a_{15} \log[x_2] \log[x_3] + a_{16} \log[x_2] \log[x_4] + a_{17} \log[x_2] \log[x_5] + a_{18} \log[x_3] \log[x_4] + a_{19} \log[x_3] \log[x_5] + a_{20} \log[x_4] \log[x_5]) / (b_0 + b_1 \log[x_1] + b_2 \log[x_2] + b_3 \log[x_3] + b_4 \log[x_4] + b_5 \log[x_5] + b_6 \log[x_1] \log[x_1] + b_7 \log[x_2] \log[x_2] + b_8 \log[x_3] \log[x_3] + b_9 \log[x_4] \log[x_4] + b_{10} \log[x_5] \log[x_5] + b_{11} \log[x_1] \log[x_2] + b_{12} \log[x_1] \log[x_3] + b_{13} \log[x_1] \log[x_4] + b_{14} \log[x_1] \log[x_5] + b_{15} \log[x_2] \log[x_3] + b_{16} \log[x_2] \log[x_4] + b_{17} \log[x_2] \log[x_5] + b_{18} \log[x_3] \log[x_4] + b_{19} \log[x_3] \log[x_5] + b_{20} \log[x_4] \log[x_5])$

Two new hybrid regression models are also proposed in this study. These regression model formulas are given in Table 2.

**Table 2.** Hybrid regression model types

Model Name	Nomenclature	Formula
Hybrid	H(FOLN+SON)	$(a_0 + a_1 \text{Log}[x_1] + a_2 \text{Log}[x_2] + a_3 \text{Log}[x_3] + a_4 \text{Log}[x_4] + a_5 \text{Log}[x_5]$ $+ a_6 + a_7 x_1 + a_8 x_2 + a_9 x_3 + a_{10} x_4 + a_{11} x_5$ $+ a_{12} x_1 x_1 + a_{13} x_2 x_2 + a_{14} x_3 x_3 + a_{15} x_4 x_4$ $+ a_{16} x_5 x_5 + a_{17} x_1 x_2 + a_{18} x_1 x_3 + a_{19} x_1 x_4$ $+ a_{20} x_1 x_5 + a_{21} x_2 x_3 + a_{22} x_2 x_4 + a_{23} x_2 x_5$ $+ a_{24} x_3 x_4 + a_{25} x_3 x_5 + a_{26} x_4 x_5)$
Hybrid	H(FOLN*L)	$(a_0 + a_1 \text{Log}[x_1] + a_2 \text{Log}[x_2] + a_3 \text{Log}[x_3] + a_4 \text{Log}[x_4] + a_5 \text{Log}[x_5])$ $* (a_6 + a_7 x_1 + a_8 x_2 + a_9 x_3 + a_{10} x_4 + a_{11} x_5)$

## 2.2 Optimization

Optimization involves refining a system or process to achieve the best possible outcome. This process entails adjusting input variables to minimize or maximize the output of a function, often referred to as the cost function, objective function, or fitness function. The goal is to optimize these inputs to achieve the best possible performance of the system [17].

### 2.2.1. Differential evolution

Differential Evolution (DE) is a population-based optimization algorithm particularly effective for solving complex, high-dimensional optimization problems. DE begins by initializing a population of candidate solutions, iteratively refining them across generations by exploiting differences (differentials) between solutions within the population. In each generation, new candidate solutions are generated through a mutation process, which typically involves selecting three random individuals to create differential vectors. These vectors are then combined with an existing solution to propose a new candidate. A crossover operation further enhances solution diversity, while a selection process ensures that only improved solutions are retained. One of DE's key strengths is its ability to reach globally optimal solutions without requiring gradient information, making it highly suitable for applications in engineering and scientific research. Its relatively low sensitivity to parameter settings also contributes to its widespread use in various optimization tasks [18].

### 2.2.2. Nelder-mead

The Nelder-Mead algorithm is a widely used direct search optimization technique that operates without the need for gradient information, making it suitable for optimizing non-differentiable or complex objective functions. The algorithm maintains a simplex—a geometric shape formed by  $n+1$  vertices in an  $n$ -dimensional space—and iteratively refines it using four main operations: reflection, expansion, contraction, and shrinkage. Through these operations, the simplex adjusts its shape, size, and orientation dynamically, allowing it to navigate the objective function landscape effectively and converge towards a local optimum. By adapting to the contours of the objective function, the Nelder-Mead algorithm demonstrates flexibility and robustness, making it particularly valuable for challenging optimization tasks where traditional gradient-based methods may be infeasible or ineffective [19].

### 2.2.3. Random search

The Random Search algorithm is a stochastic optimization technique that contrasts with deterministic methods, such as Branch and Bound or Interval Analysis, by relying on random sampling rather than systematic exploration of the search space. Unlike gradient-based or small-step methods that risk converging to local optima, Random Search samples candidate solutions across the entire search domain, thereby increasing its likelihood of identifying a global optimum, especially in multimodal objective functions. This characteristic makes Random Search particularly advantageous for problems where the objective function contains multiple peaks or valleys. The algorithm's simplicity and adaptability allow it to explore complex search landscapes without gradient information, though its efficiency can be enhanced by combining it with local refinement strategies to ensure both global exploration and local exploitation of high-potential regions [20].

### 2.2.4. Simulated annealing

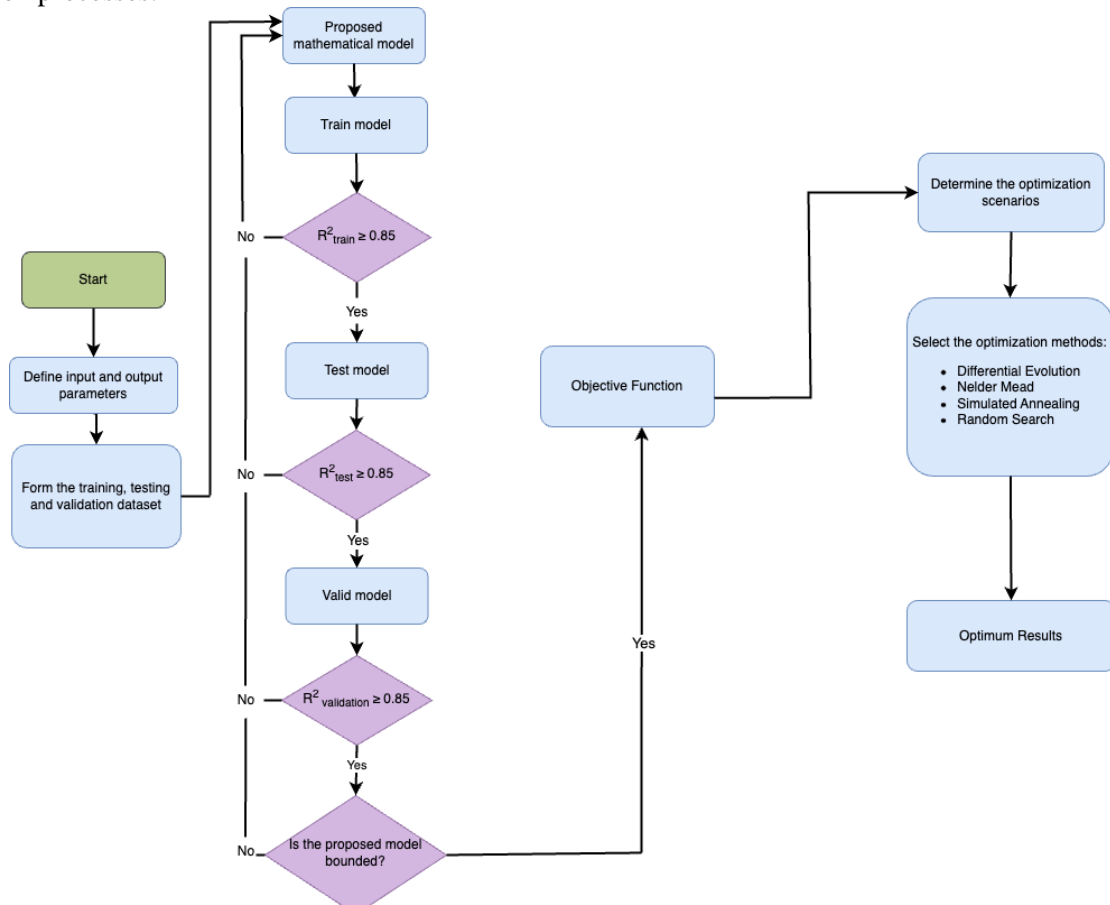
Simulated Annealing (SA) is a widely adopted optimization technique within random search methods, inspired by the physical annealing process. In this process, a metal is heated to a high temperature and then gradually cooled, allowing its atomic structure to settle into a state of lower energy, resulting in a tougher and more stable material. In the context of optimization, the SA algorithm mimics this annealing process to enable solutions to escape local minima and explore the search space more broadly in pursuit of a global optimum. Initially, the algorithm accepts a wide range of solutions, including those that may increase the objective function, which helps it to traverse diverse regions of the search landscape. As the "temperature" parameter decreases, the acceptance of higher-energy solutions becomes less likely, guiding the algorithm towards a stable and optimal solution. This dynamic makes SA particularly effective for solving complex, multimodal optimization problems, as it balances global exploration and local refinement [21].

### 2.3. Problem definition

The main aim of this study is to identify the optimal design parameters to minimize the temperature of the IGBT. The study involves several steps:

- Data Selection and determination of design variables and output parameters: Data was taken from the reference study conducted by Pourfattah Farzad et al. [15]. The design variables included the coolant flow rate, the height and width of the first pin-fin attached to the heatsink, the distance from the vortex generator, the distance from the coolant path surface to the vortex generator. The output parameter is selected as IGBT temperature.
- Model Selection: Fourteen regression models were utilized, and their validity was assessed by checking the  $R^2$  values and boundedness criteria. Models are considered successful when they achieve  $R^2$  values greater than 0.85 and have realistic maximum and minimum outputs for engineering applications.
- Optimization: The model, which successfully met the model assessment and boundedness control criteria, was optimized using four optimization methods (DE, NM, RS, SA) to obtain optimal results, which were then compared with one another.

The flow chart in Figure 1 provides a detailed description of the steps taken in the mathematical modeling and optimization processes.



**Figure 1.** The flowchart regarding mathematical modeling and optimization process

### 2.3.1 Optimization scenarios

Three scenarios with varying constraints on the design parameters were established to determine the optimal solution.

#### Scenario 1

In the first scenario, the search space was continuous. The intervals for the design variables are as follows:  $1.203 \leq x_1$  (L/min)  $\leq 4.497$ ,  $505.56 \leq x_2$  ( $\mu\text{mm}$ )  $\leq 783.33$ ,  $2.215 \leq x_3$  (mm)  $\leq 2.985$ ,  $512.96 \leq x_4$  ( $\mu\text{mm}$ )  $\leq 1187.04$ ,  $0.46 \leq x_5$  ( $\mu\text{mm}$ )  $\leq 251.91$

#### Scenario 2

For this scenario, the search space of some design variables ( $x_1$ ,  $x_2$ ,  $x_4$ ,  $x_5$ ) was considered as integer. The intervals for the design variables are as follows:  $1.203 \leq x_1$  (L/min)  $\leq 4.497$ ,  $505.56 \leq x_2$  ( $\mu\text{mm}$ )  $\leq 783.33$ ,  $2.215 \leq x_3$  (mm)  $\leq 2.985$ ,  $512.96 \leq x_4$  ( $\mu\text{mm}$ )  $\leq 1187.04$ ,  $0.46 \leq x_5$  ( $\mu\text{mm}$ )  $\leq 251.91$   $\{x_1, x_2, x_4, x_5\} \in \text{Integers}$

#### Scenario 3

In the third scenario, all design parameters were taken as only certain specific values determined in the experimental set. For this study, each design variable had 24 different levels. Due to implementing regression models for 24 different levels taking time regarding optimization, each design parameter's level is chosen as four specified values: minimum, maximum, middle, and the design parameters that performed the best results in the experimental study. The design parameters and their level values are; coolant flow rate ( $x_1$ )  $\in \{1.203, 2.850, 3.103, 4.497\}$ , height of the first pin-fin ( $x_2$ )  $\in \{650.00, 772.22, 783.33, 1505.56\}$ , distance from the vortex generator ( $x_3$ )  $\in \{2.215, 2.600, 2.659, 2.985\}$ , width of the first pin-fin ( $x_4$ )  $\in \{512.96, 850.00, 1005.56, 1187.04\}$ , distance of the vortex generator from surface ( $x_5$ )  $\in \{0.46, 113.89, 134.26, 251.91\}$ .

**Table 3.** Design points regarding with input and output parameters [15]

Run Order	coolant flow rate ( $x_1$ ) (L/min)	height of the first pin-fin ( $x_2$ ) ( $\mu\text{mm}$ )	distance from the vortex generator ( $x_3$ ) (mm)	width of the first pin-fin ( $x_4$ ) ( $\mu\text{mm}$ )	distance of the vortex generator from surface ( $x_5$ ) ( $\mu\text{mm}$ )	IGBT temperature ( $^{\circ}\text{C}$ )
1	2.723	650.00	2.748	642.6	81.67	82.68
2	1.203	527.78	2.748	979.63	47.22	106.65
3	2.977	683.33	2.778	953.70	80.93	80.42
4	1.457	672.22	2.896	746.30	40.06	98.18
5	4.370	605.56	2.215	1109.26	83.27	75.3
6	2.850	738.89	2.837	512.96	147.96	76.73
7	4.243	594.44	2.393	824.07	220.80	74.6
8	2.090	516.67	2.274	850.00	134.26	92.22
9	1.837	750.00	2.926	1187.04	119.44	85.66
10	2.343	761.11	2.511	668.52	76.98	80.65
11	1.710	583.33	2.867	1057.41	217.59	95.13
12	2.597	783.33	2.363	927.78	0.46	78.9
13	3.483	627.78	2.452	1083.33	161.30	78.28
14	2.217	572.22	2.481	798.15	251.91	88.34
15	3.863	550.00	2.719	875.93	35.00	77.77
16	3.230	727.78	2.600	564.81	99.81	76.03
17	1.583	638.89	2.689	1161.11	164.51	95.71
18	1.330	705.56	2.630	590.74	55.62	97.9
19	4.117	561.11	2.807	772.22	180.43	76.16
20	4.497	538.89	2.333	616.67	158.83	74.88
21	1.963	505.56	2.956	538.89	82.41	80.72
22	3.610	716.67	2.244	720.37	110.19	74.58
23	3.103	772.22	2.659	1005.56	113.89	74.36
24	3.737	694.44	2.985	1031.48	91.67	76

### 3. Results and Discussion

This study established a mathematical relationship between design parameters (coolant flow rate (x1), height of the first pin-fin (x2), distance from the vortex generator (x3), width of the first pin-fin (x4) and distance of the vortex generator from surface (x5)), and output parameter (IGBT temperature). The goal was to identify the values of these design parameters that minimize the IGBT temperature using the most effective model.

Table 4 presents the performance of various neuro-regression models in terms of their  $R^2$  values (for training, testing, and validation phases) and their boundedness check (maximum and minimum values). The  $R^2$  values during training are notably high across all models, with some achieving values close to 1.0. This suggests that the models exhibit a strong fit to the training data. However, such high values may indicate potential overfitting, where the model may not generalize well to unseen data.

In the testing and validation phase, several models yield negative  $R^2$  values (e.g., LR: -0.541896, SOTNR: -8.17925), implying poor performance and possibly inverse predictions relative to the data trend. Particularly in the SOTNR model, this discrepancy may indicate substantial overfitting.

The maximum and minimum values across the models reveal that some models produce extreme bounds (e.g., the minimum value for SOTNR:  $1.78333 \times 10^{10}$ ), indicating that these models may generate highly varied or extreme outputs. This wide prediction range points to a tendency toward volatility in some models.

The FOLN model demonstrates commendable performance across multiple evaluation criteria, particularly in terms of its  $R^2$  values and boundedness. The model achieves high  $R^2$  values in the training (0.99805), testing (0.996986), and validation phases (0.99921), indicating a consistently strong fit and predictive capability across different data subsets. Such uniformly high  $R^2$  values suggest that the FOLN model not only learns the training data effectively but also generalizes well to unseen data, avoiding overfitting issues commonly observed in other models.

Regarding boundedness, the FOLN model maintains a prediction range with maximum and minimum values of 105.094 and 65.7673, respectively. This bounded range suggests a stable prediction behavior. The FOLN model's boundedness further supports its robustness, as it operates within a controlled range, contrasting with models that exhibit high variance in output values.

Among the models in Table 4, the FOLN model only demonstrates a balance between model fitness and prediction stability. For this reason, it is selected as an objective function in the optimization process to minimize IGBT temperature.

**Table 4.** Result of the Neuro-Regression Models in Terms of  $R^2$  and Boundedness

Model	$R^2$ Training	$R^2$ Testing	$R^2$ Validation	Max	Min
L	0.997574	0.880102	0.952265	105.559	60.8362
LR	0.999749	-0.541896	0.903765	$\infty$	$\infty$
SON	1.	0.384109	-1.01635	125.852	36.4097
SONR	0.999493	0.83447	0.45114	182.503	$-3.5922 \times 10^9$
FOTN	0.999301	0.519079	0.91273	110.303	58.3124
FOTNR	0.999854	-1.0012	0.83765	$4.64839 \times 10^6$	$3.82897 \times 10^6$
SOTN	0.999845	-0.268708	-0.0235936	108.566	44.027
SOTNR	0.999936	-8.17925	-16.0617	$1.70706 \times 10^{15}$	$1.78333 \times 10^{10}$
<b>FOLN</b>	<u>0.99805</u>	<u>0.996986</u>	<u>0.99921</u>	<u>105.094</u>	<u>65.7673</u>
FOLNR	0.999662	-1.85905	-2.08456	$1.95554 \times 10^7$	$2.77389 \times 10^6$
SOLN	1.	-1.39353	-1.38238	339.14	154.896
SOLNR	0.999875	-0.192595	-0.355909	$4.01139 \times 10^7$	33.6077
<b>H (FOLN+SON)</b>	1.	-0.305234	0.314625	118.633	30.6786
<b>H (FOLN*L)</b>	0.999701	0.446418	0.937696	112.363	50.9868

Table 5 presents the results of optimization scenarios for the FOLN model, with a focus on achieving minimum Insulated-Gate Bipolar Transistor (IGBT) temperatures across various optimization algorithms: DE, SA, RS, and NM. The findings indicate the effectiveness and stability of the FOLN model in identifying optimal designs under diverse conditions.

In scenario 1, across all algorithms (MDE, MSA, MRS, MNM), the minimum IGBT temperature achieved is consistent at  $65.7673^\circ\text{C}$ , with the suggested design values for parameters x1 to x5 remaining identical. This outcome indicates a convergence across optimization methods toward a common design that minimizes temperature.

When the search space of some design variables (x1, x2, x4, x5) is considered an integer in scenario 2, slight variations appear between algorithms. For instance, MDE and MSA yield a minimum temperature of 68.4517°C, while MRS and MNM result in slightly higher temperatures of 68.7299°C and 71.365°C, respectively. The recommended parameter values exhibit minor differences by algorithms, suggesting some sensitivity in the model’s design variable recommendations depending on the optimization technique.

In Scenario 3, under all design parameters are taken as only certain specific values determined in the experimental set, four algorithms consistently converge to the minimum IGBT temperature of 65.7673°C with same design parameters. This convergence among the algorithms suggests that the optimal result has been attained. A comparison with the experimental results from the reference study supports this inference. While the experimentally obtained minimum temperature was 74.36°C, the present study achieves a significantly lower minimum temperature of 65.7673°C through modeling and optimization.

In conclusion, this consistency reinforces the FOLN model’s suitability for applications requiring precise thermal management within defined parameter boundaries.

**Table 5.** Results of optimization problems for FOLN model considering minimum IGBT temperature.

Objective Function	Scenario Number	Constrains	Optimization Algorithm	Minimum IGBT Temperature (°C)	Suggested Design
FOLN	1	$1.203 \leq x1 \leq 4.497,$ $505.56 \leq x2 \leq 783.33$ $2.215 \leq x3 \leq 2.985$ $512.96 \leq x4 \leq 1187.04$ $0.46 \leq x5 \leq 251.91$	DE	65.7673	x1 -> 4.497, x2 -> 783.33, x3 -> 2.985, x4 -> 512.96, x5 -> 0.46
			SA	65.7673	x1 -> 4.497, x2 -> 783.33, x3 -> 2.985, x4 -> 512.96, x5 -> 0.46
			RS	65.7673	x1 -> 4.497, x2 -> 783.33, x3 -> 2.985, x4 -> 512.96, x5 -> 0.46
			NM	65.7673	x1 -> 4.497, x2 -> 783.33, x3 -> 2.985, x4 -> 512.96, x5 -> 0.46
	2	$1.203 \leq x1 \leq 4.497,$ $505.56 \leq x2 \leq 783.33$ $2.215 \leq x3 \leq 2.985$ $512.96 \leq x4 \leq 1187.04$ $0.46 \leq x5 \leq 251.91$ {x1, x2, x4, x5} ∈ Integers	DE	68.4517	x1 -> 4, x2 -> 783, x3 -> 2.985, x4 -> 513, x5 -> 1
			SA	68.4517	x1 -> 4, x2 -> 783, x3 -> 2.985, x4 -> 513, x5 -> 1
			RS	68.7299	x1 -> 4, x2 -> 778, x3 -> 2.92473, x4 -> 513, x5 -> 1
			NM	71.365	x1 -> 4, x2 -> 707, x3 -> 2.985, x4 -> 826, x5 -> 211
	3	x1 = 1.203    x1 = 2.850    x1 = 3.103    x1 = 4.497, x2 = 505.56    x2 = 650.00    x2 = 772.22    x2 = 783.33, x3 = 2.215    x3 = 2.600    x3 = 2.659    x3 = 2.985, x4 = 512.96    x4 = 850.00    x4 = 1005.56    x4 = 1187.04, x5 = 0.46    x5 = 113.89    x5 = 134.26    x5 = 251.91	DE	65.7673	x1 -> 4.497, x2 -> 783.33, x3 -> 2.985, x4 -> 512.96, x5 -> 0.46
			SA	65.7673	x1 -> 4.497, x2 -> 783.33, x3 -> 2.985, x4 -> 512.96, x5 -> 0.46
			RS	65.7673	x1 -> 4.497, x2 -> 783.33, x3 -> 2.985, x4 -> 512.96, x5 -> 0.46
			NM	65.7673	x1 -> 4.497, x2 -> 783.33, x3 -> 2.985, x4 -> 512.96, x5 -> 0.46



#### 4. Conclusion

This study highlights the critical role of thermal management in optimizing power electronics systems, specifically those employing IGBT components. By focusing on minimizing the IGBT temperature as the objective function, the optimization framework employed key design variables such as coolant flow rate (x1), height of the first pin-fin (x2), distance from the vortex generator (x3), width of the first pin-fin (x4) and distance of the vortex generator from surface (x5). Standard methods that use limited regression models often ignore nonlinear effects and are ineffective for optimizing thermal management for cooling systems of power electronics modules. This study introduces a new way to model the relation between cooling system design parameters and IGBT temperature by combining artificial neural networks (ANN) with regression techniques. This approach, called neuro-regression, selects the best models from linear, rational, logarithmic, polynomial, trigonometric, and hybrid types based on criteria  $R^2$  and boundedness check. The FOLN neuro-regression model emerged as the most effective in achieving a balance between high predictive accuracy and model boundedness across training, testing and validation datasets.

The results indicate that when the FOLN model was selected as the objective function, the Differential Evolution, Simulated Annealing, Random Search, and Nelder-Mead algorithms found the minimum IGBT temperature to be 65.7673°C. This temperature is significantly lower than the minimum temperature of 74.36°C reported in experimental studies.

This outcome suggests that the FOLN model is particularly well-suited for applications necessitating precise and robust thermal control. Moreover, the consistency observed across different optimization algorithms emphasizes the result's robustness, as each algorithm converged to the same minimum temperature. This convergence validates the model's efficacy for thermal management in power electronics.

Future studies may further explore applying the FOLN model across a broader range of conditions to enhance predictive performance and thermal management strategies in advanced electronics systems.

#### Statements & Declarations

##### Competing Interests

“The authors declare that they have no known competing financial interests or personal relationships that could have appeared to influence the work reported in this paper. The authors have no relevant financial or non-financial interests to disclose.”

##### Conflict of Interest

“The authors declare that they have no conflict of interest.”

##### Author Contribution

Melih Savran, Conceptualization, Methodology, Software, Validation, Formal Analysis, Investigation, Data curation, Project Administration, Writing – Original Draft, Writing – Review & Editing, Visualization; Ece Nur Yüncü, Conceptualization, Software, Validation, Formal Analysis, Investigation, Data curation, Writing – Original Draft, Visualization; Levent Aydın, Conceptualization, Methodology, Software, Validation, Formal Analysis, Investigation, Data curation, Supervision, Project Administration, Writing – Review & Editing;

##### Availability of Data and Materials

The authors confirm that the data supporting the fundings of this study are available within the article.

##### Ethical Approval

All authors have previously approved this paper and judged that there is no ethical infringement.

##### Consent to Participate and Publish

All authors would like to declare that they have approved their participation and consent about the publication in this journal.

## References

- [1] F. Wang, J. Cao, Z. Ling, Z. Zhang, and X. Fang, "Experimental and simulative investigations on a phase change material nano-emulsion-based liquid cooling thermal management system for a lithium-ion battery pack," *Energy*, vol. 207, p. 118215, 2020.
- [2] R. Sabbah, R. Kizilel, J. R. Selman, and S. Al-Hallaj, "Active (air-cooled) vs. passive (phase change material) thermal management of high power lithium-ion packs: Limitation of temperature rise and uniformity of temperature distribution," *Journal of Power Sources*, vol. 182, no. 2, pp. 630-638, 2008.
- [3] Y. Song and B. Wang, "Survey on reliability of power electronic systems," *IEEE Transactions on Power Electronics*, vol. 28, no. 1, pp. 591-604, 2013.
- [4] H. Zou, W. Wang, G. Zhang, F. Qin, C. Tian, and Y. Yan, "Experimental investigation on an integrated thermal management system with heat pipe heat exchanger for electric vehicle," *Energy Conversion and Management*, vol. 118, pp. 88-95, 2016.
- [5] H. Lambate, S. Nakanekar, and S. Tonapi, "Thermal characterization of the IGBT modules used in hybrid electric vehicles," in *Fourteenth Intersociety Conference on Thermal and Thermomechanical Phenomena in Electronic Systems (ITherm)*, 2014, pp. 1086-1091.
- [6] J. He, V. Mehrotra, and M. C. Shaw, "Thermal design and measurements of IGBT power modules: Transient and steady state," in *Conference Record of the 1999 IEEE Industry Applications Conference. Thirty-Fourth IAS Annual Meeting*, vol. 2, 1999, pp. 1440-1444.
- [7] Z. Huang, T. An, F. Qin, Y. Gong, Y. Dai, and P. Chen, "Effect of the thermal contact resistance on the heat dissipation performance of the press-pack IGBT module," in *2022 23rd International Conference on Electronic Packaging Technology (ICEPT)*, 2022, pp. 1-4.
- [8] A. S. Bahman, K. Ma, and F. Blaabjerg, "A lumped thermal model including thermal coupling and thermal boundary conditions for high-power IGBT modules," *IEEE Transactions on Power Electronics*, vol. 33, no. 3, pp. 2518-2530, 2018.
- [9] F. Dong, Y. Feng, Z. Wang, and J. Ni, "Effects on thermal performance enhancement of pin-fin structures for insulated gate bipolar transistor (IGBT) cooling in high voltage heater system," *International Journal of Thermal Sciences*, vol. 146, p. 106106, 2019.
- [10] R. V. Rao and V. K. Patel, "Thermodynamic optimization of cross flow plate-fin heat exchanger using a particle swarm optimization algorithm," *International Journal of Thermal Sciences*, vol. 49, no. 9, pp. 1712-1721, 2010.
- [11] S. M. Lee and K. Y. Kim, "Multi-objective optimization of arc-shaped ribs in the channels of a printed circuit heat exchanger," *International Journal of Thermal Sciences*, vol. 94, pp. 1-8, 2015.
- [12] M. Mishra, P. K. Das, and S. Sarangi, "Optimum design of crossflow plate-fin heat exchangers through genetic algorithm," *International Journal of Heat Exchangers*, vol. 5, no. 2, pp. 379-402, 2004.
- [13] M. Mishra and P. K. Das, "Thermoeconomic design-optimisation of crossflow plate-fin heat exchanger using genetic algorithm," *International Journal of Exergy*, vol. 6, no. 6, pp. 237-252, 2009.
- [14] W. T. Han, L. H. Tang, and G. N. Xie, "Performance comparison of particle swarm optimization and genetic algorithm in rolling fin-tube heat exchanger optimization design," in *Proceedings of the ASME Summer Heat Transfer Conference*, 2008, pp. 5-10.
- [15] F. Pourfattah and M. Sabzpooshani, "On the thermal management of a power electronics system: Optimization of the cooling system using genetic algorithm and response surface method," *Energy*, vol. 232, p. 120951, 2021.
- [16] İ. Polatoğlu, L. Aydın, B. Ç. Nevruz, and S. Özer, "A novel approach for the optimal design of a biosensor," *Analytical Letters*, vol. 53, no. 9, pp. 1428-1445, 2020.
- [17] R. L. Haupt and S. E. Haupt, *Practical Genetic Algorithms*, 2nd ed., John Wiley & Sons, 2004.
- [18] M. Bilal, M. Pant, H. Zaheer, L. Garcia-Hernandez, and A. Abraham, "Differential Evolution: A review of more than two decades of research," *Engineering Applications of Artificial Intelligence*, vol. 90, p. 103479, 2020.
- [19] S. S. Fan, Y. C. Liang, and E. Zahara, "A genetic algorithm and a particle swarm optimizer hybridized with Nelder-Mead simplex search," *Computers & Industrial Engineering*, vol. 50, no. 4, pp. 401-425, 2006.
- [20] M. Savran and L. Aydın, "Natural frequency and buckling optimization considering weight saving for hybrid graphite/epoxy-sitka spruce and graphite-flax/epoxy laminated composite plates using stochastic methods," *Mech Adv Mater Struct*, vol. 30, no. 13, pp. 2637-2650, 2023. <https://doi.org/10.1080/15376494.2021.1875390>.
- [21] M. Savran, L. Aydın, A. Ayaz, and T. Uslu, "A new strategy for manufacturing, modeling, and optimization of 3D printed polylactide based on multiple nonlinear neuro regression analysis and stochastic optimization methods," *Proc Inst Mech Eng Part E J Process Mech Eng*, 2024. <https://doi.org/10.1177/09544089241272909>.

## Appendix

**Table 6.** Full form of fitted models given in Table 4 for IGBT temperature minimization

<b>L</b>	$Y = 134.038 - 7.87285 x_1 - 0.0350987 x_2 - 3.84166 x_3 + 0.0052639 x_4 - 0.00609842 x_5$
<b>LR</b>	$Y = (-56968.1 + 9808.6 x_1 - 1.10747 x_2 + 13287.9 x_3 + 13.2699 x_4 - 61.2884 x_5) / (-701.641 + 127.019 x_1 - 0.0295636 x_2 + 160.471 x_3 + 0.169625 x_4 - 0.731965 x_5)$
<b>SON</b>	$Y = 143.16 - 26.1213 x_1 - 0.266133 x_1^2 + 0.205037 x_2 + 0.0164272 x_1 x_2 - 0.000215368 x_2^2 - 18.958 x_3 + 3.16074 x_1 x_3 + 0.0316068 x_2 x_3 - 6.10672 x_3^2 - 0.0526772 x_4 - 0.00426426 x_1 x_4 - 0.000106816 x_2 x_4 + 0.0345683 x_3 x_4 + 6.44028 * 10^{-6} x_4^2 - 0.272179 x_5 + 0.0666481 x_1 x_5 - 0.000173512 x_2 x_5 - 0.118958 x_3 x_5 + 0.000406931 x_4 x_5 + 0.000532985 x_5^2$
<b>SONR</b>	$Y = (0.999998 + 1.00097 x_1 + 1.01073 x_1^2 + 1.18699 x_2 + 1.59093 x_1 x_2 + 8.07745 x_2^2 + 0.998744 x_3 + 0.999468 x_1 x_3 + 1.09109 x_2 x_3 + 0.992578 x_3^2 + 1.41276 x_4 + 2.01127 x_1 x_4 + 0.653814 x_2 x_4 + 1.19341 x_3 x_4 - 2.6552 x_4^2 + 0.973255 x_5 + 0.983887 x_1 x_5 +$

	$\frac{26.9875 x^2 x^5 + 0.545355 x^3 x^5 + 32.2588 x^4 x^5 + 8.93178 x^5 x^2}{(1.0129 + 1.03098 x^1 + 0.564842 x^1^2 - 10.509 x^2 + 16.1285 x^1 x^2 - 0.0135889 x^2^2 + 1.13578 x^3 + 1.33672 x^1 x^3 + 0.959568 x^2 x^3 + 1.6953 x^3^2 - 28.6864 x^4 - 7.66247 x^1 x^4 + 0.141888 x^2 x^4 - 3.19273 x^3 x^4 - 0.034077 x^4^2 + 7.05449 x^5 + 15.0197 x^1 x^5 + 0.586826 x^2 x^5 + 50.0506 x^3 x^5 + 0.032533 x^4 x^5 - 0.0767997 x^5^2)}$
<b>FOTN</b>	$Y = -15.9216 + 7.46982 \cos[x_1] - 0.53304 \cos[x_2] - 87.6255 \cos[x_3] + 3.45563 \cos[x_4] - 5.72392 \cos[x_5] + 8.78887 \sin[x_1] - 2.18344 \sin[x_2] + 55.7862 \sin[x_3] + 0.443394 \sin[x_4] - 0.265043 \sin[x_5]$
<b>FOTNR</b>	$Y = (-0.730588 + 3.91013 \cos[x_1] + 1.92254 \cos[x_2] + 2.31712 \cos[x_3] + 0.309367 \cos[x_4] + 1.69828 \cos[x_5] + 4.18388 \sin[x_1] + 2.25224 \sin[x_2] + 0.0157777 \sin[x_3] + 4.34981 \sin[x_4] + 3.5477 \sin[x_5]) / (-0.0652036 + 0.0572789 \cos[x_1] + 0.0230747 \cos[x_2] - 0.0244729 \cos[x_3] - 0.000143375 \cos[x_4] + 0.0243629 \cos[x_5] + 0.0523375 \sin[x_1] + 0.0282922 \sin[x_2] + 0.0261709 \sin[x_3] + 0.0541034 \sin[x_4] + 0.0437403 \sin[x_5])$
<b>SOTN</b>	$Y = -47.3662 + 13.0646 \cos[x_1] - 63.5034 \cos[x_1]^2 + 2.56338 \cos[x_2] - 47.6764 \cos[x_2]^2 - 436.355 \cos[x_3] - 85.6018 \cos[x_3]^2 + 3.40693 \cos[x_4] - 57.3509 \cos[x_4]^2 - 3.14863 \cos[x_5] - 63.4364 \cos[x_5]^2 + 1.11839 \sin[x_1] - 60.9763 \sin[x_1]^2 - 0.102909 \sin[x_2] - 45.5672 \sin[x_2]^2 + 41.4972 \sin[x_3] + 139.775 \sin[x_3]^2 + 1.48926 \sin[x_4] - 45.5672 \sin[x_4]^2 + 1.84843 \sin[x_5] - 72.5961 \sin[x_5]^2$
<b>SOTNR</b>	$Y = (5.06164 + 0.909237 \cos[x_1] + 1.62149 \cos[x_1]^2 + 4.89529 \cos[x_2] + 0.713049 \cos[x_2]^2 - 3.13559 \cos[x_3] + 5.15071 \cos[x_3]^2 + 7.27358 \cos[x_4] + 0.649242 \cos[x_4]^2 + 8.37598 \cos[x_5] + 0.962236 \cos[x_5]^2 + 9.14495 \sin[x_1] + 4.44015 \sin[x_1]^2 - 1.426 \sin[x_2] + 9.76099 \sin[x_2]^2 + 1.76466 \sin[x_3] + 0.910924 \sin[x_3]^2 - 3.48212 \sin[x_4] + 9.76099 \sin[x_4]^2 + 4.65247 \sin[x_5] + 5.0994 \sin[x_5]^2) / (-0.470253 + 0.0554459 \cos[x_1] + 0.303031 \cos[x_1]^2 + 0.0464251 \cos[x_2] - 0.0150328 \cos[x_2]^2 - 0.831807 \cos[x_3] - 0.398776 \cos[x_3]^2 + 0.0428974 \cos[x_4] + 0.0192471 \cos[x_4]^2 + 0.145946 \cos[x_5] + 0.197603 \cos[x_5]^2 + 0.127331 \sin[x_1] + 0.226716 \sin[x_1]^2 - 0.00448057 \sin[x_2] + 0.0552804 \sin[x_2]^2 - 0.834104 \sin[x_3] + 0.928524 \sin[x_3]^2 - 0.0553128 \sin[x_4] + 0.0552804 \sin[x_4]^2 - 0.0048026 \sin[x_5] + 0.332144 \sin[x_5]^2)$
<b>FOLN</b>	$Y = 205.583 - 21.4628 \log[x_1] - 15.1221 \log[x_2] - 8.89083 \log[x_3] + 0.497759 \log[x_4] + 0.21157 \log[x_5]$
<b>FOLNR</b>	$Y = (-3235.77 + 668.333 \log[x_1] - 13.0274 \log[x_2] + 734.962 \log[x_3] + 472.026 \log[x_4] - 211.339 \log[x_5]) / (-38.355 + 8.63867 \log[x_1] - 0.687913 \log[x_2] + 8.60961 \log[x_3] + 6.16418 \log[x_4] - 2.65231 \log[x_5])$
<b>SOLN</b>	$Y = 515.076 - 89.4979 \log[x_1] - 12.2692 \log[x_1]^2 + 155.474 \log[x_2] + 40.8295 \log[x_1] \log[x_2] + 10.3467 \log[x_2]^2 - 1127.92 \log[x_3] + 35.5538 \log[x_1] \log[x_3] + 119.958 \log[x_2] \log[x_3] - 58.5787 \log[x_3]^2 - 53.6643 \log[x_4] - 37.03 \log[x_1] \log[x_4] - 40.0015 \log[x_2] \log[x_4] + 98.8895 \log[x_3] \log[x_4] + 1.13931 \log[x_4]^2 + 5.57893 \log[x_5] + 11.8181 \log[x_1] \log[x_5] - 50.2312 \log[x_2] \log[x_5] - 55.7305 \log[x_3] \log[x_5] + 53.933 \log[x_4] \log[x_5] - 1.0972 \log[x_5]^2$
<b>SOLNR</b>	$Y = (-0.335514 - 0.927137 \log[x_1] - 1.39658 \log[x_1]^2 - 8.13074 \log[x_2] - 11.7381 \log[x_1] \log[x_2] - 61.8455 \log[x_2]^2 - 0.248695 \log[x_3] - 0.482518 \log[x_1] \log[x_3] - 7.36435 \log[x_2] \log[x_3] - 0.316384 \log[x_3]^2 + 12.6517 \log[x_4] + 11.259 \log[x_1] \log[x_4] + 75.8271 \log[x_2] \log[x_4] + 12.8017 \log[x_3] \log[x_4] + 97.142 \log[x_4]^2 - 7.39312 \log[x_5] - 9.95039 \log[x_1] \log[x_5] - 56.4527 \log[x_2] \log[x_5] - 7.26286 \log[x_3] \log[x_5] + 42.0738 \log[x_4] \log[x_5] - 40.8657 \log[x_5]^2) / (-0.000139906 + 2.58557 \log[x_1] - 1.09347 \log[x_1]^2 - 6.58669 \log[x_2] - 17.3051 \log[x_1] \log[x_2] - 12.3961 \log[x_2]^2 + 0.64661 \log[x_3] - 25.9045 \log[x_1] \log[x_3] - 18.4784 \log[x_2] \log[x_3] + 14.4391 \log[x_3]^2 - 8.55544 \log[x_4] + 23.5569 \log[x_1] \log[x_4] + 17.504 \log[x_2] \log[x_4] - 6.39485 \log[x_3] \log[x_4] + 0.668614 \log[x_4]^2 - 2.53994 \log[x_5] - 3.78552 \log[x_1] \log[x_5] + 20.5227 \log[x_2] \log[x_5] + 36.3018 \log[x_3] \log[x_5] - 23.9271 \log[x_4] \log[x_5] - 0.375019 \log[x_5]^2)$
<b>H (FOLN+SON)</b>	$Y = 90.1139 + 0.442761 x^1 - 2.29539 x^1^2 + 0.0746876 x^2 + 0.0156922 x^1 x^2 - 0.000131892 x^2^2 - 3.71033 x^3 + 6.16279 x^1 x^3 + 0.0250583 x^2 x^3 - 4.80624 x^3^2 + 0.0418714 x^4 - 0.0109804 x^1 x^4 - 0.0000555408 x^2 x^4 + 0.0225954 x^3 x^4 - 0.0000395405 x^4^2 - 0.11773 x^5 + 0.0438845 x^1 x^5 - 0.000164884 x^2 x^5 - 0.139816 x^3 x^5 + 0.000415728 x^4 x^5 + 0.000370671 x^5^2 - 39.8574 \log[x_1] + 8.66859 \log[x_2] - 26.899 \log[x_3] - 5.74012 \log[x_4] - 0.201207 \log[x_5]$
<b>H (FOLN*L)</b>	$Y = (37.0467 + 2.02402 x^1 - 0.044288 x^2 - 6.24084 x^3 - 0.0585529 x^4 - 0.000602118 x^5) (-19.2654 + 0.230805 \log[x_1] + 1.18382 \log[x_2] + 0.215925 \log[x_3] + 1.41784 \log[x_4] + 0.000906653 \log[x_5])$

# Time Series Prediction of Heart Rate Using Deep Learning Models

Emir EVCİL<sup>1\*</sup>

## Abstract

Cardiovascular diseases are among the leading causes of mortality worldwide and represent a significant global health burden, affecting millions of individuals each year. Early diagnosis of these diseases is critical not only for improving patient survival rates but also for ensuring the economic sustainability of healthcare systems. Heart rate values serve as essential biological indicators, providing important insights into cardiovascular health and offering potential utility in early diagnosis. In this study, conducted a comprehensive time series analysis to predict the next 5-minute heart rate values based on a 3-minute segment of pulse data collected from healthy individuals. Employed four deep learning models—Recurrent Neural Network (RNN), Bidirectional Long Short-Term Memory (BI-LSTM), Gated Recurrent Unit (GRU), and Long Short-Term Memory (LSTM)—to analyze the constructed dataset. The predictive performances of these models were rigorously compared using the Root Mean Square Error (RMSE) metric, which serves as a reliable measure of accuracy in regression tasks. Findings indicate that deep learning techniques, particularly LSTM and its variants, hold significant promise for enhancing the accuracy of heart rate predictions. This study underscores the potential of these advanced methodologies in the early diagnosis of cardiovascular diseases, aiming to offer new perspectives for the development of clinical decision support systems that could ultimately improve patient outcomes and optimize healthcare delivery.

**Keywords:** *BiLSTM; GRU; Heart Rate Prediction; LSTM; RNN; Time Series.*

## 1. Introduction

Cardiovascular diseases (CVDs) represent one of the most pressing health challenges globally, contributing significantly to morbidity and mortality rates. The World Health Organization (WHO) estimates that CVDs account for approximately 32% of all global deaths, making early diagnosis and intervention crucial for improving patient outcomes [1]. These conditions often manifest with subtle symptoms, and traditional diagnostic methods may not always provide timely insights. Thus, there is an urgent need for innovative approaches to facilitate early detection.

Heart rate, as a vital physiological parameter, serves as an essential indicator of cardiovascular health. Accurate and continuous monitoring of heart rate values can provide valuable information that aids in the early diagnosis of potential cardiovascular issues [2]. Recent advancements in machine learning, particularly in deep learning techniques, offer promising avenues for analyzing complex time series data, such as heart rate measurements [3].

The heart rate signal is a typical time series [4]. A time series is a series of data points indexed in chronological order. Effective forecasting of time series enables better use of available information for analysis and decision making. Its wide range of applications includes but is not limited to clinical medicine, financial forecasting, traffic flow forecasting, human behaviour forecasting and other fields [5]. Unlike other predictive modelling tasks, time series increases the complexity of sequence dependencies between input variables. Therefore, how to build a suitable predictive model for real-time forecasting tasks by fully exploiting the complex sequence dependencies is an important issue [6].

Importantly, time series data, including heart rate values, often exhibit non-stationarity, meaning their statistical properties, such as the mean and variance, can change over time. This non-stationary behavior can pose significant challenges for modeling and forecasting, as many traditional time series analysis techniques assume stationarity and may not perform well when this assumption is violated [7]. For instance, physiological signals can be influenced by a variety of factors, including emotional states, physical activity levels, and underlying health conditions, all of which can lead to fluctuations in heart rate that reflect complex dynamics rather than stable patterns. Consequently, failing to account for non-stationarity in these signals can result in inaccurate predictions and misinterpretations of the underlying physiological processes. Therefore, developing methodologies that can effectively handle non-stationary time series data is crucial for enhancing the reliability of cardiovascular health monitoring and interventions.

\*Corresponding author

Emir EVCİL<sup>\*</sup>; Izmir Katip Celebi University, Faculty of Engineering and Architecture, Electrical Electronic Engineering Department, Türkiye; e-mail: [190403008@ogr.ikcu.edu.tr](mailto:190403008@ogr.ikcu.edu.tr);  0009-0004-4089-6638

Using machine learning methods, non-linear prediction models can be built based on large amounts of historical data. In fact, through repeated training iterations and learning approaches, machine learning models can obtain more accurate predictions than traditional statistical-based models. Typical methods include tree-based ensemble learning methods such as support vector regression or kernel-based classification and artificial neural multiagent (ANN) with strong nonlinear function approximation, and gradient-enhanced regression or decision tree. However, since the above method lacks efficient handling of sequence dependencies between input variables, its effectiveness in time series forecasting tasks is limited [8].

Recurrent neural networks (RNN) are often regarded as the most efficient method of time series forecasting. In fact, RNN is an artificial neural network in which the nodes are connected in a loop and the internal state of the network can exhibit dynamic timing behaviour. However, as the length of the process time series increases, problems such as gradient disappearance often arise during the training of RNNs using conventional activation functions such as tanh or sigmoid functions, which limit the prediction accuracy of RNNs. The Long and Short Term Memory Unit (LSTM) is based on a simple RNN that solves the memory and forgetting problems by adding some multi-threshold gates. Therefore, LSTM and Gated Loop Unit (GRU) address to some extent the limited ability to deal with long-term dependencies. These methods have been successfully applied to many sequential learning problems such as machine translation [9]. In general, LSTM is considered one of the most advanced methods for dealing with time series forecasting problems. Inspired by cognitive neuroscience, some researchers have incorporated attentional mechanisms into the encoding-decoding framework [10]. Attentional mechanisms can better select input sequences and encode semantics in long-term memory to improve the information processing capabilities of neural multimodality. Recently, attention mechanisms have been widely used and perform well in many different types of deep learning tasks, such as image captioning, visual question answering and speech recognition. Specifically, most research work [11] is usually done by adding an attention layer to the encoding-decoding framework.

Thanks to the effective performance of the LSTM model in time series analysis, it is widely applied to establish long-term relationships in heart rate data and abstract high-dimensional features. For example, Haijun et al. proposed an LSTM-BiLSTM-Att model that combines LSTM, BiLSTM, and an attention mechanism with a fully connected neural network, achieving a notable RMSE of 1.729 in heart rate prediction tasks [12]. Similarly, Haowei et al. emphasized the superiority of transformer-based models for predicting cardiovascular health data, as compared to traditional methods like ARIMA and Prophet, as well as other deep learning approaches [13]. Additionally, Staffini et al. (2022) demonstrated the challenges in predicting heart rate time series due to their nonlinear and non-stationary nature. Their study compared three different prediction models: Autoregressive Model, Long Short-Term Memory Network, and Convolutional Long Short-Term Memory Network, highlighting that the Autoregressive Model consistently outperformed the others across different environments, achieving an average absolute error of 2.069, which was better than the results from the LSTM and ConvLSTM models [14].

By demonstrating the effectiveness of deep learning models in predicting heart rate values, this research contributes to the existing body of knowledge in e-health and underscores the transformative potential of machine learning technologies in enhancing cardiovascular care. Specifically, in this study, we collect data from a pulse oximeter at 1-second intervals, enabling a comprehensive time series analysis. We evaluate the performance of various models, including LSTM, BiLSTM, RNN, and GRU, to assess their effectiveness in capturing the underlying patterns of heart rate variability. The subsequent sections of this paper will detail methodology, present the results of analysis, and discuss the implications of findings for future research and clinical practice.

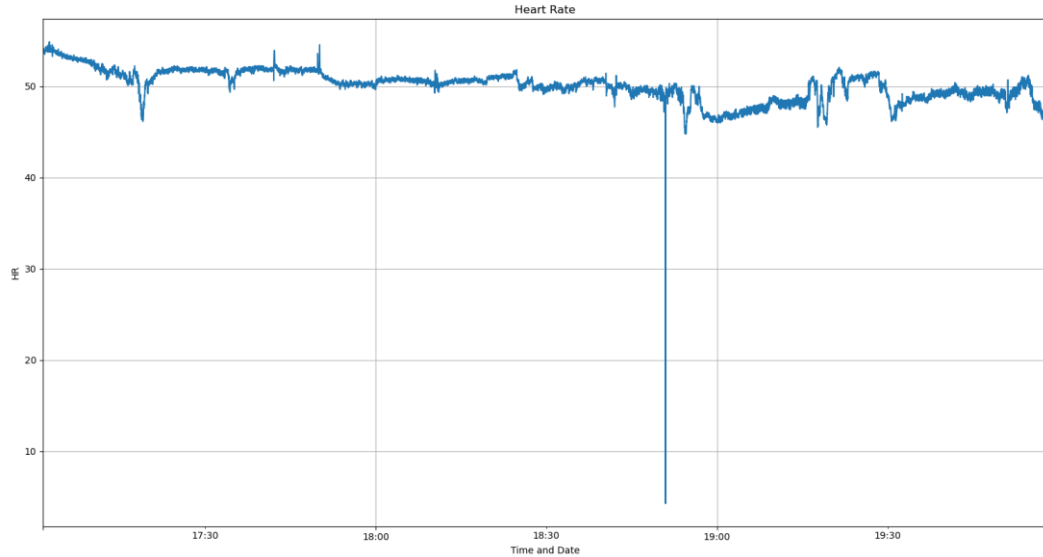
## 2. Method

In this study, proposed a deep learning-based system aimed at predicting heart rate five minutes in advance, leveraging the strengths of neural network architectures to handle the temporal dependencies inherent in heart rate data. This predictive model is designed to assist healthcare professionals by providing real-time insights into the patient's cardiovascular trends, which could be instrumental for early intervention in clinical settings. The model development is structured through five essential stages to ensure accurate and reliable predictions:

- Data collection
- Data pre-processing
- Data splitting
- Training and optimization of the model

## 2.1. Data Collection

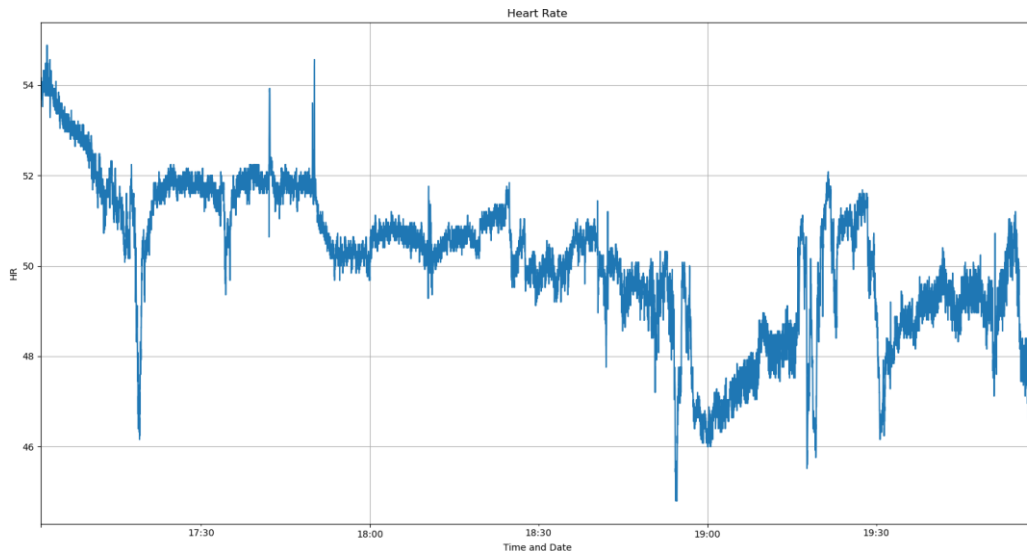
Heart rate data was collected continuously over a 3-hour period using a MAX30100 pulse oximeter sensor connected to a Raspberry Pi 3 microcontroller. Measurements were taken at a rate of one sample per second, resulting in a dataset of 10,597 entries.



**Figure 1.** Time - heart rate plot of the data prepared as a result of 3 hours of measurement.

## 2.2 Data Pre-Processing

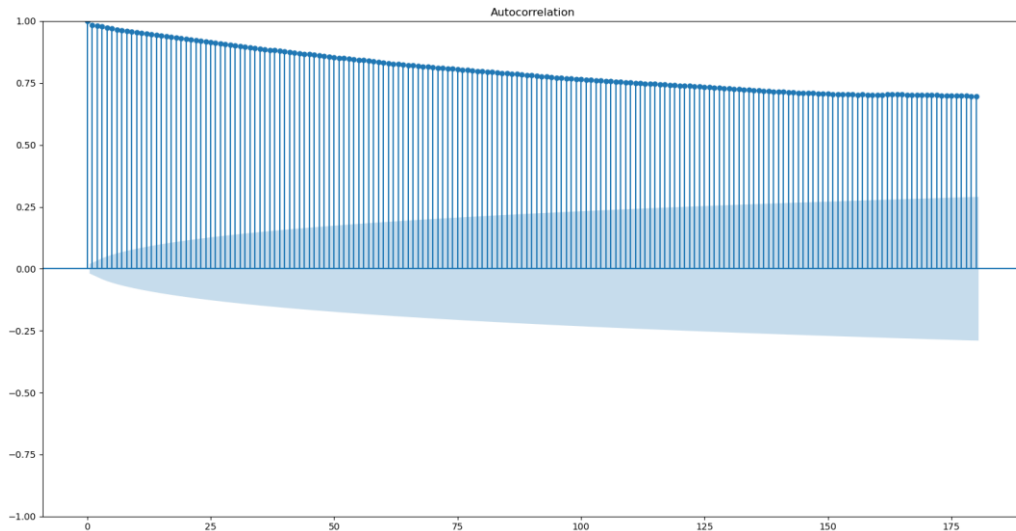
During the data collection process, several outliers were identified, likely resulting from sensor measurement errors. To enhance the quality and reliability of the dataset, these outliers were addressed by replacing them with the average of several preceding and succeeding data points. This technique, known as local averaging, effectively smooths the data and mitigates the impact of random measurement errors on the overall analysis. By using the average of surrounding values, preserved the continuity of the data, avoiding the introduction of arbitrary values that could skew the results. Additionally, local averaging helps reduce noise, dampening the influence of extreme values that could distort trends within the dataset. Ultimately, cleaning the dataset by replacing outliers ensures that the model is trained on more representative data, leading to improved predictions and generalization when applied to unseen data. This approach enhances the robustness of the predictive model by ensuring that the underlying trends in the heart rate data are accurately represented, thereby contributing to more reliable outcomes in the analysis.



**Figure 2.** Data after outliers are removed.

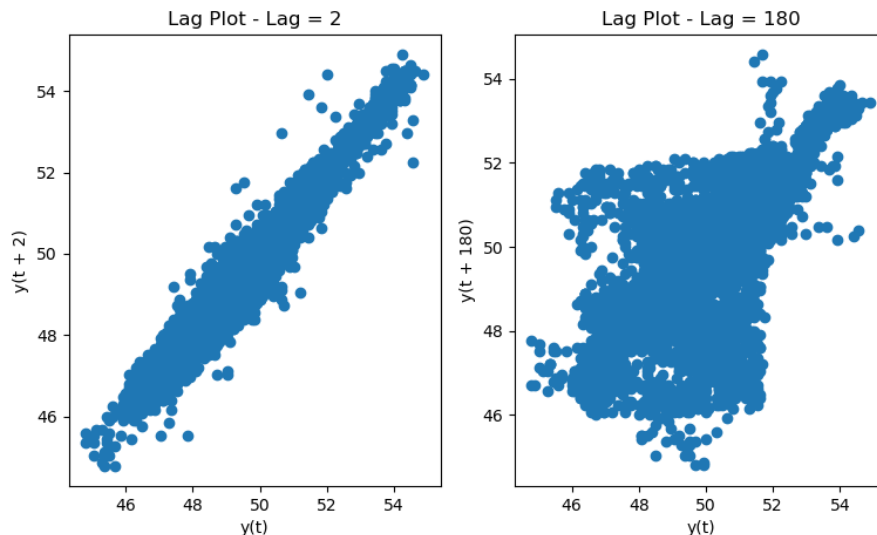
In time series analysis, an autocorrelation graph is utilized to examine the degree to which current values in a dataset are related to their past values over different time lags. Autocorrelation measures the correlation of a time series with its own past values, helping to identify patterns, trends, and the presence of seasonality in the data. By plotting the autocorrelation function (ACF) against various lag values, researchers can discern whether past observations significantly influence future observations, which is essential for understanding the underlying structure of the data.

In this study, an autocorrelation graph was constructed to investigate the relationship between pulse values over time. Specifically, with a lag value of 180 seconds, the autocorrelation coefficient was calculated to be 75%. This substantial correlation suggests a strong relationship between pulse values at this lag, indicating that changes in pulse readings are significantly influenced by preceding values. Such a correlation value is compelling for time series analysis, as it implies the potential for predictive modeling based on historical data.



**Figure 3.** Autocorrelation plot for lag = 180.

A lag plot is a graphical tool used in time series analysis to visualize the relationship between observations at different time lags. By plotting the values of a time series against their lag values, it is possible to assess whether a linear or nonlinear relationship exists. In this study, lag plots were created for both lag values of 2 and 180 seconds. A clear linear trend was observed in both plots, indicating that the current pulse values are closely related to their past values at these particular lags. This linearity indicates that the data exhibits a persistent structure over time, reinforcing previous findings of nonstationarity. The presence of a linear trend in the lag plots serves as a secondary indicator of the nonstationarity of the data set, alongside the gradual decrease in the autocorrelation plot, since stationary data usually exhibit a more random distribution without a clear directional trend.



**Figure 4.** Lag plot of heart rate values for different lag values.

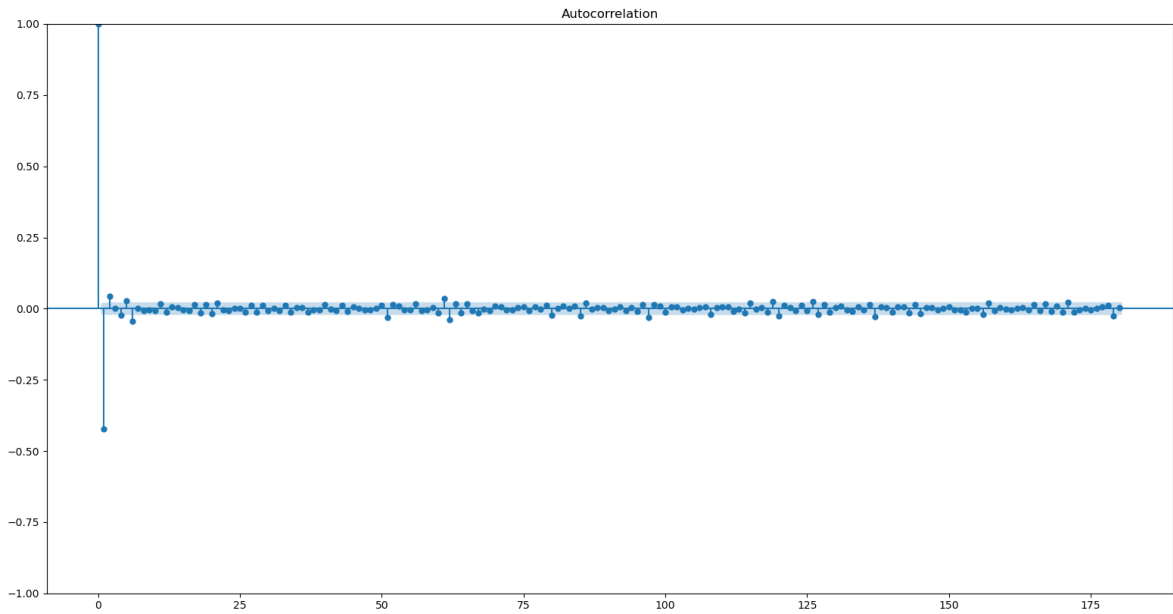
The statistical properties of stationary data remain constant over time, which means that key characteristics such as mean, variance, and autocorrelation do not exhibit changes as time progresses. Stationarity is a fundamental assumption in time series analysis, as many statistical methods and models rely on this condition for their validity. When working with non-stationary data, statistical inference can yield misleading results, leading to incorrect conclusions about the underlying processes.

To transform non-stationary data into stationary data, various techniques can be employed [15], with differencing being one of the most commonly used methods. This process involves subtracting the previous observation from the current observation, effectively removing trends and seasonality from the data.

$$Y_t = X_t - X_{t-1} \tag{1}$$

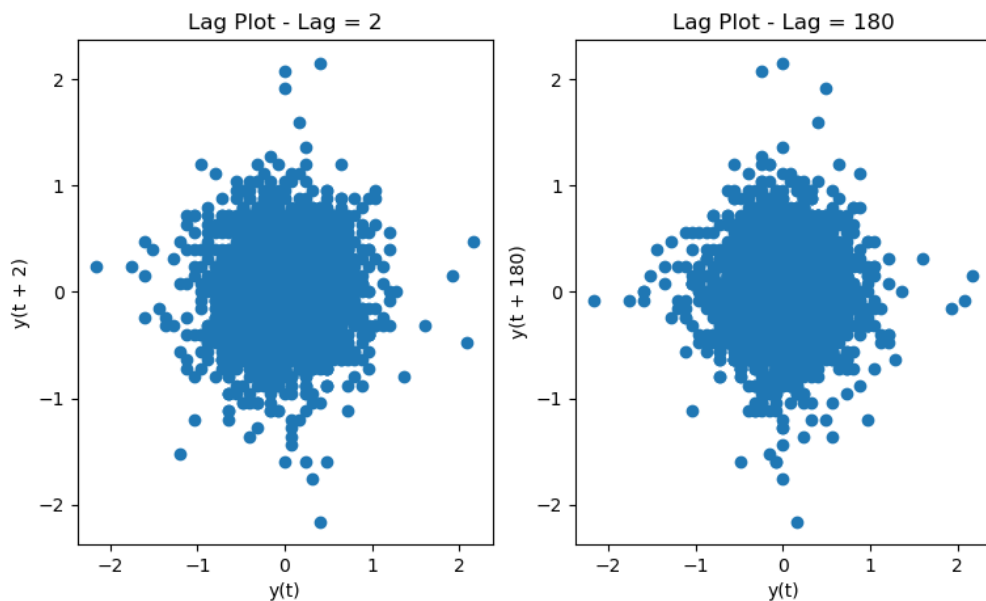
The differencing method was applied during the data preprocessing phase to achieve stationarity.

Figures 5 and 6 show lag autocorrelation plots after the data are made stationary by differencing.



**Figure 5.** Autocorrelation plot of stationary data.

A rapid decrease observed in the autocorrelation graph reveals that the data set has a constant mean and variance, and this situation allows for more reliable results to be obtained in the modeling processes.



**Figure 6.** Lag plot of stationary heart rate values for different lag values.



The analyses performed on the autocorrelation graph and lag plot obtained after the differencing process show that the data has been made stationary. The random distribution of the lag plot shows that the dependency in the time series data has disappeared and the observations have become independent of each other.

After differencing, heart rate data were scaled to the range  $[-1, 1]$  using the “minmax scaler”. Observation windows were defined so that time series data could be processed with the deep learning models. Observation window is a vector representing the data observed during a certain time period.

$$X_t = (x_{t-1}, x_{t-2}, \dots, x_{t-n}) \quad (2)$$

In this context, the model utilizes these historical data points to predict the heart rate measurement at  $t+300$  seconds, which represents a 5-minute prediction horizon. By leveraging the temporal dependencies captured in the observation windows, the deep learning models can effectively forecast future heart rate values. This capability is crucial for providing timely insights in real-time healthcare applications, where swift decision-making can significantly impact patient outcomes.

**Table 1.** 180-seconds observation and 300-seconds target Windows.

Time (hours)	t-180	t-179	t-178	t-177	...	t+296	t+297	t+298	t+299
17:04:29	0.08	0.00	0.08	-0.08	...	-0.24	0.24	0.08	-0.24
17:04:30	0.00	0.08	-0.08	-0.08	...	0.24	0.08	-0.24	0.24
17:04:31	0.08	-0.08	-0.08	0.32	...	0.08	-0.24	0.24	-0.08
...	...	...	...	...	...	...	...	...	...
19:53:03	-0.32	-0.56	0.80	-0.16	...	0.16	-0.72	0.16	0.24
19:53:04	-0.56	0.80	-0.16	-0.80	...	-0.72	0.16	0.64	-0.24
19:53:05	0.80	-0.16	-0.44	0.06	...	-0.32	0.16	-0.72	0.64

As shown in Table 1, after being subjected to differencing, the data scaled to the range  $[-1, 1]$  were shifted and prepared as 180-second observation (X) and 300-second target (Y) windows. As a result of these operations, the data preprocessing process was completed.

### 2.3. Data Splitting

As this study involves time series data, the dataset is split in chronological order, without shuffling to preserve the temporal structure. Specifically:

- 75% of the data is used for training,
- 10% for validation,
- 15% for testing.

This method ensures that the model is evaluated on unseen future data, which mimics real-world forecasting scenarios. The training **set** is used to update the model's parameters, while the validation set is reserved for assessing overfitting tendencies and for techniques like early stopping. Finally, the **test set** is used exclusively to evaluate the model's performance on data that was not involved in any part of the training process.

### 2.4. Training and Optimization Of the Models

To comprehensively evaluate the performance of various deep learning architectures in processing the prepared pulse rate data, developed models using Gated Recurrent Unit (GRU), Bidirectional Long Short-Term Memory (Bi-LSTM), Long Short-Term Memory (LSTM), and Recurrent Neural Network (RNN) architectures. Each model was configured with varying depths, specifically utilizing 1 to 3 layers, to determine the optimal structure for data analysis.

To fine-tune the hyperparameters of these models, employed the grid search method [16], focusing on two critical parameters: the number of neurons within each layer and the dropout rate. The grid search configuration was set up as follows:

- The number of neurons in each layer was systematically varied from 10 to 200, increasing in increments of 10.
- The dropout rate was adjusted between 0.1 and 0.5, with increments of 0.1.

The use of Grid Search allowed us to comprehensively explore the hyperparameter space, facilitating the identification of the most effective configurations for each architecture. For the training of the models, employed

the Adam optimizer, known for its efficiency and ability to adaptively adjust the learning rate. The models were trained for 50 epochs with a patience level of 5 for early stopping, preventing overfitting by monitoring the validation loss. Setted the learning rate to 0.0001 and used a batch size of 1, which enabled the model to update weights more frequently and respond dynamically to the training data.

### 3. Results

In the evaluation of the performance of each model, utilized the root mean square error (RMSE) metric, defined mathematically as follows:

$$\text{RMSE} = \sqrt{\frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2} \quad (3)$$

RMSE was chosen as the primary evaluation metric for several reasons. Firstly, it provides a clear measure of the average magnitude of the prediction errors, allowing for a direct interpretation of how well the model's predictions align with the actual values. By squaring the errors, RMSE emphasizes larger discrepancies, making it particularly sensitive to outliers, which is important in healthcare applications where accurate pulse rate measurements are critical.

Secondly, RMSE has the same units as the target variable (in this case, heart rate), which facilitates intuitive understanding of the model's performance. This characteristic makes it easier for practitioners and stakeholders to assess the clinical relevance of the model's predictions.

The performance of each model was assessed based on the root mean square error (RMSE) for a lag of 3 minutes and a forecast horizon of 5 minutes. The results, summarized in Table 2, show that there are some variations in the predictive capabilities of the different architectures and configurations employed in this study.

**Table 2.** *The result of models for lag of 3 minutes and forecasting 5 minutes.*

Model	Number of Layers	Number of Neurons	Dropout Rate	RMSE
RNN	1	100	0.3	2.121
	2	[40,80]	[0.4,0.2]	2.041
	3	[40,120,80]	[0.5,0.3,0.5]	1.971
LSTM	1	70	0.5	2.071
	2	[120,100]	[0.4,0.3]	1.881
	3	[160,120,40]	[0.3,0.5,0.3]	1.839
GRU	1	150	0.2	2.008
	2	[140,50]	[0.5,0.5]	2.045
	3	[70,120,70]	[0.3,0.5,0.3]	1.940
Bi-LSTM	1	70	0.2	2.035
	2	[70,120]	[0.5,0.5]	2.173
	3	[120,140,110]	[0.3,0.5,0.3]	2.223

The results show that the RNN model with three layers achieved the best performance among all configurations, with an RMSE of 1.971, while the one-layer RNN had the worst performance with an RMSE of 2.121. For the LSTM model, the three-layer configuration also performed best, resulting in an RMSE of 1.839, whereas the two-layer LSTM produced the highest RMSE of 2.071. In the case of the GRU model, the two-layer setup yielded the best performance with an RMSE of 1.940, while the three-layer configuration had the least favorable outcome, resulting in an RMSE of 2.045. The Bi-LSTM model demonstrated the highest performance with a one-layer setup, achieving an RMSE of 2.035, while the three-layer configuration recorded the highest RMSE of 2.223.

### 4. Conclusions

In this study, comprehensively evaluated the performance of various deep learning architectures (i.e. Long Short-Term Memory (LSTM), Bidirectional LSTM (BiLSTM), Gated Recurrent Unit (GRU), and Recurrent Neural Networks (RNN)) in predicting heart rate data by collecting 10,597 pulse values for three hours. Findings revealed that the LSTM model configured with three layers demonstrated superior predictive capabilities, achieving the lowest Root Mean Square Error (RMSE) of 1.839. This result not only underscores the effectiveness of LSTM networks in handling time series data but also highlights their potential to significantly enhance the accuracy of heart rate predictions, which is critical for the development of advanced e-health monitoring systems.

The implications of these findings are far-reaching, as accurate heart rate predictions can lead to timely interventions and improved patient outcomes, particularly in remote monitoring scenarios where immediate clinical assessment may not be possible. The integration of reliable predictive models in e-health systems could enable healthcare providers to make informed decisions swiftly, thereby optimizing patient care and resource allocation.

However, it is important to acknowledge certain limitations within this study. The dataset utilized was confined to a specific population and time duration, which may limit the generalizability of the results across diverse demographic groups and clinical conditions. Future research could address this limitation by incorporating a broader range of datasets that include various age groups, health statuses, and environmental conditions. Additionally, the integration of complementary physiological signals—such as blood pressure, oxygen saturation, and even physical activity levels—could further enhance the predictive power of these models, providing a more holistic view of patient health.

## Declaration of Interest

The authors declare that there is no conflict of interest.

## References

- [1] World Health Organization, “Cardiovascular Diseases (CVDs),” Available: [https://www.who.int/news-room/fact-sheets/detail/cardiovascular-diseases-\(cvds\)](https://www.who.int/news-room/fact-sheets/detail/cardiovascular-diseases-(cvds)).
- [2] H. Rahman, M. U. Ahmed, and S. Begum, “Vision-based remote heart rate variability monitoring using camera,” in *Internet of Things (IoT) Technologies for HealthCare: 4th Int. Conf., HealthyIoT 2017, Angers, France, Oct. 24-25, 2017, Proc.\**, 2018, pp. 10–18.
- [3] A. Esteva, A. Robicquet, B. Ramsundar, V. Kuleshov, M. DePristo, K. Chou, and J. Dean, “A guide to deep learning in healthcare,” *Nat. Med.\**, vol. 25, no. 1, pp. 24–29, 2019.
- [4] M. Oyeleye, T. Chen, S. Titarenko, and G. Antoniou, “A predictive analysis of heart rates using machine learning techniques,” *Int. J. Environ. Res. Public Health\**, vol. 19, no. 4, pp. 2417, 2022.
- [5] S. Makridakis, E. Spiliotis, and V. Assimakopoulos, “M5 accuracy competition: Results, findings, and conclusions,” *Int. J. Forecast.\**, vol. 38, no. 4, pp. 1346–1364, 2022.
- [6] R. J. Hyndman, *Forecasting: Principles and Practice\**. OTexts, 2018.
- [7] R. B. Govindan, A. N. Massaro, N. Niforatos, and A. Du Plessis, “Mitigating the effect of non-stationarity in spectral analysis—An application to neonate heart rate analysis,” *Comput. Biol. Med.\**, vol. 43, no. 12, pp. 2001–2006, 2013.
- [8] Y. Qin, D. Song, H. Chen, W. Cheng, G. Jiang, and G. W. Cottrell, “A dual-stage attention-based recurrent neural network for time series prediction,” in *Proc. 26th Int. Joint Conf. Artif. Intell.\**, 2017, pp. 2627–2633.
- [9] K. Cho, B. van Merriënboer, Ç. Gülçehre, D. Bahdanau, F. Bougares, H. Schwenk, and Y. Bengio, “Learning phrase representations using RNN encoder–decoder for statistical machine translation,” in *Proc. 2014 Conf. Empir. Methods Nat. Lang. Process.\**, 2014, pp. 1724–1734.
- [10] D. Bahdanau, K. Cho, and Y. Bengio, “Neural machine translation by jointly learning to align and translate,” *arXiv preprint arXiv:1409.0473*, 2014.
- [11] Y. Liang, S. Ke, J. Zhang, X. Yi, and Y. Zheng, “Geoman: Multi-level attention multi-order for geo-sensory time series prediction,” in *Proc. 27th Int. Joint Conf. Artif. Intell.\**, 2018, pp. 3428–3434.
- [12] H. Lin, S. Zhang, Q. Li, Y. Li, J. Li, and Y. Yang, “A new method for heart rate prediction based on LSTM-BiLSTM-Att,” *Measurement\**, vol. 207, p. 112384, 2023.
- [13] H. Ni et al., “Time Series Modeling for Heart Rate Prediction: From ARIMA to Transformers,” *arXiv preprint arXiv:2406.12199*, 2024.
- [14] A. Staffini, T. Svensson, U. I. Chung, and A. K. Svensson, “Heart rate modeling and prediction using autoregressive models and deep learning,” *Sensors\**, vol. 22, no. 1, p. 34, 2021.
- [15] R. Salles, K. Belloze, F. Porto, P. H. Gonzalez, and E. Ogasawara, “Nonstationary time series transformation methods: An experimental review,” *Knowl.-Based Syst.\**, vol. 164, pp. 274–291, 2019.
- [16] J. Bergstra and Y. Bengio, “Random Search for Hyper-Parameter Optimization,” *J. Mach. Learn. Res.\**, vol. 13, pp. 281–305, 2012.

# Artificial Intelligence Applications in Drug Discovery and Research

Seyma MINTAS<sup>1</sup>, Canan SEVİMLİ-GUR<sup>2,\*</sup>

## Abstract

Existing drug treatments may be inadequate for all the prevalent and emerging diseases that people face every day. Therefore, the discovery and development of new drugs is an inevitable necessity to protect human health and treat diseases. Discovering a new drug involves many steps, such as selecting the drug with therapeutic effect from a large number of active compounds, determining the ADMET properties of the drug, and conducting clinical studies to determine its toxicity and side effect profile. It is a costly and time-consuming process. According to the California Biomedical Research Association, it takes an average of 12 years and \$359 million to get a drug from the lab to the patient. With the increase in digitalization in the field of health, as in every field, scientists have resorted to artificial intelligence (AI) to solve the problem of cost and time. Pharmaceutical manufacturing companies have made major investments and developed numerous AI-based algorithms to be used in different stages of drug discovery. With the use of these algorithms in drug discovery and research, the money and time spent has decreased and efficiency has increased. This mini-review discusses AI applications in drug discovery and research.

**Keywords:** *Medicine, drug research, artificial intelligence, ADMET, algorithm, Big Data.*

## 1. Introduction

Although there is no universally accepted definition of the concept of artificial intelligence (AI) to date, it is broadly described as the simulation of human intelligence by computers. It refers to software and systems that enable computers to perform actions that require human intelligence. It is a multidisciplinary field [1].

The data-intensive environment referred to as 'Big Data' necessitates the acquisition, integration, and analysis of vast datasets to address complex medical and scientific challenges. It is essential to utilize AI techniques to uncover meaningful hidden patterns. Machine learning (ML) and deep learning (DL) are the two most commonly used subfields/domains of AI in drug discovery. ML and DL algorithms are applied to Big Data to extract meaningful and actionable insights from complex and heterogeneous datasets. [2, 3].

Scanning Big Data in the pharmaceutical field with AI techniques has significantly shortened the drug discovery and development process. With the shortening of the process, drug discovery has become more economical.

AI algorithms leverage existing data to enhance analysis and evaluation, use from the identification of a drug candidate to the manufacturing process in the industry. Therefore, before the synthesis and experimental evaluation of the drug molecule, AI-driven analysis has an important place in evaluating the effectiveness and efficiency of drug candidates against the desired disease.

Commonly employed DL methods in drug discovery and development can be listed as convolutional neural network, recurrent neural network, restricted boltzmann machine, autoencoders, artificial neural networks, fuzzy expert system approach [1, 4].

## 2. Application of Artificial Intelligence Techniques to Drug Development Steps

After successfully undergoing preclinical testing, a drug candidate proceeds to clinical trials, which are conducted in three phases: Phase 1 focuses on drug safety evaluation in a small group of healthy volunteers; Phase 2 examines drug efficacy and safety in a small group of patients with the target disease; and Phase 3 involves large-scale efficacy and safety studies in diverse patient populations. Successful completion of clinical trials is mandatory for FDA evaluation and subsequent market approval. Failures at any of these three stages render the drug development process inefficient, reduces the value of the investment, and increases the cost. The two main reasons behind the high failure rates are incorrect patient selection and inefficient monitoring during trials. The integration of AI applications into drug discovery has significantly enhanced clinical trial success rates by mitigating issues such as improper patient selection and suboptimal monitoring.

\*Corresponding author

Seyma MINTAS; Izmir Katip Celebi University, Department of Basic Pharmaceutical Sciences, Faculty of Pharmacy, Türkiye; e-mail: [seymamintas1@icloud.com](mailto:seymamintas1@icloud.com);  0009-0000-0340-8552

Canan SEVİMLİ-GUR\*; Izmir Katip Celebi University, Department of Basic Pharmaceutical Sciences, Faculty of Pharmacy, Türkiye; e-mail: [canan.sevimli.gur@ikcu.edu.tr](mailto:canan.sevimli.gur@ikcu.edu.tr);  0000-0002-2210-5925

AI models also analyze relevant parameters such as toxicity, side effects, etc., thus increasing the success rate and thus reducing the cost of clinical trials. In drug discovery and development, AI can generally be used in steps such as cell classification and sorting, calculation of the properties of small molecules, organic compound synthesis of computational tools, design of novel compounds, development of experiments and prediction of the three-dimensional structure of target molecules.

## 2.1. Artificial intelligence applications in literature screening

Nowadays, literature review has become one of the cornerstones of academic and scientific research. However, increasing information density makes it difficult for researchers to carry out this process more efficiently and effectively. In this context, AI technologies optimize literature review processes and provide significant advantages to researchers.

Literature review was a time-consuming process that was usually carried out manually in the past. In traditional searches, researchers could scan a limited number of databases and had difficulty summarizing or analyzing the results. However, the support of AI with technologies such as natural language processing (NLP), ML, and data mining has created a revolutionary transformation in this field.

AI-based systems can scan millions of articles, reports, and data sets in seconds and provide researchers with comprehensive and customized results. These systems are not only used for scanning the literature, but also for tasks such as summarizing, thematic analysis, and trend identification. AI stands out with the following application areas in literature review:

**a. Intelligent Search Engines and Databases:** Platforms such as Google Scholar, PubMed and Semantic Scholar offer the most appropriate results for users' search terms with AI-based algorithms. These systems have the ability to analyze users' past search behavior and provide suggestions.

**b. Automatic Summarization:** AI systems can analyze long articles, summarize them and quickly present the main idea of the article to the researcher. For example, important sections of articles can be identified and summarized in plain language using NLP techniques.

**c. Thematic and Textual Analysis:** ML can classify large numbers of studies under certain themes or concepts. This allows researchers to easily organize related literature and identify critical trends.

**d. Citation and Relationship Analysis:** AI systems can determine which studies are more influential and gaps in the literature by examining citation relationships between studies. For example, it is possible to determine the main studies and subfields in a field with methods such as "Citation Network Analysis".

**e. Systematic Review and Meta-Analysis:** AI-powered tools allow systematic reviews and meta-analyses to be conducted more quickly and objectively. These tools greatly simplify the processes of filtering data, assessing quality, and summarizing results.

AI technologies will be increasingly used in literature review processes and will increase the quality of research. In particular, explainable AI and DL techniques will make systems in this field more user-friendly and reliable. In addition, it will be possible to provide researchers with a more personalized experience by developing AI solutions specific to different disciplines.

In the literature review, AI applications provide great contributions to scientific progress by accelerating academic research. The development of innovative and ethical systems that are suitable for the needs of researchers will ensure that these technologies are used more widely and effectively [13, 14].

## 2.2. Use of artificial intelligence in the preclinical stage

Preclinical studies are a critical phase in which safety and efficacy data are evaluated, forming the basis of drug development and biomedical research. This process is usually costly, time-consuming, and data-intensive. However, in recent years, the use of AI technologies in this field offers great potential to accelerate processes and increase accuracy.

At the preclinical stage, AI is used in areas such as analyzing various types of biomedical data, evaluating drug candidates, and optimizing experiments. The main applications of AI include methods such as ML, DL, and NLP. These technologies offer effective solutions for both in vitro and in vivo studies. The subsequent sections provide a detailed examination of the principal applications of AI in the preclinical stage, including its pivotal contributions to drug candidate identification, experimental design optimization, and other critical processes.

**a. Identification and Optimization of Drug Candidates:** AI can predict new drug candidates that are suitable for biological targets by analyzing large molecular data sets. For example, it is used in ligand-based screening to predict the biological activity of chemical compounds. Additionally, target identification is facilitated through molecular dynamics simulations and bioinformatics analyses, which enable the identification of potential target

proteins and pathways. Moreover, structure-activity relationship (SAR) analysis, powered by machine learning algorithms, accelerates drug optimization by predicting relationships between chemical structure and biological activity.

**b. Toxicity and Safety Assessment:** AI is developing alternative models to reduce toxicity testing in laboratory animals. For instance, *in silico* toxicity prediction leverages AI to predict safety profiles by correlating chemical properties with biological effects. Additionally, the combination of AI with organ-on-a-chip technology, which integrates microfluidic systems, enables more accurate modeling of the toxic responses of human organs.

**c. Cell Culture and Experiment Optimization:** AI-based algorithms analyze processes such as cell proliferation, differentiation, and gene expression in *in vitro* experiments. For example, deep learning methods are employed in cell imaging to classify, count, and analyze cell phenotypes. Furthermore, AI optimizes experimental parameters in experimental design, leading to significant savings in both time and cost.

**d. Studies on Animal Models:** AI enables more efficient analysis of data from animal models. For instance, AI algorithms are used in behavioral analysis to evaluate neurological disorders or the effects of drugs by monitoring animal behavior. Additionally, AI facilitates the analysis of physiological data collected from sensors, allowing for continuous monitoring of biomarkers such as heart rate and respiration.

**e. Biomarker Development:** AI plays an important role in the discovery of new biomarkers for disease diagnosis and treatment targets. Validation of biomarkers can be done more quickly and accurately by analyzing genomic, proteomic and metabolomic data.

The role of AI in preclinical studies is expanding. Explainable AI, in particular, will enable these technologies to gain wider acceptance by providing transparency in regulatory processes. In addition, the integrated use of AI with large data sets obtained from biobanks will enable more sensitive and effective drug development processes.

In preclinical studies, AI offers faster, more cost-effective, and more ethically sustainable solutions than traditional methods. However, in order for these technologies to realize their full potential, challenges such as data quality, algorithm reliability, and regulation must be overcome. AI will continue to play an important role in scientific progress as a critical tool shaping the future of biomedical research [15,16].

### 2.3 Use of artificial intelligence peptide synthesis and small molecule design

Advancements in artificial intelligence have significantly transformed peptide synthesis and small molecule design, two fundamental areas in modern drug development. These innovations are particularly evident in the work of Yan et al. (2020), who developed a Convolutional Neural Networks (CNN)-based platform for the identification of antimicrobial peptides (AMPs) [5].

Peptide synthesis and small molecule design are fundamental to modern drug development processes. However, these processes often involve time-consuming, costly, and complex computational steps. AI is accelerating research processes and increasing accuracy by offering innovative approaches to the design of peptides and optimization of small molecules.

Peptides are short amino acid sequences that play a key role in natural biological processes. Especially in the design and production of therapeutic peptides, AI offers significant advantages with the following applications:

**a. AI in Peptide Design:** AI is used in the design of peptides with desired biological properties. For instance, it enables sequence prediction by identifying amino acid sequences that exhibit targeted biological effects, with deep learning algorithms often used to predict peptides with high binding affinity. Additionally, AI facilitates the optimization of artificial peptides through molecular dynamics simulations and algorithms that suggest modifications to improve the stability and solubility of peptides.

**b. Optimization of Peptide-Synthesis Pathways:** AI can enhance efficiency by optimizing sequencing and synthesis protocols. For example, in Solid Phase Peptide Synthesis (SPPS), AI is used to optimize reaction conditions, increasing efficiency and reducing byproduct formation. Additionally, AI-supported models facilitate reagent prediction by selecting appropriate reagents for chemical reactions and optimizing synthesis conditions.

**c. Development of Therapeutic Peptides:** In the development of therapeutic peptides, AI plays a crucial role by predicting binding affinity to biological targets, analyzing immunogenicity, and optimizing pharmacokinetic properties. Similarly, in the design of small molecules, AI significantly contributes to drug discovery processes by facilitating the creation of molecular structures that interact effectively with biological targets. AI scans the vast chemical space and identifies molecules with high biological activity. For example, Generative Adversarial Networks (GANs) are used to generate new chemical structures, scanning millions of potential small molecules to produce candidates suitable for biological targets. Additionally, AI models predict structure-activity relationships (SAR), guiding the design process by revealing connections between molecular structures and biological activity. AI also optimizes chemical synthesis processes by predicting the outcomes of reactions. Tools such as retrosynthesis analysis automatically determine synthesizable pathways for complex molecules, with

examples including platforms like ASKCOS and Synthia. Furthermore, AI aids in reagent selection by predicting the most suitable reagents for chemical reactions, streamlining synthesis pathways and improving efficiency. In preclinical studies, AI predicts the pharmacokinetic and toxicological profiles of molecules through ADMET (Absorption, Distribution, Metabolism, Elimination, and Toxicity) analysis. This helps guide safer molecule design and reduces the risk of adverse effects. Additionally, machine learning models predict toxic side effects, enabling the creation of chemical structures with improved safety profiles.

AI is supported by numerous advanced tools and technologies that facilitate peptide synthesis and small molecule design. For example, AlphaFold, a deep learning-based tool, has revolutionized protein structure prediction and protein-peptide interaction analysis, providing unprecedented accuracy in understanding molecular dynamics. Schrödinger Maestro offers robust capabilities for small molecule pharmacophore modeling and molecular docking, enabling precise predictions in drug design workflows. Similarly, open-source platforms like DeepChem assist in chemical modeling and design, while ChemProp leverages machine learning for molecular property prediction, further enhancing the efficiency and precision of these processes.

As algorithms and big data technologies continue to evolve, the role of AI in peptide synthesis and small molecule design will expand even further. Explainable AI and hybrid modeling approaches are expected to provide more reliable and practical solutions, ensuring the broader applicability of AI-driven methods in drug discovery. Additionally, the integration of AI algorithms with high-performance computing (HPC) systems is anticipated to significantly enhance the speed and accuracy of molecular simulations, enabling faster and more effective drug development.

Despite its transformative potential, realizing the full benefits of AI in peptide synthesis and small molecule design requires addressing several challenges. Issues such as data quality, algorithm reliability, and regulatory compliance must be carefully managed to ensure the safe and effective deployment of these technologies in drug discovery and development [17, 18].

#### **2.4. Use of AI to define drug dosage and drug delivery efficiency**

Optimization of drug dosage and delivery systems is of critical importance in modern medicine and pharmacology. In order to provide effective and safe treatment, drugs must be administered in the correct dosage and delivered to biological targets with high efficiency. In these processes, AI offers significant innovations with data analytics, predictive modeling and optimization algorithms.

Drug dosage varies according to the physical and biological characteristics of the individual. While traditional dosage determination methods are usually based on clinical trials and data from the general population, AI can make more precise and personalized dosage calculations with individual patient data.

AI plays a crucial role in personalized dosage optimization by analyzing individual patient characteristics such as genetic profile, age, weight, gender, and comorbidities. Traditional methods often fail to capture these individual differences, but AI provides a solution through advanced modeling techniques. Pharmacokinetic and pharmacodynamic (PK/PD) models powered by AI can precisely simulate the absorption, distribution, metabolism, and excretion of drugs, enabling tailored dosage recommendations. Furthermore, AI leverages pharmacogenetic data to assess the impact of genetic variations on drug metabolism and provides personalized dosage recommendations based on genetic factors, enhancing treatment accuracy and safety.

In addition to personalized optimization, AI supports dynamic dosage management by utilizing real-time patient monitoring. AI-powered closed-loop systems can dynamically adjust dosages in response to biological feedback, ensuring optimal therapeutic outcomes while minimizing toxicity risks. These systems also help in therapeutic index optimization by balancing the potential toxic effects of drugs with their therapeutic benefits. Moreover, ML algorithms analyze historical clinical data to predict side effect profiles and identify safe dose ranges, further enhancing patient safety and improving treatment precision.

AI is also revolutionizing drug delivery systems by ensuring that drugs are delivered to their targets at the right dose to maximize therapeutic effect. In targeted drug delivery systems, AI algorithms optimize the binding efficiency of nanoparticles to biological targets, enhancing the precision of nanotechnology-based delivery methods. Additionally, ML-based ligand-selection models identify the most suitable ligands for surface modifications of drug delivery systems, further improving their targeting capabilities.

AI also contributes to the development of controlled release systems, which are designed to provide sustained and controlled drug release over time. It models drug release profiles to maintain therapeutic levels effectively and uses DL algorithms to evaluate the biocompatibility and biodegradability of materials employed in these systems. This ensures both safety and efficiency in drug delivery. Furthermore, AI simulates and analyzes how drugs distribute within the body, predicting their concentrations in specific tissues. Pharmacological modeling

tools predict drug delivery times and concentrations at target sites, while molecular dynamics simulations analyze drug behavior in both intracellular and extracellular environments, enhancing delivery efficiency.

The integration of AI with various advanced technologies has significantly improved drug dosage and delivery systems. DL algorithms are extensively used in PK/PD analyses and toxicity predictions, while Bayesian optimization is employed to refine controlled release systems and analyze clinical data. Reinforcement learning has proven effective in optimizing dynamic dosage management and drug delivery system performance. Additionally, retrosynthesis analysis supported by AI determines ideal chemical pathways for drug synthesis and carrier system design, streamlining drug development processes and enhancing efficiency. The use of AI in drug dosage and delivery systems has the potential to provide more effective, faster and more economical treatment methods in the future. In particular, explainable AI models will increase the trust in AI in clinical decision-making processes. In addition, smart drug delivery systems integrated with biosensors will maximize treatment effectiveness with real-time monitoring and dosage adjustments.

AI is a powerful tool that supports personalized medicine and precision treatment approaches in defining drug dosage and delivery efficiency. The speed, accuracy and cost advantages offered by AI will reshape drug development and application processes. However, data management, regulation and ethical issues need to be addressed in order for these technologies to be implemented effectively and safely [19, 20, 21].

## 2.5. Use of artificial intelligence in predicting bioactive agents and monitoring drug release

Development of bioactive agents and monitoring of drug release processes are key to pharmaceutical research and clinical applications. Conducting these processes using traditional methods can be time-consuming and costly. AI offers powerful tools for predicting bioactive molecules, optimizing drug release mechanisms, and real-time monitoring by analyzing large datasets.

Bioactive agents are molecules that provide therapeutic effects by exhibiting specific interactions in biological systems. AI plays an important role in identifying these molecules and predicting their properties.

AI plays a transformative role in discovering new molecules by predicting bioactive agents through the analysis of large molecular databases. DL algorithms are particularly effective in identifying potential new molecules by learning the relationships between chemical structures and biological activity. Additionally, AI-powered quantum chemistry simulations assist in the selection of effective molecules by predicting their energy profiles and interaction potentials at the molecular level, enabling more accurate predictions in early-stage drug discovery.

Predicting the ADMET properties of bioactive molecules is another critical application of AI in drug development. ML models analyze the physicochemical properties of molecules to predict their ADMET profiles with high accuracy. Furthermore, feature-selection algorithms identify promising candidates by evaluating the effects of molecular structures on toxicity and bioavailability, streamlining the drug development pipeline.

AI also contributes significantly to modeling ligand-target interactions, simulating the interactions of bioactive agents with biological targets to predict mechanisms of action. Docking simulations powered by AI identify molecules with strong interactions by analyzing their molecular binding energies, while tools such as AlphaFold predict protein structures and ligand binding sites, facilitating target-oriented molecule discovery.

In drug delivery systems, AI provides innovative approaches for the design, optimization, and monitoring of controlled release mechanisms. AI-based simulations and data analytics optimize drug release profiles by improving timing and dosage control, ensuring that therapeutic levels are maintained consistently. DL algorithms play a key role in material design by evaluating the biodegradability and biocompatibility properties of drug carrier materials, further enhancing the efficiency of release systems.

AI also enables real-time monitoring of drug release processes by analyzing sensor data. For example, biosensor integration allows for immediate assessment of treatment efficacy by measuring the amount of drug released and biomarker levels in the patient. Closed-loop control systems, powered by AI, dynamically adjust release mechanisms based on real-time data, optimizing the release process to ensure maximum therapeutic benefit.

Furthermore, AI enhances the simulation and visualization of drug distribution within the body. Pharmacokinetic modeling tools optimize release strategies by predicting how drugs are distributed across tissues and how they access their targets. Molecular dynamics analyses add another layer of insight by simulating the intracellular and extracellular release dynamics of drugs, providing a deeper understanding of their behavior in biological systems.

A variety of AI tools and algorithms support these advancements in bioactive agent discovery and drug delivery systems. Machine learning techniques, including regression and classification algorithms, are extensively used to analyze molecular properties and release profiles. Deep learning models are particularly effective in predicting the biological activity and ADMET profiles of molecular structures. Bayesian optimization is employed



to refine drug delivery system designs, while molecular simulation software such as GROMACS and AutoDock, combined with AI-powered pharmacokinetic tools, provide robust platforms for simulation and analysis. These technologies collectively enhance the precision, efficiency, and scalability of drug discovery and delivery processes.

The use of AI in the prediction of bioactive agents and drug release monitoring processes will enable the development of more effective and sensitive treatment methods in the future. Hybrid AI systems will be one of the cornerstones of personalized medicine by being integrated with sensor technologies and bioinformatics tools. In particular, explainable AI will improve clinical decision processes and increase the reliability of AI in the healthcare sector [22, 23, 24].

## 2.6. Use of artificial intelligence in identifying adverse drug reactions

Adverse drug reactions (ADRs) are a term used to describe unexpected and often harmful effects of medications. According to the World Health Organization, ADRs pose a significant burden to the global healthcare system and are estimated to account for 5-10% of hospitalizations. In this context, AI has emerged as a powerful tool for faster, more accurate, and proactive identification and management of ADRs.

AI plays a pivotal role in the identification and prevention of ADRs by utilizing advanced technologies such as big data analytics, ML, DL, and natural language NLP. These technologies enable the comprehensive analysis of vast and complex datasets, offering innovative solutions to enhance patient safety and improve pharmacovigilance systems.

Through data analysis and association identification, AI collects and examines extensive data from sources such as electronic health records, side effect reports, genetic profiles, and even social media. By analyzing these datasets, AI algorithms can uncover previously undetected drug-reaction relationships, providing critical insights into the causes and patterns of ADRs. This capability is particularly valuable for identifying subtle or rare reactions that might otherwise go unnoticed.

AI also contributes significantly to the prediction and risk assessment of ADRs. DL models, in particular, have shown exceptional capability in predicting ADRs for individual patients by analyzing genetic predispositions, environmental factors, and medical histories. These models support the development of personalized treatment plans by identifying potential risks associated with specific drugs before they are administered, thereby improving patient outcomes and minimizing harm.

In the area of ADR monitoring and reporting, AI-driven NLP systems automatically extract and report ADR information from clinical notes, scientific literature, and social media content. Platforms such as Google Health AI and IBM Watson exemplify the application of these technologies, enabling the efficient detection and communication of ADR-related information to healthcare professionals and regulatory authorities.

A range of AI tools supports these efforts, including ML models and algorithms like Support Vector Machines (SVM), Random Forests, and advanced techniques such as XGBoost, which offer high accuracy in ADR classification and prediction across large datasets. Similarly, NLP tools such as SpaCy and BERT are used to extract structured information from unstructured text, with systems like the FDA's Sentinel System leveraging these technologies for pharmacovigilance.

DL applications, including Convolutional Neural Networks (CNNs) and Recurrent Neural Networks (RNNs), further enhance ADR detection by analyzing genetic data and clinical images. Transformer models, such as BERT and GPT, extend this capability by learning from vast datasets, enabling accurate predictions and insights into ADR mechanisms.

AI-powered pharmacovigilance platforms provide an additional layer of support for ADR identification and prevention. For example, VigiBase, the World Health Organization's global pharmacovigilance database, employs AI algorithms to analyze ADR reports from across the globe. Similarly, the FDA's Adverse Event Reporting System (FAERS) integrates AI to track and analyze ADR data efficiently, improving regulatory oversight and patient safety.

By leveraging these advanced technologies, AI significantly enhances the efficiency and accuracy of ADR identification, prediction, and prevention, contributing to safer drug development and more effective pharmacovigilance practices.

AI has the potential to revolutionize ADR identification and prevention. In the future, Prediction accuracy will increase with more data sets and advanced algorithms. ADRs can be monitored instantly with real-time data from IoT devices. Integration of regulatory processes with AI will accelerate drug safety assessments. AI has become an integral part of pharmacovigilance studies and plays a critical role in improving patient safety and treatment success [25, 26, 27].

## 2.7. Prediction of protein folding and protein-protein interactions

Machine and statistical learning approaches such as K-nearest neighbor algorithm, Naive Bayesian algorithm, Support Vector Mechanism, Artificial Neural Networks and Random Forest are used to predict inhibition in PPIs.

## 2.8. Structure-based and ligand-based virtual screening

ML models such as PARASHIFT, HEX, USR and ShaPE algorithms have been developed for LBVS. Tools such as MTiOpenScreen, FlexX-Scan, CompScore, PlayMolecule BindScope, GeauxDock and ENRI have been developed for SBVS.

## 2.9. QSAR modeling and drug redesign

Various algorithms and tools such as VEGA platform, QSAR-Co, FL-QSAR, Meta-QSAR, Transformer-CNN, Cloud 3D-QSAR have been developed for QSAR modeling.

## 2.10. Prediction of mode of action and toxicity of compounds

Different web-based tools such as LimTox, pkCSM, admetSAR and Toxtree can be given as examples.

## 2.11. Identification of molecular pathways and polypharmacology

Text mining-oriented databases such as DisGeNET, STITCH, STRING are widely used to detect gene-disease relationships, drug-target relationships and molecular pathways, respectively.

## 2.12. Application of AI in de novo drug design

In de novo drug design, AI plays a critical role in designing molecules from scratch, optimizing their properties, and ensuring they possess targeted biological activity. Compared to traditional methods, the use of AI in this field offers faster and more cost-effective solutions, streamlining the drug discovery process.

AI algorithms are employed to explore the vast chemical space, enabling the creation of entirely new molecules tailored to specific needs. Among the most prominent methods, Generative Adversarial Networks (GANs) utilize a competitive framework between two neural networks to generate novel molecular structures. Similarly, Variational Autoencoders (VAEs) learn the patterns of chemical structures and use this knowledge to recreate and design new, similar molecules. Reinforcement Learning (RL) further enhances the process by improving the alignment of designed molecules with desired target properties.

These advanced AI techniques enable the generation of unique chemical structures that are optimized for specific biological targets, making de novo drug design a powerful and efficient approach to discovering innovative therapeutic solutions.

## 2.13. Example of neurodegenerative diseases in the use of AI

AI is playing an increasingly significant role in drug discovery, as demonstrated by numerous advancements and applications in the field. For example, Ponzoni et al. (2019) combined decision tree algorithms, quantitative association rules, and hierarchical clustering to identify potential risk genes associated with Alzheimer's disease through gene expression profiling of patient and control samples. They further utilized protein-protein interaction networks, autoencoders, and support vector machines (SVMs) to predict novel target genes associated with Parkinson's disease [6].

The pharmaceutical industry has heavily invested in AI-based applications to streamline drug discovery and development. Xie et al. (2018) developed a model capable of predicting drug-target interactions using transcriptome data from the L1000 database of the Integrated Network-Based Cellular Signatures Library, achieving an impressive accuracy of 98% [7, 8]. Deep learning (DL) methods have also been applied to identify nutraceuticals with anti-aging and anti-cancer properties that mimic the effects of FDA-approved drugs like metformin and rapamycin without causing adverse effects. By mapping the gene-level pathways of these compounds and analyzing over 800 potential alternatives from the LINCS dataset, DL classifiers predicted safer compounds for further development [8].

Prominent technology companies have also contributed to this revolution in drug discovery. Microsoft developed an AI-based system called "Hanover" to assist in identifying optimal cancer treatments by analyzing extensive medical datasets [9]. Similarly, IBM, in collaboration with Pfizer, introduced IBM Watson, a cloud-based platform designed to accelerate drug discovery by analyzing patient-specific data such as medical lab

reports and identifying potential relationships across datasets. The platform also enables personalized treatment plans by engaging patients in dynamic, data-driven interactions with doctors [9].

Koneksa Health has taken a unique approach with its AI-driven software, which integrates mobile and wearable devices to streamline clinical trials by analyzing biomarker data. Biomarkers, which are substances or data indicative of diseases and measurable in body fluids such as blood or urine, are collected and analyzed by this software, expediting the drug development process. The software allows for efficient data sharing with healthcare providers and pharmaceutical companies conducting clinical trials.

Pharmaceutical companies are making substantial investments in AI to improve efficiency and reduce costs. GlaxoSmithKline, for instance, allocated \$43 million to Exscientia, an AI-driven company based in Scotland, to expedite drug development while reducing costs by up to 75% [10]. AstraZeneca is collaborating with Berg to develop biomarkers and treatments for neurological diseases, further illustrating the industry's reliance on AI for innovation [10].

Historical contributions to AI in drug discovery also highlight its foundational importance. Corwin Hansch, known as the "father of computer-aided molecular design," pioneered the use of AI algorithms for predicting the physicochemical properties and biological activities of drug compounds. His work enabled the detailed prediction of chemical structures and their pharmacological properties, laying the groundwork for modern AI-based drug discovery methods [11].

SwissADME has also emerged as a significant resource in the field. This free, user-friendly web tool evaluates the pharmacokinetics and medicinal chemistry properties of small molecules, integrating robust computational methods to facilitate global accessibility and open scientific collaboration [12].

These advancements demonstrate the transformative potential of AI in drug discovery, offering faster, more cost-effective, and highly precise solutions for addressing complex medical challenges.

### 3. Conclusion and Discussion

Although AI seems to be creating and transforming the future of healthcare, it has yet to produce substantial efficacy in certain areas. The lack of any FDA-approved drugs developed solely using AI serves as a significant indicator of its current limitations. This situation is due to both the shortcomings of AI applications compared to humans and the disadvantages that the pharmaceutical industry itself provides to drug discovery.

A key distinction between AI and human intelligence lies in the latter's ability to empathize, a trait AI fundamentally lacks. The most important difference is the ability to empathize. So while various forms of AI have surpassed human performance, they lack higher-level background knowledge and are not as capable of forming relationships as the human brain. They can be trained to do a single task.

The pharmaceutical sector faces significant barriers to AI implementation, including limited open data sharing, inconsistent data formatting, and a disproportion between available 'data' and actionable 'information'. However, overcoming these obstacles is most possible with the development of AI applications. Advancements in ML algorithms and the adoption of DL approaches have significantly enhanced the accuracy and precision of AI applications. However, overcoming these obstacles is most possible with the development of AI applications. The deepening of AI and its transformation into ML algorithms; later on, with the development of these algorithms and the widespread use of the DL approach, the accuracy and precision in AI applications have increased significantly. Considering all these developments, we can say that there is high hope for the elimination of the handicaps we mentioned.

Drug discovery is a complex and lengthy process that often spans many years and costs billions of dollars. However, artificial intelligence (AI) has the potential to transform and accelerate traditional methods. By offering significant innovations at various stages of drug discovery, from molecule design to clinical trials, AI is reshaping drug development processes.

One of the most notable impacts of AI is the shortening of drug development times. Traditional processes typically take 10-15 years, but with AI, this timeline is expected to be reduced by 30-50%. This acceleration will enable faster drug delivery to the market, allowing patients earlier access to treatments. Furthermore, AI is advancing personalized medicine by analyzing genetic variations to develop individual treatment plans. For instance, it can design patient-specific molecules for cancer immunotherapies.

AI also enables the design of challenging molecules by creating new compounds that can target biologically difficult targets, such as protein-protein interactions. This innovation allows researchers to address problems that are difficult to solve with traditional chemistry. Moreover, AI significantly reduces research and development costs, potentially cutting development expenses, which range from \$2-3 billion, by 20-40%.

The success rate of clinical trials may also improve through more accurate predictions of toxicity, side effects, and efficacy profiles of clinical candidates. Additionally, AI's efficiency in drug discovery could increase access to affordable medicines in low-income areas, further broadening the reach of life-saving treatments.

Considering these developments, there is substantial hope for overcoming many of the existing challenges in drug discovery. While persistent obstacles remain, there is little doubt that AI will bring transformative changes to the field in the near future.

### Declaration of Interest

The authors declare that there is no conflict of interest.

### Acknowledgements

No financial support or external contribution was provided for this study.

### Author Contributions

This article was produced from the Graduation Project titled "Artificial Intelligence Applications in Drug Discovery and Research, İzmir Katip Çelebi University, Faculty of Pharmacy, 2023" prepared by Şeyma Mintaş under the supervision of Prof. Dr. Canan SEVİMLİ GÜR. During the preparation of the article, the literature review was done by Şeyma Mintaş, and the writing and publishing processes of the article were done by Canan Sevimli Gür.

### References

- [1] A. N. Ramesh, C. Kambhampati, J. R. T. Monson, P. J. Drew "Artificial intelligence in medicine" *Ann R Coll Surg Engl*, Vol.86 no.5, Sep., pp.334–338, 2004.
- [2] M.K.Tripathi, A. Nath, T.P.Singh, A.S. Ethayathulla, P. Kaur, "Evolving scenario of big data and Artificial Intelligence (AI) in drug discovery", *Molecular Diversity*, vol. 25, no.3, pp. 1440-1446, 2021.
- [3] R. Gupta, D. Srivastava, M. Sahu, S. Tiwari, R.K. Ambasta, P. Kumar, "Artificial intelligence to deep learning: machine intelligence approach for drug discovery", *Molecular Diversity*, vol.25, no.3, pp.1315-1360, 2021.
- [4] M Coşkun, Ö Yıldırım, A Uçar, Y Demır, "An Overview Of Popular Deep Learning Methods", *European Journal of Technique*, Vol. 7, noç 2, pp. 165 – 176, 2017.
- [5] J. Yan, P. Bhadra, A. Li, P. Sethiya, L. Qin, H.K. Tai, K.H. Wong, S.W.I Siu, "Deep-AmPEP30: improve short antimicrobial peptides prediction with deep learning", *Molecular Therapy-Nucleic Acids*, vol. 20, pp.882-894, 2020.
- [6] I. Ponzoni, V. Sebastián-Pérez, M.J. Martínez, C. Roc. "QSAR classification models for predicting the activity of inhibitors of beta-secretase (BACE1) associated with Alzheimer's disease" *Scientific reports*, vol. 9, Article number: 9102, 2019.
- [7] L. Xie, S. He, X. Song, X. Bo, Z. Zhang, "Deep learning-based transcriptome data classification for drug-target interaction prediction", *BMC genomics*, vol. 19, article number 667, 2018.
- [8] İ. N. Çelik, F. K. Arslan, R. Tunç, İ. Yıldız, "İlaç Keşfi ve Geliştirilmesinde Yapay Zekâ", *Journal of Faculty of Pharmacy of Ankara University*, vol. 45, Issue: 2, pp.400 - 427, 2021.
- [9] P. Agrawal, "Artificial Intelligence in Drug Discovery and Development", *Journal of*
- [10] S Büyükgöze, E Dereli. "Dijital sağlık uygulamalarında yapay zeka", VI. Uluslararası Bilimsel ve Mesleki Çalışmalar Kongresi-Fen ve Sağlık, no:4, 2019.
- [11] S. Hochreiter, G. Klambauer, M. Rarey, "Machine learning in drug discovery" *Journal of Chemical Information and Modeling*, vol. 58, Issue 9, pp. 1723 – 1724, 2018.
- [12] A. Daina, O. Michielin, V. Zoete, "SwissADME: a free web tool to evaluate pharmacokinetics, drug-likeness and medicinal chemistry friendliness of small molecules", *Scientific reports*, 7, Article number: 42717, 2017.
- [13] Y Zhang, S Liang, Y Feng, Q Wang, F Sun, S Chen, "Automation of literature screening using machine learning in medical evidence synthesis: a diagnostic test accuracy systematic review protocol", *Systematic reviews*, vol. 11, article number 11, 2022.
- [14] Y Feng, S Liang, Y Zhang, S Chen, "Automated medical literature screening using artificial intelligence: a systematic review and meta-analysis" *Journal of the American Medical Informatics Association*, vol. 29, 8, pp. 1425–1432, August 2022.
- [15] RSK Vijayan, J Kihlberg, JB Cross, V Poongavanam, "Enhancing preclinical drug discovery with artificial intelligence" *Drug discovery today*, vol. 27, Issue 4, 967-984, 2022.
- [16] A Khadela, S Popat, J Ajabiya, D Valu. "AI, ML and other bioinformatics tools for preclinical and clinical development of drug products, *Bioinformatics Tools for Pharmaceutical Drug Product Development*, Chapter. 12, 2023.
- [17] B Lewandowski, G De Bo, JW Ward, M Pappmeyer, "Sequence-Specific Peptide Synthesis by an Artificial Small-Molecule Machine", *Science*, vol 339, 6116, pp. 189-193, 2013.
- [18] M Goles, A Daza, G Cabas-Mora, "Peptide-based drug discovery through artificial intelligence: towards an autonomous design of therapeutic peptides", *Briefings in Bioinformatics*, Volume 25, Issue 4, July 2024.
- [19] EA Poweleit, AA Vinks, T Mizuno, "Artificial intelligence and machine learning approaches to facilitate therapeutic drug management and model-informed precision dosing" *Therapeutic drug monitoring*, 45(2):p 143-150, April 2023.

- [20] KK Mak, YH Wong, MR Pichika, “Artificial intelligence in drug discovery and development”, *Drug Discovery and Evaluation: Safety and Pharmacokinetic Assays*, Springer Nature, pp 1461–1498.
- [21] KS Vidhya, A Sultana, N Kumar, H Rangareddy, “Artificial intelligence's impact on drug discovery and development from bench to bedside”, *Cureus*, vol. 22, 15(10), 2023.
- [22] P Hassanzadeh, F Atyabi, R Dinarvand, “The significance of artificial intelligence in drug delivery system design”, *Advanced drug delivery reviews*, vol. 151–152, pp. 169-190, 2019.
- [23] LK Vora, AD Gholap, K Jetha, RRS Thakur, HK Solanki, “Artificial intelligence in pharmaceutical technology and drug delivery design”, *Pharmaceutics*, vol. 15(7), pp. 1916, 2023
- [24] AI Visan, I Negut, “Integrating Artificial Intelligence for Drug Discovery in the Context of Revolutionizing Drug Delivery” *Life*, vol.14(2), pp.233, 2024.
- [25] S Yang, S Kar, “Application of artificial intelligence and machine learning in early detection of adverse drug reactions (ADRs) and drug-induced toxicity”, *Artificial Intelligence Chemistry*, Vol. 1 (2), pp. 10001, 2023.
- [26] GL Martin, J Jouganous, R Savidan, A Bellec, “Validation of artificial intelligence to support the automatic coding of patient adverse drug reaction reports, using nationwide pharmacovigilance data”, *Drug Safety*, vol. 45, pp 535–548, 2022.
- [27] A Syrowatka, W Song, MG Amato, D Foer, “Key use cases for artificial intelligence to reduce the frequency of adverse drug events: a scoping review”, *The Lancet Digital Health*, vol. 4, Iss. 2, pp.137-e148, 2022.

# Artificial Intelligence Based Customer Risk Classification for Receivables Management of Businesses

Şaban Can TİRYAKİ<sup>1</sup>, Adnan KAVAK<sup>2\*</sup>

## Abstract

This study is carried out with the aim of developing and implementing artificial intelligence-based receivables management systems for businesses. A model is created to predict customers' debt payment situations. In the study, invoice data of a company named QF\_CARIRAPOR is utilized. The features table is created in Apache druid and risk scoring label is made manually according to set rules. Then, various machine learning models such as XGBoost, Random Forest are implemented on MindsDB platform. The classified risk score is visualized with the Streamlit user interface using the results created in MindsDB. Among the applied models, XGBoost has resulted in the highest classification accuracy of 98.8 %. The findings reveal the potential to increase the effectiveness of receivables management processes by applying machine learning models.

**Keywords:** *Mindsdb; Risk Classification; Receivables Management; Xgboost.*

## 1. Introduction

Receivables management is a critical function for businesses, as it directly impacts cash flow, financial stability, and long-term viability. Despite its importance, many organizations face significant challenges in efficiently managing their accounts receivable. Key issues include delayed payments, increased risk of bad debts, and the inability to accurately predict customer payment behavior. Traditional approaches to receivables management, which often rely on manual processes, are becoming increasingly inadequate. These methods struggle to handle the complexities of modern business environments, which are characterized by high transaction volumes, diverse customer profiles, and rapidly changing market dynamics.

Delayed payments can lead to cash flow shortages, forcing businesses to rely on costly financing options. Unmanaged credit risk increases the likelihood of bad debt write-offs and weakening financial health. For small-sized and medium-sized enterprises (SMEs), in particular, these issues can pose existential threats, as they typically operate with narrower financial margins than larger corporations. Complex network theory has offered an innovative approach to assessing credit risk by examining debt and credit relationships in financial systems [1]. Focusing on the use of complex relationship models in the assessment of credit risk, the study in [2] examined the interactions between customer behavior and future payment habits. It was aimed at improving the credit scoring system by using customer segmentation and behavior analysis. Fuzzy rule-based systems are used to manage uncertainties in credit risk assessment. In [3], it was shown that accurate predictions were made by analyzing customer credit history with fuzzy logic. The study in [4] provides a systemic review of the recent studies, identifying trends in credit scoring using a fixed set of questions.

In [5], a literature survey was conducted to systematically review statistical and machine learning models in credit scoring, to identify limitations in literature, to propose a guiding machine learning framework, and to point to emerging directions. Support Vector Machines (SVM) is applied to predict systemic risk in the complex and interconnected realm of financial markets [6]. Deep neural network model was designed to predict high-risk behaviors in financial traders by analyzing vast amounts of transaction data such as Global Insider Trading data [7]. However, there is a need for advanced, data-driven solutions to enhance the accuracy, efficiency, and adaptability of receivables management systems. Specifically, businesses require tools that can classify customer risk more effectively, enabling proactive strategies to mitigate defaults and optimize cash flow. Leveraging artificial intelligence (AI) and machine learning (ML) offers a promising pathway to address these challenges, providing businesses with the ability to process complex data, predict customer behaviors, and implement dynamic, real-time risk management strategies.

The utilization of artificial intelligence (AI) in receivables management may not only enhance the accuracy of risk assessments but also empower businesses to adopt proactive measures aimed at mitigating financial losses. For example, AI-driven systems can provide tailored recommendations for credit terms or automate follow-up

\*Corresponding author

Şaban Can TİRYAKİ<sup>1</sup>; AYASOFYAZILIM Bilişim A.Ş., Bilişim Vadisi, Gebze, Kocaeli, Türkiye & Kocaeli University, Faculty of Engineering, Computer Engineering Department, İzmit, Kocaeli, Türkiye; e-mail: [can.tiryaki@ayasofyazilim.com](mailto:can.tiryaki@ayasofyazilim.com);  0009-0006-2765-5551

Adnan KAVAK<sup>2\*</sup>; Kocaeli University, Faculty of Engineering, Computer Engineering Department, İzmit, Kocaeli, Türkiye, e-mail: [akavak@kocaeli.edu.tr](mailto:akavak@kocaeli.edu.tr);  0000-0001-5694-8042

schedules based on the predicted risk profiles of individual customers. The contribution of this study is to develop an AI-based framework for customer risk classification in receivables management of businesses and demonstrate feasibility of implementation of various machine learning models on MindsDB platform based on real data.

## 2. Dataset

In this study, the invoice information list belonging to a company called QF CARIRAPOR was used (See Table 1). The QF\_CARIRAPOR data set contains invoice information of 1000 customers. The data set contains information such as when customers paid their debts, whether they paid late, on time, and how many days it took to pay. There is a number of invoice information for each customer in the data set. This information in QF CARIRAPOR was combined with a dataset called Top Customer, which includes different information about customers belonging to the same company. Both datasets include the following attributes for each customer: *Paid Invoice*, *Total Paid Invoices*, *Sum Amount Paid Invoices*, *Total Invoices Late*, *Sum Amount Late Invoices*, *Total Outstanding Invoices*, *Total Outstanding Late*, *Sum Total Outstandings*, *Sum Late Outstanding*, *Average Days Late*, *Average Days Outstanding Late*. The merging process was carried out using both data sets. Since the common column in both data sets is the CARDREF column, the merging process was done based on these values. As a result of this merging process, a dataset called Features Table was created. Attributes containing redundant information was eliminated and features were created as shown in the columns of Table 2.

Table 1. A Sample of QF\_CARIRAPOR dataset

LOGID	LOGICALREF	CARDREF	PROJECTREF	PROCDATE	DATE_
954689	18021816	91033	2020-08-12	2020-08-12	2020-08-12
			00:00:00	00:00:00	00:00:00
954690	20485106	91033	2021-06-09	2021-06-09	2021-06-09
			00:00:00	00:00:00	00:00:00
954691	10440444	91033	2018-04-18	2018-04-25	2023-07-05
			00:00:00	00:00:00	00:00:00

Table 2. A Sample of Features Table in QF\_CARIRAPOR dataset

CARDREF	Paid Invoice	Total Paid Invoices	Sum Amount Paid Invoices	Total Invoice Late
8	1	347	5329577000125605	570
10	1	369	573816600039994867	712
19	1	434	884015000000407552	568
34	1	422	447293200030855488	490
36	1	451	38215310003937	559

After combining the datasets, each customer is labelled manually with a risk score according to

$$S = T_{ADL} + P_{ADOS} + C_{TINV}, \quad (1)$$

where  $T_{ADL}$  is the average late day score,  $P_{ADOS}$  is the points received for average amount of overdue debt, and  $C_{TINV}$  is the score for the number of invoices paid on time, which are expressed as;

$$T_{ADL}(t) = \begin{cases} 627 & t = 0 \\ 470 & 1 \leq t \leq 15 \\ 313.5 & 16 \leq t \leq 30 \\ 156.75 & 31 \leq t \leq 45 \\ 0 & t \geq 46 \end{cases} \quad (2)$$

$$P_{ADOS} = 627 \left( 1 - \left( \frac{\text{Amount of unpaid invoices}}{\max(\text{Amount of unpaid invoices})} \right) \right) \quad (3)$$

$$C_{TINV} = \left( \frac{\text{Number of paid invoices on time}}{\text{Total number of invoices}} \right) \quad (4)$$

### 3. Implementation of Customer Risk Classification

Firstly, various Scikit-learn models [8] in Python were applied on the dataset created. While applying the Scikit-learn models, the dataset was divided into training and test sections as 80% and 20%, respectively. A part of the code is given in Figure 1.

```
import pandas as pd
from sklearn.model_selection import train_test_split
from sklearn.preprocessing import StandardScaler, LabelEncoder
from sklearn.metrics import classification_report, confusion_matrix

label_encoder = LabelEncoder()
df['Risk Group'] = label_encoder.fit_transform(df['Risk Group'])

X = df.drop(columns=['Risk Group'])
y = df['Risk Group']

X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)
```

Figure 1. Implementation using Scikit-learn models

Simulation screenshots of various models are given in Figures 2 through 6. As seen in Figure 2, the dataset was first labeled according to the Risk group, according to the calculated risk score values. It was labeled with a numerical result called calculated risk score in Apache Druid. Then, according to these numerical results, risk levels were categorized as *very high risk*, *high risk*, *medium risk*, *low risk*, and *no risk*. The Decision Trees model resulted in classification with 91% accuracy (Figure 2). Random Forest classifier is generally a powerful and flexible ensemble learning method widely used in machine learning. In our problem, it resulted in overall accuracy of 90% (Figure 3). Unfortunately, KNN and SVM classifiers resulted in poor performance with accuracy of 59% and 60%, respectively when our QF\_CARIRAPOR dataset was applied (Figure 4 and 5). The poorer performance of KNN and SVM likely stems from their sensitivity to the characteristics of the financial dataset (e.g., imbalance, noise, and high dimensionality) and their reliance on parameter tuning and preprocessing, which tree-based models handle more gracefully. In addition, KNN assumes that similar instances (in terms of feature values) belong to the same class. If the QF\_CARRAPOR data has complex patterns or nonlinear relationships between features, KNN might fail to capture them. SVM assumes the existence of a clear margin between classes, which might not hold for financial data. The best classification performance is obtained by XGBoost classifier model resulting in accuracy of 93% (Figure 6). The XGBoost excels at modeling the complex interactions in financial data due to its tree-based structure and gradient boosting framework, which builds successive trees to reduce residual errors from previous ones.

Decision Trees- Confusion Matrix:

31	2	0	0	3
2	123	6	0	0
0	3	14	0	0
0	0	0	8	1
0	0	0	1	6

Decision Trees- Classification Report:

	precision	recall	F1-score	support
0	0.94	0.86	0.90	36
1	0.96	0.94	0.95	131
2	0.70	0.82	0.76	17
3	0.89	0.89	0.89	9
4	0.60	0.86	0.71	7
accuracy			<b>0.91</b>	200
macro avg	0.82	0.87	0.84	200
weighted avg	0.92	0.91	0.91	200

Figure 2. Simulation screenshot of Decision Trees using Scikit-learn models



## Random Forest - Confusion Matrix:

32	3	0	0	1
3	123	5	0	0
0	6	11	0	0
0	0	0	9	0
1	0	0	1	5

## Random Forest - Classification Report:

	precision	recall	F1-score	support
0	0.89	0.89	0.89	36
1	0.93	0.94	0.94	131
2	0.69	0.65	0.67	17
3	0.90	1.00	0.95	9
4	0.83	0.71	0.77	7
accuracy			<b>0.90</b>	200
macro avg	0.85	0.84	0.84	200
weighted avg	0.90	0.90	0.90	200

Figure 3. Simulation screenshot of Random Forest using Scikit-learn models

## KNN - Confusion Matrix:

7	28	0	0	1
13	112	0	3	3
1	16	0	0	0
1	8	0	0	0
0	6	0	1	0

## KNN - Classification Report:

	precision	recall	F1-score	support
0	0.32	0.19	0.24	36
1	0.66	0.85	0.74	131
2	0.00	0.00	0.00	17
3	0.00	0.00	0.00	9
4	0.00	0.00	0.00	7
accuracy			<b>0.59</b>	200
macro avg	0.20	0.21	0.20	200
weighted avg	0.49	0.59	0.53	200

Figure 4. Simulation screenshot of KNN using Scikit-learn models

SVM - Confusion Matrix:

0	36	0	0	0
0	131	0	0	0
0	17	0	0	0
0	9	0	0	0
0	7	0	0	0

SVM - Classification Report:

	precision	recall	F1-score	support
0	0.00	0.00	0.00	36
1	0.66	1.00	0.79	131
2	0.00	0.00	0.00	17
3	0.00	0.00	0.00	9
4	0.00	0.00	0.00	7
accuracy			<b>0.66</b>	200
macro avg	0.13	0.20	0.16	200
weighted avg	0.43	0.66	0.52	200

Figure 5. Simulation screenshot of SVM using Scikit-learn models

XGBoost - Confusion Matrix:

30	4	0	1	1
0	127	3	0	1
0	3	14	0	0
0	0	0	8	1
1	0	0	0	6

XGBoost - Classification Report:

	precision	recall	F1-score	support
0	0.97	0.83	0.90	36
1	0.95	0.97	0.96	131
2	0.82	0.82	0.82	17
3	0.89	0.89	0.89	9
4	0.67	0.86	0.75	7
accuracy			<b>0.93</b>	200
macro avg	0.86	0.87	0.86	200
weighted avg	0.93	0.93	0.93	200

Figure 6. Simulation screenshot of XGBoost using Scikit-learn models

Apache Druid [9] and MindsDB [10] have been used for evaluating real-time processing of QF CARİRAPOR dataset. Apache druid is an open-source data storage and analysis platform used for big data analytics. Druid can process and query data in real time. Apache druid also can pull large amounts of data from many different data sources. It can store the retrieved data in a scalable way and then provide rapid access for real-time analysis. Druid's internal architecture is built to provide these fast query and analysis capabilities. MindsDB is an automatic machine learning database that can connect to multiple data sources. It helps to make predictions using the data in the database. The aim of MindsDB is to make data analysis and prediction tasks simple and accessible. The connection process between Apache Druid and MindsDB takes place. In the first stage, a database called "druid\_datasource" is created using the dataset described above. It is stated that this database will use the druid data engine. In the next stage, the information required to connect to the druid data engine is given. Which port is on, which path and scheme is used. MindsDB inherently supports various machine learning algorithms, such as XGBoost, Random Forest, Neural Networks, etc. This allows users to choose the most suitable model option. The platform has the ability to automatically select the most suitable machine learning model based on the dataset, making the model selection process easy. As shown in Figure 7, a folder called druid\_datasource has been created in MindsDB. In this folder, the tables we uploaded to Apache Druid can be seen in MindsDB.

```

1 CREATE DATABASE druid_datasource
2 WITH
3     engine = 'druid',
4     parameters = {
5         "host": "localhost",
6         "port": 8888,
7         "path": "/druid/v2/sql/",
8         "scheme": "http"
9     };
10
11 SELECT * FROM druid_datasource.TopCustomer;
    
```

Figure 7. Druid data source connection in MindsDB

In MindsDB implementation, the table called FeaturesTable in the files section was selected. Using this table, a model called “tahsilet\_sonuc” was created. The desired from this model is to predict the numerical value called TahsiletSkor. It is a command that shows the performance of the applied models. The implementation results on MindsDB real-time database platform is shown in Figure 8. As can be seen, the XGBoost model implementation on MindsDB resulted in the highest accuracy of 98.8%, as in the simulated Scikit-learn models in Python. A sample test case using Streamlit [11] interface using the Collection result model in in MindsDB is given in Figure 9 and 10. As seen in Figure 10, the risk group has been successfully determined based on the information entered by the user.

	name	performance	training_time	selected	accuracy_functions
1	Neural	0.961	55.93	0	['r2_score']
2	XGBoostMixer	0.988	1.01	1	['r2_score']
3	Regression	0.88	0.32	0	['r2_score']
4	RandomForest	0.971	0.52	0	['r2_score']

Figure 8. Classification performance of various models on MindsDB platform applied on QF CARİRAPOR dataset

### Login Information

CARDREF

0

Telephone No

5464513546464

E-mail

PaidInvoice

1.00 - +

TotalPaidInvoices

7.00 - +

SumAmountPaidInvoices

9.00 - +

TotalInvoiceLate

4.00 - +

Figure 9. Streamlit user interface applied on QF CARİRAPOR test dataset

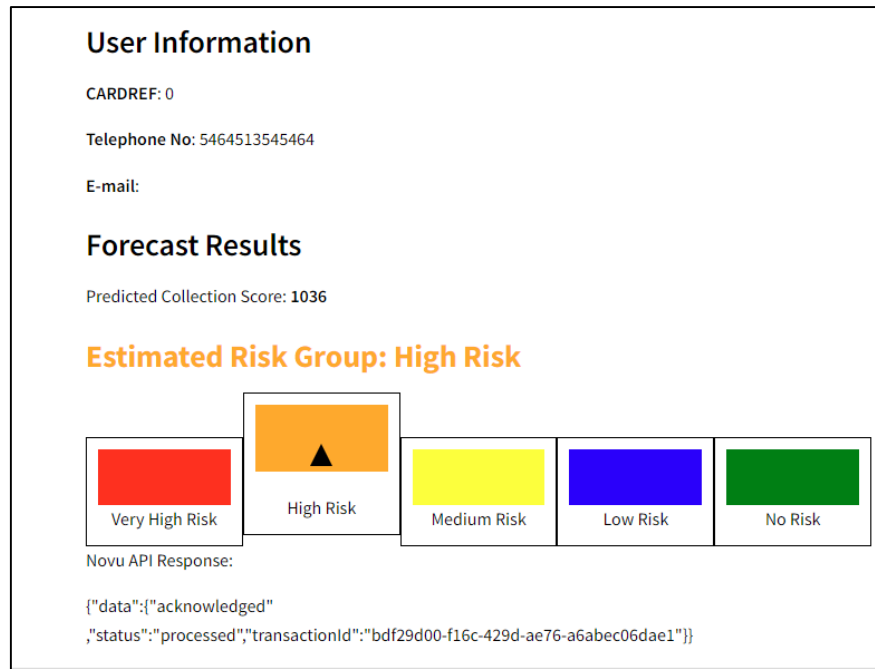


Figure 10. Risk classification when XGBoost is applied on MindsDB using QF CARIRAPOR test dataset.

## 5. Conclusion

In this study, it was aimed to develop and implement artificial intelligence-based receivables management for businesses. A model was created to determine the risk level of a customer whether this customer may pay debts in time or not. For this purpose, a dataset called QF\_CARIRAPOR was utilized in our study, which contains data of 1000 customers. This dataset was then uploaded to the Apache Druid environment. It is labeled with numbers between 0 and 1900 called Apache Druid Collection Score. Then, we applied classification models such as XGBoost, Random Forest, KNN, SVM, and Decision Trees to the dataset. The model with the highest accuracy rate was named Collection Result into the MindsDB environment and created a model named Tahsilet\_Sonuc. The XGBoost model resulted in the highest classification accuracy in both the Scikit-learn simulation (93% accuracy) and the MindsDB real-time database implementation (98.8% accuracy).

## Acknowledgements

This study was carried out within the scope of TAHSILET project conducted by AYASOFYAZILIM A.Ş. The project is supported by Tübitak Teydeb 1507 program.

## References

- [1] H.Lam, "Analyzing the Measures of Credit Risk on Financial Corporation and It's Impact on Profitability," International Journal of Research in Vocational Studies (IJRVOCAS), vol. 3, no. 1, pp. 64-70, 2023.
- [2] N. Wilson, B. Summers, R. Hope, "Using payment behaviour data for credit risk modelling," International Journal of the Economics of Business; vol. 7, no. 3, pp. 33-346, 2000.
- [3] J. Reyes, J. Perez, and S. Ake, "Credit risk management analysis: An application of fuzzy theory to forecast the probability of default in a financial institution," Contaduría y Administración, vol. 69, no. 1, pp. 18-211, 2024.
- [4] A. Markov, Z. Seleznyova, and V. Lapshin, "Credit scoring methods: Latest trends and points to consider," The Journal of Finance and Data Science, vol. 8, pp. 180-201, 2022.
- [5] X. Dastile, T. Celik, and M. Potsane, "Statistical and machine learning models in credit scoring: A systematic literature survey," Applied Soft Computing, vol. 91, 106263, 2000.
- [6] Q. Zhou, "Predicting Systemic Risk in Financial Markets Using Machine Learning," Transactions on Economics Business and Management Research vol. 8, pp. 455-460, 2024.
- [7] K. Xu, Y. Wu, Z. Li, R. Zhang, and Z. Feng, "Investigating Financial Risk Behavior Prediction Using Deep Learning and Big Data," International Journal of Innovative Research in Engineering and Management (IJIREM), vol. 11, no. 3, pp. 77-81, 2024.
- [8] Scikit-learn Machine Learning in Python, <https://scikit-learn.org/>
- [9] <https://druid.apache.org/>
- [10] MindsDB-Platform for Building AI, <https://docs.mindsdb.com/>
- [11] <https://streamlit.io/>

# Production of a Cost Effective Microstrip Antenna Operating at 2.4 GHz and 5 GHz

Burak DOKMETAS <sup>1\*</sup>

## Abstract

In this study, the aim was to design a microstrip antenna operating in the ISM band using three-dimensional printer technology. The substrate material of the proposed antenna was produced using ABS filament with a 70% infill percentage. The dielectric constant of the produced substrate material was measured experimentally in the laboratory and the Reel-Imaginer graph of the material was extracted. According to the obtained dielectric constant value, the antenna was designed step by step in the simulation environment. Then, the antenna was produced and return loss ( $S_{11}$  parameters) were measured using spectrum analyzer. It was observed that the simulation and measurement results were consistent. The measurement results shows that the produced antenna operates at 2.4 GHz and 5 GHz. This study demonstrates the simplicity of antenna production with a low-cost production technique, 3D printing.

**Keywords:** *Microstrip Antenna; 3D printer; ABS filament; 2.4 GHz ;Wifi.*

## 1. Introduction

Nowadays, with the widespread use of wireless communication, interest in applications especially in the WiFi frequency band is increasing. Popular areas such as the Internet of Things and fifth generation communication applications encourage researchers to solve problems that may be encountered in these areas [1,2]. The most commonly preferred antenna type in communication is the microstrip antenna [3]. Their compact structure and flexible design allow them to be used frequently in microwave applications. The biggest disadvantage of traditional microstrip antennas is their production technique. They are often produced in factories and workshops with a device called LPKF. These devices are not the kind of equipment that researchers can easily obtain in terms of cost and size. Another disadvantage of the traditional production method is that the substrate materials are mostly imported from abroad and are costly. In addition, this production technique has certain limitations and the design of complex structures is not possible with these devices [3].

Researchers who are looking for an innovative production technique propose 3D printing technology. With this technology, complex 3D antenna structures can be produced quickly and easily. Due to such advantages, the use of 3D printers in antenna production has attracted great interest from both academic and industrial researchers [4]. In addition, the filament material used as raw material in 3D printers is economically feasible. Researchers have successfully produced various microwave elements such as, slotted array antenna, SIW (substrate integrated waveguide) antennas, dielectric lenses, Ku band horn antenna and folded microstrip antenna using this new production technology [5-10]. In recent years, in addition to mobile phones, many technological devices such as tablets and smart watches have been connected to the internet. As the number of devices that can connect to the internet increases every day, the need for WiFi is also increasing.

All devices that support WiFi technology connect to the local network, via wireless access points. The connection is made at 2.4 GHz or 5 GHz radio frequency depending on the IEEE 802.11 protocol, which is supported by wireless access points and the device. 2.4 GHz offers a wide range of signal, whereas 5 GHz has a shorter range. Researchers have done various studies for the ISM band. A 5 cm x 5 cm microstrip antenna with dual band operation at 2.45/5.8 GHz was designed [11]. In another article, a fractal slot antenna with an operating frequency between 2.5GHz and 4.5GHz was designed [12].

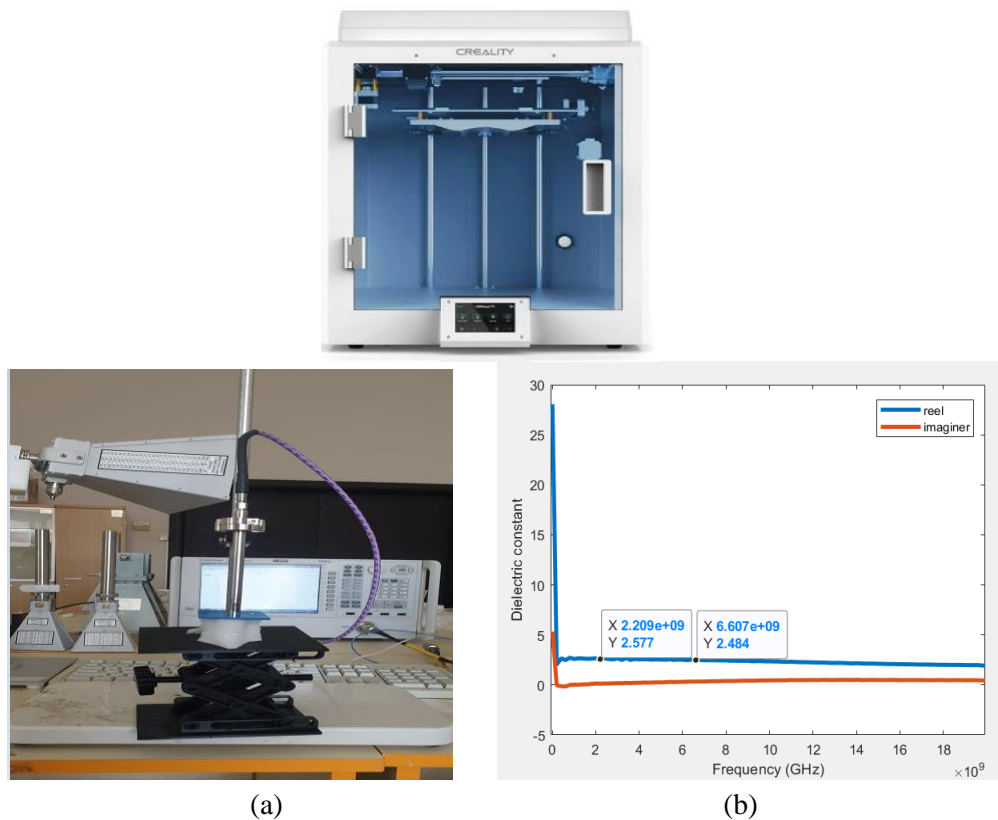
Within the scope of this study, a microstrip antenna was designed and manufactured for the wireless communication frequency band. The study started with determining the dielectric constant of the ABS filament to be used as raw material. The proposed design was produced in 3D printer and the performance of the resulting prototype was measured with spectrum analyzer.

\*Corresponding author

BURAK DOKMETAS: Kafkas University, Faculty of Engineering and Architecture, Electrical Electronic Engineering Department, Türkiye; e-mail: [burakd@kafkas.edu.tr](mailto:burakd@kafkas.edu.tr);  0000-0001-5900-6691

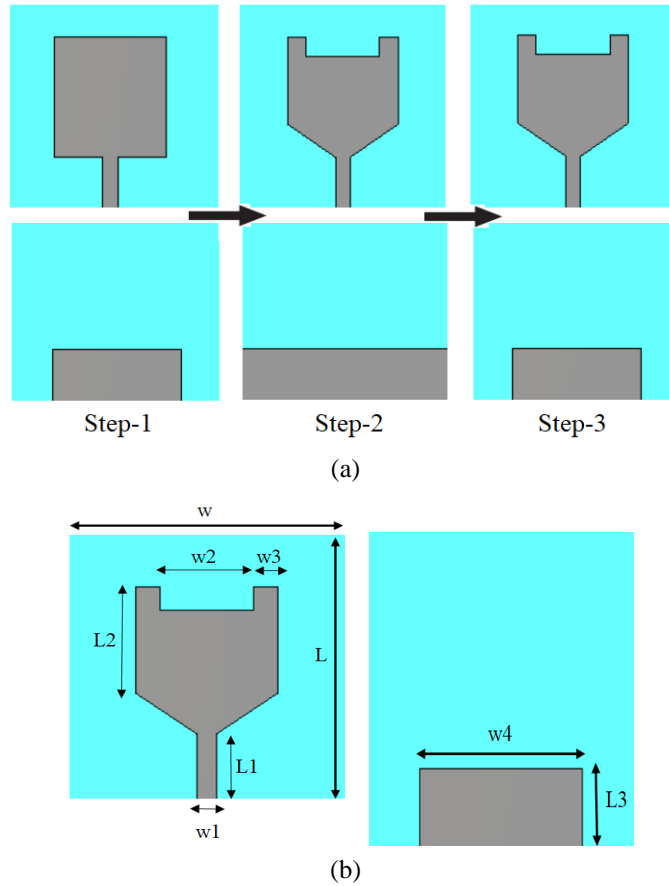
## 2. Substrate Production Process and Antenna Design

In this section, firstly the fabrication and measurement process of the antenna substrate material will be explained. ABS filament, which is prominent with its durability, was preferred as the substrate material. The recommended substrate was produced using a Creality Cr 5 Pro-H type 3D printer with dimensions of  $60 \times 60 \times 1.6 \text{ mm}^3$ . The infill percentage of the filament was set to 70% using the software program of the 3D printer. During the printing phase, the bed temperature of the 3D printer was set to  $110 \text{ }^\circ\text{C}$  and the nozzle temperature was set to  $240 \text{ }^\circ\text{C}$ . The printing process was completed in approximately 50 minutes. After the production process, the measurement setup given in Figure 1(a) was established to determine the dielectric value of the substrate material. The substrate was placed between two conductors and connected to the measurement device. The dielectric values of water and air were used as references during the calibration phase. The graph in Figure 1(b) was obtained by taking measurements at every 100 MHz intervals in the frequency range between 100 MHz and 18 GHz. According to Figure 1(b), the real part of the dielectric constant was measured as  $\epsilon_r = 2.5$  on average, especially between 2 GHz and 6 GHz.



**Figure 1.** Images of the Substrate Production Process (a) Measurement set-up (b) Graph of the measured dielectric constant

As a result of the measurements, the dielectric constant of the substrate material was determined. After this stage, the focus was on the antenna design using the simulation software. As a result of the measurements, the dielectric constant of the substrate material was determined. After this stage, the focus was on the antenna design using the simulation program. The steps given in Figure 2(a) were carried out in order to obtain the desired frequency band range. The front and back views of the final design of the antenna are given in Figure 2(b). The final dimensions optimized in the simulation program are summarized in Table 1.



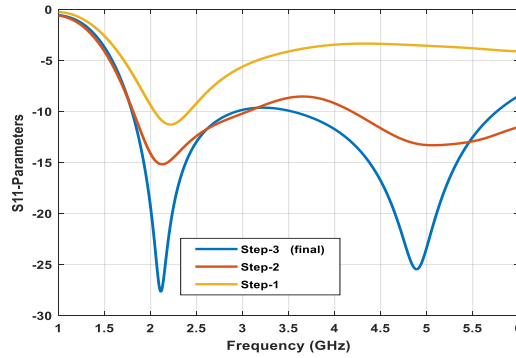
**Figure 2.** Images of the Proposed Microstrip Antenna (a) Design steps (b) Layout of the design

**Table 1.** Optimized dimensions of the final antenna (mm)

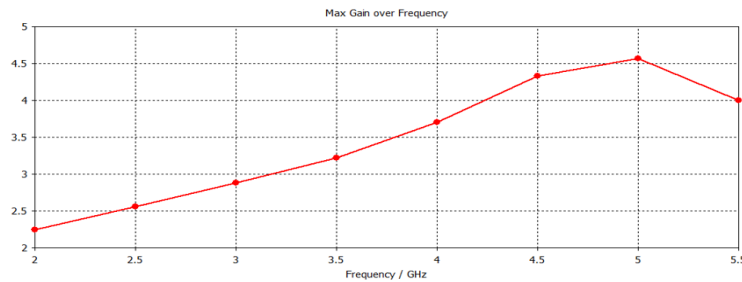
<b>W</b>	60	<b>L</b>	60
<b>W1</b>	4	<b>L1</b>	15
<b>W2</b>	20	<b>L2</b>	23.6
<b>W3</b>	5	<b>L3</b>	15
<b>W4</b>	35		

### 3. Simulation and Measurement Results

In this section, firstly the simulation results of the designed antenna will be analyzed. In Figure 3(a), the return loss ( $S_{11}$ -parameters) graph obtained in the design phase of the antenna is given comparatively. In the first phase, the antenna resonates only at 2.4 GHz. In the second phase, a second resonance is produced at 5 GHz and the bandwidth increases. In the final design, the antenna provides the desired bandwidth between 1.8 GHz and 5.8 GHz. In Figure 3(b), the maximum gain over frequency graph of the antenna is given. The antenna gives a gain of 2.5 dBi at 2.4 GHz and 4.5 dBi at 5 GHz. The radiation loss of the antenna at 2.4 GHz and 5 GHz is given in Figures 4(a) and 4(b), respectively.

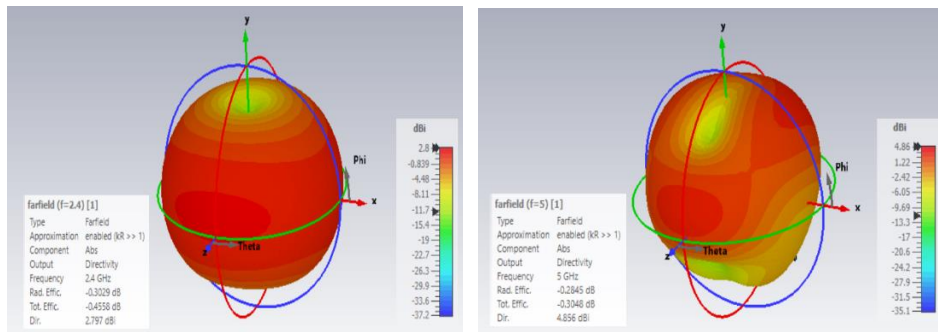


(a)



(b)

Figure 3. Simulation results (a) Return loss (S<sub>11</sub>), (b) Max. Gain over frequency



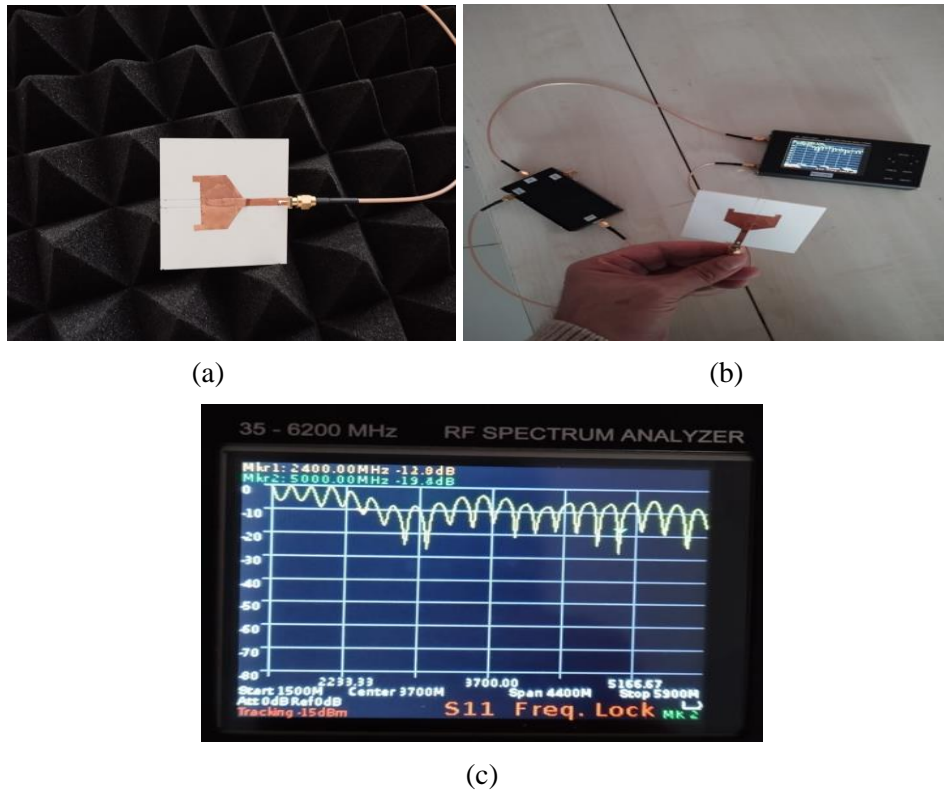
(a)

(b)

Figure 4. Simulated 3D radiation pattern (a) at 2.4 GHz (b) 5 GHz

Once the desired efficiency was obtained from the simulation results, the production and measurement phase of the antenna was started. 35 micron thick adhesive copper tape was used for the conductive parts of the antenna. The measurements of the produced antenna were made with the ARINST SSA-TG R2 spectrum analyser device. This device is an easy portable panoramic spectrum analyzer (RF multimeter) designed to display signal spectrums in the frequency range from 35 to 6200 MHz. One of the features of this device is the presence of an internal tracking generator, which makes it possible to measure Standing Wave Ratio (SWR) and return loss of the antenna. The picture of the manufactured antenna is given in Figure 5(a). The measurement setup using the Radio Frequency (RF) bridge and spectrum analyzer is given in Figure 5(b). The measurement results obtained from the spectrum analyzer are given in Figure 5(c). The tracking input attenuation value on the device is set to -10 dB and the generator power value is set to -15 dBm. Before the antenna measurement, the spectrum analyzer was calibrated. For this purpose, open, short and load (50 Ω) were connected to the Device Under Test (DUT) port, respectively. As seen on the measurement screen, -12 dB return loss was obtained at 2.4 GHz and -19.4 dB at 5 GHz.





**Figure 5.** Images of the measurement (a) Prototype of the antenna (b) Measurement set-up (c) Display of spectrum analyzer

#### 4. Conclusion

One of the important hardware parts of wireless communication is the antenna. It is clear that with the development of infrastructures such as the Internet of Things, the importance of next-generation antenna production techniques and studies on improving antenna performance will increase. For this purpose, the subject of the study was determined and a microstrip antenna was produced using three-dimensional printer. It was observed that the measurement results and simulation results were consistent. According to the measurement results the fabricated antenna is suitable for use in 2.4 GHz and 5 GHz applications. This study has shown that 3D printers can play an alternative role in antenna production. At the same time, this study has proven that ABS filaments can be used without any problems in ISM band applications.

#### Declaration of Interest

The authors declare that there is no conflict of interest.

#### Author Contributions

Conceptualization, BD; methodology, BD ; data generation, BD; investigation, BD; designing BD; writing—original draft preparation, BD; writing—review and editing, BD; visualization, BD; supervision, BD; project administration, BD.

#### Acknowledge

This study is supported by Kafkas University Scientific Research Project within the scope of project number 2022-FM-72.

#### References

- [1] Balanis, C.A., *Antenna Theory: Analysis and Design*. John Wiley & Sons, New Jersey, A.B.D., 2005.
- [2] Garg, R., Bhartia, P., Bahl, I., ve Ittipiboon, A., “*Microstrip Antenna Design Handbook*”. Artech House, Boston, A.B.D., 2001.
- [3] D. Cao, Y. Li and J. Wang, “A Millimeter-Wave Spoof Surface Plasmon Polaritons-Fed Microstrip Patch Antenna Array”, in *IEEE Transactions on Antennas and Propagation*, vol. 68, no. 9, pp. 6811-6815, 2020.
- [4] Olan-Nuñez KN, Murphy-Arteaga RS, “Dual- band antenna on 3D-printed substrate for 2.4/5.8 GHz ISM-band applications.” *Electronics*, 12(11):2368,2023.

- [5] Bjorgaard J, Hoyack M, Huber E, Mirzaee M, Chang YH, Noghianian S., "Design and fabrication of antennas using 3d printing". *Progress In Electromagnetics Research C*, 84:119-134, 2018.
- [6] Zhang, S., Arya, R.K., Pandey, S., et al., "3D printed planar graded index lenses", *Microw. Antennas Propagation.*, 1411-1419, 2016.
- [7] B. Dokmetas and G. O. Arican, "Design of Dual-Band SIW Antenna for Millimeter-Wave Communication", 31st Telecommunications Forum (TELFOR), Belgrade, Serbia, pp. 1-4, 2023.
- [8] Ghazali M.I.M., Karuppuswami S., Kaur, A, "3D printed air substrates for the design and fabrication of RF components", *Trans. Compon. Packag. Manuf. Technol.*, 982-989, 2017.
- [9] Jun S., Sanz-Izquierdo B., Heirons J., "Circular polarised antenna fabricated with low-cost 3D and inkjet printing equipment", *Electronic Letters.*, 370-371, 2017.
- [10] Belen MA., Güneş F, Mahouti P, Belen A, "UWB Gain Enhancement of Horn Antennas Using Miniaturized Frequency Selective Surface", *Applied Computational Electromagnetics Society Journal*, 997-1002, 2018.
- [11] B. Dokmetas, G. O. Arican and B. A. Yilmaz, "A Folded Pyramid-Shaped Microstrip Antenna with Improved Bandwidth", *IEEE International Symposium on Antennas and Propagation and INC/USNC-URSI Radio Science Meeting (AP-S/INC-USNC-URSI)*, Firenze, Italy, pp. 1273-1274, 2024.
- [12] Malakooti S.-A., Hayati, M., Fahimi, V., Afzali, B., "Generalized dual-band branch-line coupler with arbitrary power division ratios." *International Journal of Microwave and Wireless Technologies*, 1-9, 2015.
- [13] Krishna, D.D., Gopikrishna, M., Anandan, C., Mohanan, P., Vasudevan, K., "CPW-fed Koch fractal slot antenna for WLAN/WiMAX applications." *IEEE Antennas and Wireless Propagation Letters*, 7, 389-392, 2008.

# Machine Learning Approaches for Prediction of Alzheimer's Disease

Kadriye Filiz BALBAL <sup>1\*</sup>

## Abstract

Alzheimer's Disease (AD) is a disorder that significantly impacts an individual's behavior, memory, and cognitive functions, ultimately leading to a loss of independence. Early and accurate diagnosis of AD is critical to mitigating its progression and improving patient outcomes, especially as no definitive cure is currently available. This study investigates the application of machine learning algorithms to predict and diagnose AD based on patient symptoms and clinical data. The dataset used in this research includes comprehensive health information from 2,149 patients, with 35 features covering demographic, lifestyle, and medical factors, and no missing values. Seven widely recognized machine learning algorithms—KNN, GNB, SVM, DT, RF, AdaBoost, and XGBoost—were evaluated to determine their effectiveness in disease prediction. Performance was assessed using recall, precision, accuracy, and F1-score metrics, providing a robust evaluation of each model. XGBoost achieved the highest accuracy rate of 95.35%, highlighting its superior predictive capability, while KNN recorded the lowest accuracy at 75.54%. The results demonstrate the strength of machine learning algorithms, particularly ensemble methods like XGBoost, in analyzing complex clinical data for the early detection of Alzheimer's Disease. These findings underscore the critical role of machine learning in enhancing diagnostic accuracy and enabling timely interventions, which are essential for improving the quality of life for individuals at risk of Alzheimer's Disease.

**Keywords:** *Alzheimer's disease; classification; disease prediction; machine learning.*

## 1. Introduction

Alzheimer's Disease (AD) is a progressive neurodegenerative disorder that often manifests with symptoms resembling the natural effects of aging, making it challenging to distinguish in its early stages. Symptoms such as forgetfulness, confusion, or behavioral changes like stress and paranoia can overlap with other conditions or be dismissed as normal aging processes. Primarily affecting individuals aged 65 and older, AD progressively deteriorates cognitive functions, eventually impairing the ability to perform daily activities independently [1]. The disease profoundly impacts memory, thinking, and reasoning abilities, gradually diminishing the quality of life for both patients and their caregivers. Despite extensive research, there is currently no definitive cure for AD, and it remains a significant global health challenge. However, evidence suggests that early and accurate diagnosis, coupled with timely interventions, can slow the disease's progression and alleviate its symptoms [2]. By delaying the onset of severe cognitive decline, early diagnosis not only enhances patient outcomes but also reduces the emotional and financial burden on families and healthcare systems. Therefore, identifying AD in its early stages through advanced diagnostic techniques is of paramount importance, as it enables effective management strategies and improves the quality of life for affected individuals [3].

In this study, the prediction of Alzheimer's Disease (AD) status was conducted using machine learning algorithms that are widely recognized in the literature for their effectiveness in disease diagnosis and prediction. The analyses aimed to determine the presence of AD by utilizing various data types, including demographic information, medical history, cognitive and functional evaluations, lifestyle factors, clinical measurements, and patient-reported symptoms. The structure of the paper is organized as follows: Section 2 provides a review of related work, Section 3 details the dataset description and the methodology employed, Section 4 presents the results of the analyses, and Section 5 concludes the study with a discussion of the findings and their implications.

## 2. Related Works

The role of machine learning (ML) in disease prediction, early diagnosis, and prevention has been extensively emphasized in the literature, highlighting its potential to revolutionize healthcare by improving diagnostic accuracy and enabling timely interventions [4-10]. Numerous studies have explored the application of ML techniques to analyze complex medical datasets, demonstrating their versatility and effectiveness. For example, in [5], where the critical importance of early diagnosis in Alzheimer's Disease (AD) was underscored, predictions were performed using DT and GNB models, showcasing their ability to identify patterns associated with AD.

\*Corresponding author

Similarly, in [6], an audio dataset from the UCI Machine Learning Repository was employed to predict diseases using four different ML algorithms. This study demonstrated the adaptability of ML techniques to various data modalities, such as audio features, further expanding their applicability in medical diagnostics.

In [7], a broader investigation was undertaken by predicting diseases such as diabetes, breast cancer, and heart disease using DT and GNB models. The study aimed to identify the most effective algorithm for disease prediction through a comparative analysis, emphasizing the cross-domain applicability of these methods in addressing diverse medical conditions. Furthermore, [8] utilized two supervised ML algorithms for disease prediction, achieving 87% accuracy with GNB and 91% with DT, further substantiating the capability of these methods in clinical data analysis.

In [9], a comprehensive evaluation was conducted using six classification algorithms: Multi-Layer Perceptron (MLP), Logistic Regression (LR), Extremely Randomized Trees Classifier (ERT), SVM, RF, and Gradient Boosting Classifier (GBC), to classify heart disease. This study, performed on the Cleveland dataset containing 14 features, revealed that MLP and SVM achieved the highest performance with an accuracy rate of 91.7%. These findings highlighted the potential of advanced ML algorithms in achieving high precision and reliability in disease classification.

Additionally, in [10], ML algorithms were applied to hospital data collected in Andhra Pradesh, India, between 2018 and 2020, to evaluate their effectiveness in disease prediction. Among the algorithms tested, AdaBoost and KNN were identified as the most successful, demonstrating their ability to extract valuable insights from real-world clinical datasets. Collectively, these studies underscore the importance of ML in disease prediction, highlighting the role of algorithm selection, dataset characteristics, and feature engineering in optimizing model performance. Such findings reveal the transformative potential of ML techniques in healthcare, particularly in addressing diagnostic challenges and enhancing early detection strategies.

### 3. Data and Methodology

#### 3.1. Dataset

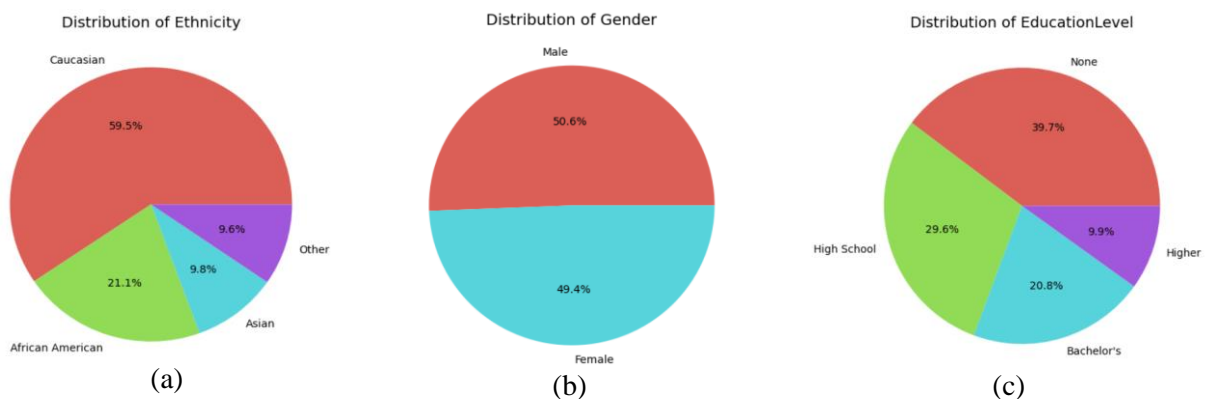
The Alzheimer's Disease dataset was obtained from Kaggle [11]. The dataset with no missing values contains a total of 2,149 rows and 35 columns. The columns in the dataset include information on individuals' demographic characteristics, lifestyle factors, health history, and cognitive status. 16 of the variables are categorical (such as gender, smoking, Alzheimer's diagnosis), and 19 are numerical (such as age, BMI, cholesterol levels). While categorical variables allow patients to be divided into groups, numerical variables allow for more detailed analyses. The dataset provides a wide range of data to understand the risk factors and symptoms of Alzheimer's disease. This diversity provides a solid basis for both descriptive analyses and more advanced statistical methods. Detailed information about the dataset is provided in Table 1.

**Table 1.** *Alzheimer's Disease Dataset Description*

Category	Subcategory	Description
Patient Identification	PatientID	Unique identifier assigned to each patient (4751-6900).
Demographic Details	Age	Patients' ages range between 60 and 90 years.
	Gender	0: Male, 1: Female.
	Ethnicity	0: Caucasian, 1: African-American, 2: Asian, 3: Other.
	Education Level	0: None, 1: High School, 2: Undergraduate, 3: Higher.
Lifestyle Factors	BMI	Body Mass Index ranging from 15 to 40.
	Smoking	0: No, 1: Yes.
	AlcoholConsumption	Weekly alcohol consumption ranging from 0 to 20 units.
	PhysicalActivity	Weekly physical activity in hours, ranging from 0 to 10.
	DietQuality	Diet quality score ranging from 0 to 10.
Medical History	SleepQuality	Sleep quality score ranging from 4 to 10.
	FamilyHistoryAlzheimer	0: No, 1: Yes (family history of Alzheimer's disease).
	Cardiovascular Disease	0: No, 1: Yes.
	Diabetes	0: No, 1: Yes.

	Depression Head Injury Hypertension	0: No, 1: Yes. 0: No, 1: Yes (history of head injury). 0: No, 1: Yes.
Clinical Measurements	SystolicBP DiastolicBP TotalCholesterol LDLCholesterol HDLCholesterol Triglycerides	Systolic blood pressure ranging from 90 to 180 mmHg. Diastolic blood pressure ranging from 60 to 120 mmHg. Total cholesterol levels ranging from 150 to 300 mg/dL. Low-density lipoprotein cholesterol levels ranging from 50 to 200 mg/dL. High-density lipoprotein cholesterol levels ranging from 20 to 100 mg/dL. Triglyceride levels ranging from 50 to 400 mg/dL.
Cognitive and Functional Assessments	MMSE FunctionalAssessment MemoryComplaints BehavioralIssues ADL	Mini-Mental State Examination score ranging from 0 to 30 (lower scores indicate cognitive impairment). Functional assessment score ranging from 0 to 10 (lower scores indicate greater impairment). 0: No, 1: Yes (presence of memory complaints). 0: No, 1: Yes (presence of behavioral issues). Activities of Daily Living score ranging from 0 to 10 (lower scores indicate greater impairment).
Symptoms	Confusion Disorientation PersonalityChanges DifficultyCompletingTasks Forgetfulness	0: No, 1: Yes. 0: No, 1: Yes. 0: No, 1: Yes. 0: No, 1: Yes. 0: No, 1: Yes.
Diagnosis Information	Diagnosis	0: No, 1: Yes (Alzheimer's diagnosis).
Confidential Information	DoctorInCharge	Information about the doctor in charge (specified as 'XXXConfid' for all patients).

Before proceeding to the analysis, the 'PatientID', and 'DoctorInCharge' columns, which were not important for the analysis, were deleted. The 'Diagnosis Information' column was determined as the target variable and the analysis continued with 32 columns. The distribution of the demographic features Gender, Ethnicity, and Educational Level in the dataset is given in Figure 1.



**Figure 1.** *Distribution of Demographic Features.*

According to Figure 1, 50.6% of the participants in the dataset are male and 49.4% are female. In terms of ethnicity, most of the participants (59.5%) are Caucasian, 21.1% are African American, 9.8% are Asian, and 9.6% are Other. In the dataset, there are 29.6% with High School education level and 20.8% with Bachelor's education level. In addition, 9.9% have Higher education level and 39.7% have neither of these education levels.

### 3.2. Methodology

This section introduces the machine learning algorithms employed for predicting AD. This study utilized seven widely recognized machine learning algorithms: DT, SVM, GNB, KNN, RF, AdaBoost, and XGBoost. Each algorithm's details and characteristics are outlined below.

K-Nearest Neighbors (KNN) is a supervised learning algorithm commonly applied to regression and classification tasks. The algorithm's performance is governed by the parameter  $k$ , which specifies the number of neighbors considered in predictions. While KNN is advantageous for its simplicity and low computational requirements, its performance deteriorates when applied to large datasets due to increased computational costs [12, 13].

Gaussian Naïve Bayes (GNB) operates on the assumption that the features within a class are distributed according to a Gaussian distribution. During training, the algorithm calculates the mean and standard deviation for each class and uses these parameters to estimate probabilities for continuous variables [14].

Support Vector Machine (SVM) is a robust algorithm particularly effective for classification problems in high-dimensional datasets. It supports both linear and nonlinear classification through the use of kernel functions, enabling adaptability to various data structures. However, SVM is computationally intensive and requires complex mathematical operations [15].

Decision Tree (DT) is a tree-based algorithm capable of handling both numerical and categorical data. It excels in capturing intricate interactions between variables and demonstrates robustness against certain levels of noise and inconsistencies in the data [16].

Random Forest (RF) is a powerful ensemble learning algorithm that constructs multiple decision trees to mitigate the high variance observed in single-tree models. It is particularly effective for handling high-dimensional datasets, providing stable and reliable classification performance [17].

Adaptive Boosting (AdaBoost) is an ensemble technique designed to improve the classification performance of weak learners, particularly in scenarios with imbalanced datasets. Its adaptive nature iteratively adjusts the weights of misclassified instances to enhance overall accuracy [18].

Extreme Gradient Boosting (XGBoost) is a high-performance gradient boosting library widely recognized for its efficiency and scalability. XGBoost is particularly effective in classification and prediction problems due to its fast processing capabilities, extensibility, and robust generalization properties [19].

The implementation of these machine learning algorithms for AD prediction was carried out using Python in the Google Colab environment. The dataset, obtained from the Kaggle platform in CSV format, was processed and analyzed using the pandas and numpy libraries. Performance evaluation metrics were imported from the sklearn library, while graphical visualizations were generated using the matplotlib and seaborn libraries. These tools provided a comprehensive framework for data preprocessing, algorithm implementation, and result interpretation in the study.

### 3.3. Metrics

To evaluate the performance of machine learning algorithms, recall, precision, F1 score, and accuracy metrics were used. These metrics are four evaluation criteria that are widely used in prediction and classification problems [20]. As seen in Equations (1)-(4), the metrics are calculated depending on the values of TN (true negative), TP (true positive), FN (false negative), and FP (false positive).

$$Accuracy = \frac{TP+TN}{TP+TN+FP+FN} \quad (1)$$

$$Precision = \frac{TP}{TP+FP} \quad (2)$$

$$Recall = \frac{TP}{TP+FN} \quad (3)$$

$$F1 \text{ score} = 2 \times \frac{Precision \times Recall}{Precision + Recall} \quad (4)$$

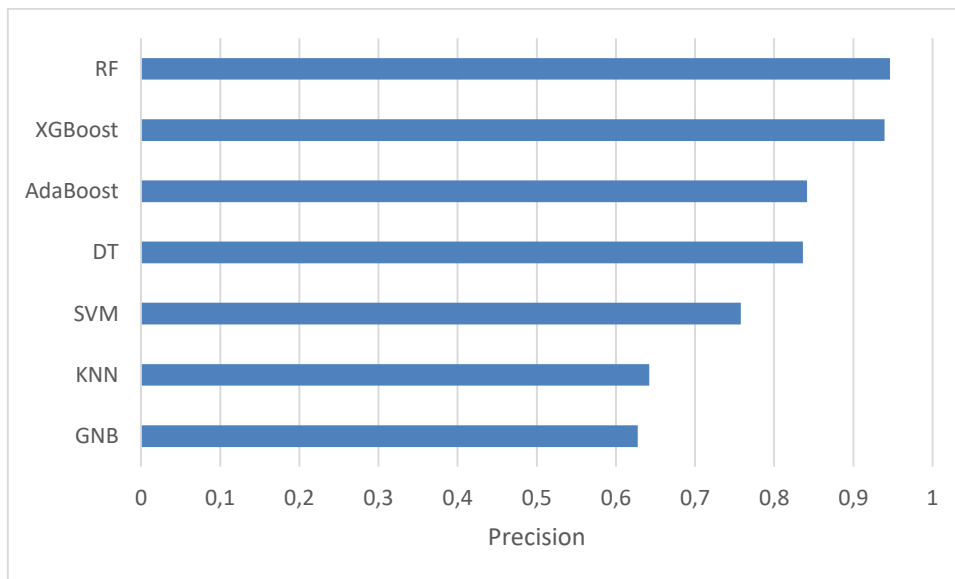
#### 4. Results

In this part of the study, the results obtained from the ML algorithms applied to the AD prediction problem are evaluated in terms of different metrics and presented. The accuracy results obtained from seven different machine learning algorithms are presented and compared in Table 2.

**Table 2.** Comparison of ML methods in accuracy.

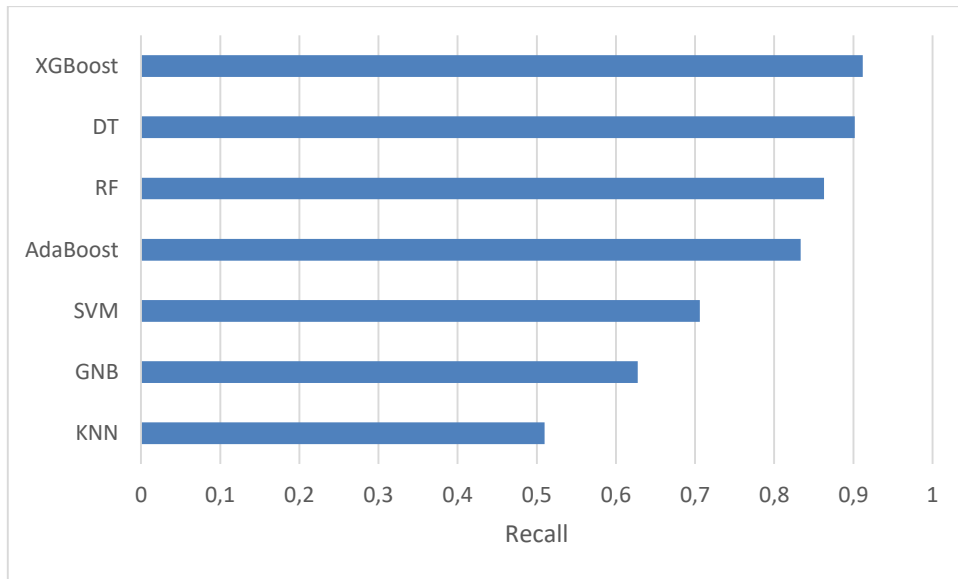
Model	Accuracy
RF	94.12
DT	91.33
SVM	83.59
KNN	75.54
AdaBoost	89.78
GNB	76.47
XGBoost	<b>95.35</b>

In Table 2, the accuracy rates of the machine learning algorithms applied in the study are compared. XGBoost stood out as the most successful model on this dataset, exhibiting the highest performance with an accuracy rate of 95.35%. While the RF algorithm closely follows XGBoost with an accuracy rate of 94.12%, AdaBoost showed a moderate performance with 89.78%. While SVM showed a relatively lower success with an accuracy rate of 83.59%, KNN and GNB showed the lowest performances with accuracy rates of 75.54% and 76.47%, respectively. The results show that the ensemble methods XGBoost and RF are strong options for more complex models.



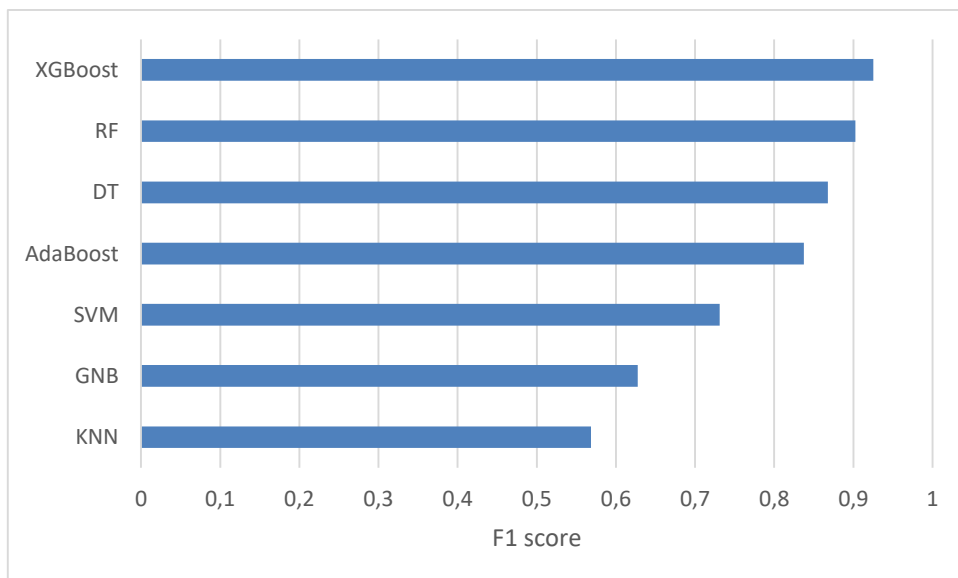
**Figure 2.** Comparison of ML methods in precision.

Figure 2 compares the precision values of different machine learning algorithms implemented in this study. Precision measures the success of a model in minimizing false positive results by expressing the ratio of true positive predictions to total positive predictions. When Figure 2 is examined, it is seen that XGBoost and RF algorithms reach high precision values and perform better than other models. This shows that these two algorithms can predict true positive results with high accuracy in Alzheimer's diagnosis. In contrast, it is understood that the precision values of KNN and GNB algorithms are relatively lower. AdaBoost and SVM provide balanced results in terms of precision, exhibiting a moderate performance. In general, these results show that XGBoost and RF are effective not only in accuracy but also in reducing false positive predictions. These results emphasize the importance of the precision metric during model selection, especially in areas such as medical diagnosis, where false positives are critical.



**Figure 3.** Comparison of ML methods in recall.

Figure 3 compares the recall values of different machine learning algorithms implemented in this study. Recall measures the rate at which a model correctly predicts all true positives and is important for minimizing false negatives (missed positives). As can be seen from the graph, XGBoost, DT, and RF algorithms show superior performance in capturing true positives with high recall values. This highlights the ability of these models to correctly identify patients in a critical diagnosis such as Alzheimer's. AdaBoost and SVM, which showed moderate performance, showed relatively good recall values, but fell behind the most successful models in this metric. KNN and GNB, on the other hand, have lower recall values, indicating that these models missed some of the positive cases. In general, high recall values are of great importance, especially in cases such as disease diagnosis, where missing positive classes is critical. These results show that XGBoost and RF are reliable models, showing a balanced performance not only in terms of accuracy and precision, but also in terms of recall.

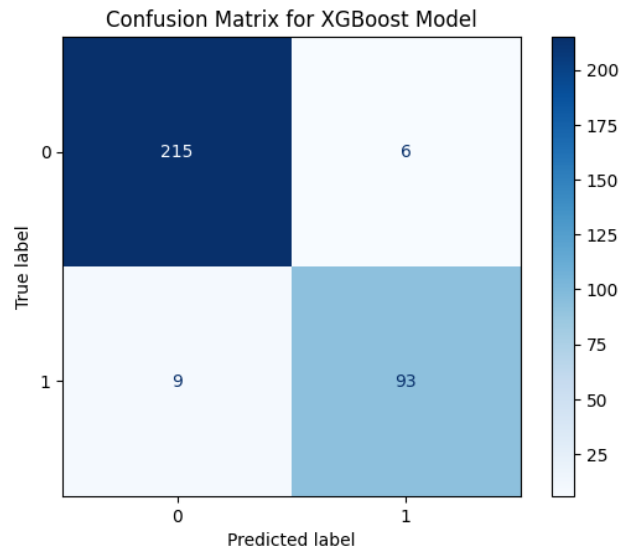


**Figure 4.** Comparison of ML methods in F1 score.

Figure 4 compares the F1-score values of different machine learning algorithms implemented in this study. Since F1-score is the harmonic mean of precision and recall metrics, it measures the success of a model in correctly predicting the positive class and its ability to deal with imbalanced classes in this process in a holistic way. According to the graphical results, XGBoost and RF achieved the highest F1-score values, which shows that both algorithms effectively detect positive classes while minimizing false positive predictions. It is clear that these models successfully maintain class balance and their overall generalization performance is superior. Although DT and AdaBoost performed reasonably in terms of F1-score, they fell behind the best performing models. On the

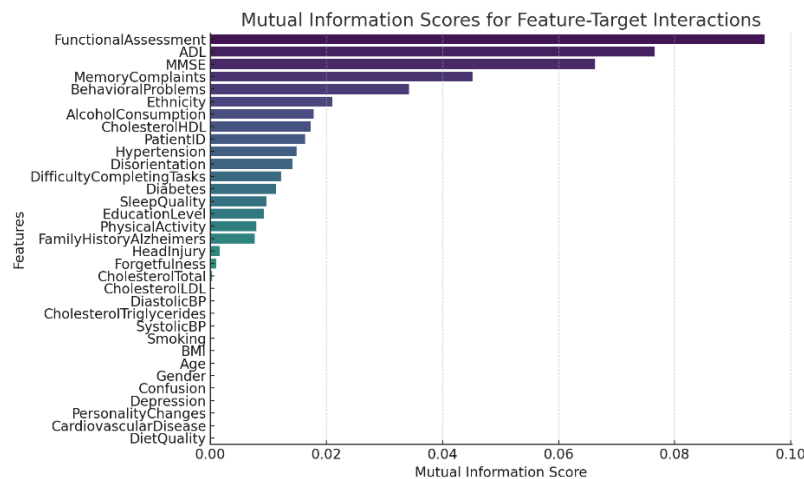


other hand, the low F1-score values of models such as KNN and GNB show that these algorithms cannot distinguish between imbalanced classes effectively enough. SVM achieved a moderate success and was relatively successful in finding a balance between precision and recall. This analysis reveals that F1-score is a critical metric and provides a comprehensive evaluation when evaluating model performance, especially on datasets where class imbalance is significant.



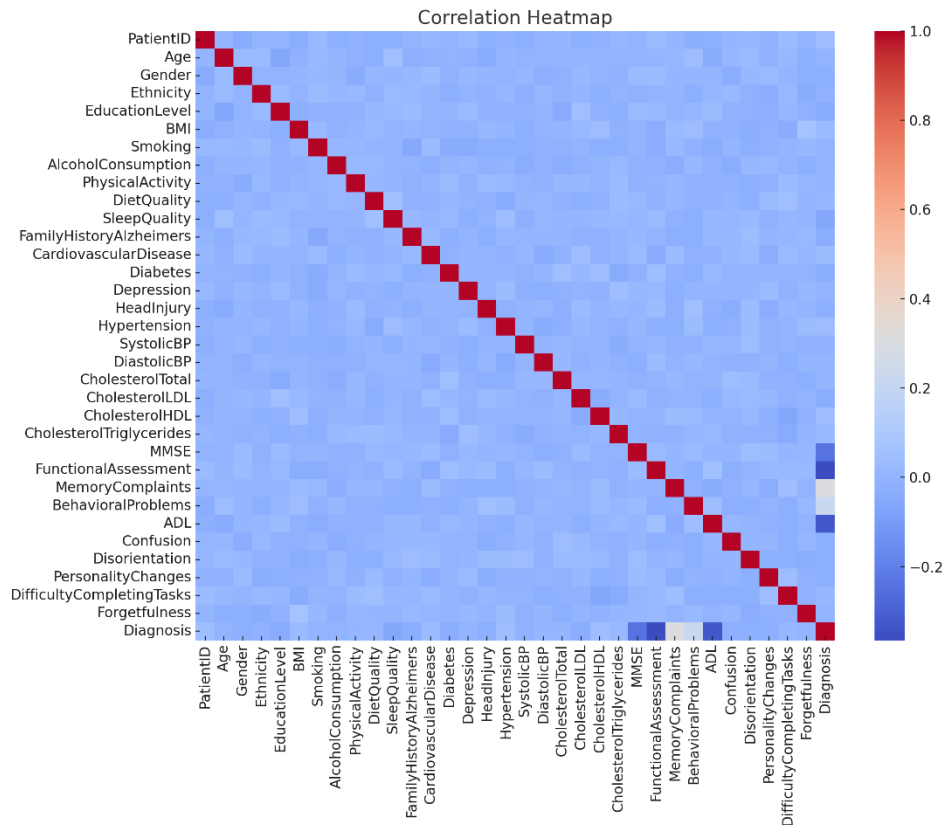
**Figure 5.** Confusion matrix for XGBoost model.

The confusion matrix in Figure 5 demonstrates the performance of the classification model made with the XGBoost algorithm. According to the matrix, the model predicted 215 true negatives and 93 true positives. However, there were 6 false positives and 9 false negatives among the incorrect predictions. This displays that the model is quite successful in correctly predicting the negative class, but makes relatively more errors for the positive class. When evaluated in general, it is clearly seen that the accuracy rate and prediction performance of the model are high.



**Figure 6.** Mutual Information Scores For Feature-Target Interactions.

The graph in Figure 6 evaluates the relationships between the target variable (Alzheimer's diagnosis) and the features using Mutual Information Score. As seen in Figure 6, FunctionalAssessment, ADL (Activities of Daily Living) and MMSE (Mini Mental State Assessment) scores are the features with the strongest relationship with the target variable. This situation reveals that cognitive and functional assessments are critical indicators for AD. Less effective features include factors such as CholesterolTotal, Smoking, BMI, and Age, which indicates that these variables provide limited information in diagnosing the disease. According to the Mutual Information Score results, cognitive status and daily life functions are among the most important factors for both risk assessment and early diagnosis.



**Figure 7.** Correlation Heatmap

The correlation heatmap in Figure 7 visualizes the relationships between numerical variables in the dataset. The graph represents the positive or negative correlation of each variable with other variables through color intensity. Specifically, red tones represent high positive correlation (close to 1), while blue tones represent negative correlation (close to -1). Overall, a heatmap is an important tool for understanding which variables are related to each other and for taking these relationships into account during modeling.

## 5. Conclusion

The findings obtained from this study indicate that machine learning algorithms have demonstrated significant success in predicting Alzheimer's Disease. The experiments were conducted using a dataset that is publicly available on the Kaggle platform, which includes comprehensive demographic and clinical features relevant to Alzheimer's diagnosis. This dataset consists of 2,149 patient records and 35 features, with no missing values, ensuring the reliability and completeness of the data for the analysis. By utilizing this dataset, the study aimed to assess the predictive performance of widely recognized machine learning algorithms frequently employed in disease prediction and classification studies in the literature.

The machine learning algorithms evaluated in this study include XGBoost, KNN, GNB, DT, RF, AdaBoost, and SVM. Each algorithm was assessed using four essential performance metrics: accuracy, precision, recall, and F1-score. The XGBoost algorithm achieved the highest accuracy, with a rate of 95.35%, demonstrating its superior ability to predict Alzheimer's Disease compared to other methods. This finding aligns with previous research, where ensemble learning methods like XGBoost have shown notable effectiveness in handling complex and multidimensional data [8, 9, 21]. In contrast, the lowest accuracy was observed in the KNN algorithm, with a rate of 75.54%. Although KNN performed relatively less effectively than the other models, it is notable that all algorithms achieved accuracy rates above a certain threshold, highlighting their overall competence in predicting Alzheimer's Disease.

The results of this study emphasize the potential of ML algorithms as powerful tools for the early diagnosis and prediction of AD. The superior performance of XGBoost, in particular, emphasizes the importance of using advanced ensemble methods to capture complex patterns and interactions in clinical data. Furthermore, the inclusion of diverse demographic and clinical features in the dataset highlights the significance of integrating multi-dimensional data in predictive modeling. These findings suggest that machine learning-driven predictions

can play a pivotal role in facilitating early detection and guiding timely intervention strategies, ultimately contributing to improved patient outcomes. The study not only confirms the viability of machine learning algorithms in medical diagnosis but also provides valuable insights into their application in addressing critical challenges in healthcare.

## References

- [1] D.M. Khan, N. Yahya, N. Kamel, I. Faye, "Automated diagnosis of major depressive disorder using brain effective connectivity and 3D convolutional neural network," *IEEE Access*, 9, pp. 8835-8846, 2021, 10.1109/ACCESS.2021.3049427
- [2] M. Sudharsan and G. Thailambal. "Alzheimer's disease prediction using machine learning techniques and principal component analysis (PCA)," *Materials Today: Proceedings* 81, pp. 182-190, 2023.
- [3] A. Association, "2019 Alzheimer's disease facts and figures", *Alzheimer's & Dementia*, 15 (3), pp. 321-387, 2019.
- [4] C.K. Gomathy and A. Rohith Naidu, "The Prediction Of Disease Using Machine Learning Techniques", *International Journal of Scientific Research in Engineering and Management (IJSREM)*, vol. 5, no. 7, 2021.
- [5] C. R. Mallela, L. R. Bhavani and B. Ankayarkanni, *Disease Prediction Using Machine Learning Techniques*, IEEE, pp. 962-966, 2021.
- [6] T.V. Sriram, M.V. Rao, G.S. Narayana, D.S.V.G. Kaladhar and T.P.R. Vital, "Intelligent Parkinson Disease Prediction Using Machine Learning Algorithms", *International Journal of Engineering and Innovative Technology (IJEIT)*, vol. 3, no. 3, September 2013.
- [7] K. Gomathi and D. Shanmuga Priya, *Multi Disease Prediction using Data Mining Techniques*, 2016.
- [8] A. Gavhane, G. Kokkula, I. Pandya and K. Devatkar, *Prediction of Heart Disease Using Machine Learning Algorithms*, 2018.
- [9] S. Arunachalam, "Cardiovascular Disease Prediction Model using Machine Learning Algorithms", *International Journal for Research in Applied Science & Engineering Technology*, vol. 8, no. VI, June 2020, ISSN 2321-9653.
- [10] A.D., Praveen, T.P., Vital, D., Jayaram and L.V. Satyanarayana, "Intelligent Liver Disease Prediction (ILDLP) System Using Machine Learning Models". *Intelligent Computing in Control and Communication. Lecture Notes in Electrical Engineering*, vol 702, 2021. Springer, Singapore. [https://doi.org/10.1007/978-981-15-8439-8\\_50](https://doi.org/10.1007/978-981-15-8439-8_50).
- [11] R.E. Kharoua, "Alzheimer's Disease Dataset", 2024. Kaggle. <https://doi.org/10.34740/KAGGLE/DSV/8668279>.
- [12] K. Taunk, S. De, S. Verma and A. Swetapadma, "A Brief Review of Nearest Neighbor Algorithm for Learning and Classification," 2019 International Conference on Intelligent Computing and Control Systems (ICCS), Madurai, India, 2019, pp. 1255-1260, doi: 10.1109/ICCS45141.2019.9065747.
- [13] I. Triguero, D. García-Gil, J. Maillo, J. Luengo, S. García and F. Herrera, "Transforming big data into smart data: An insight on the use of the k-nearest neighbors algorithm to obtain quality data", *Wiley Interdisciplinary Reviews: Data Mining and Knowledge Discovery*, 9(2), e1289, 2019.
- [14] H. Kamel, A. Dhahir and J.M. Al-Tuwajjari. "Cancer classification using gaussian naive bayes algorithm." 2019 international engineering conference (IEC). IEEE, 2019.
- [15] S. Suthaharan and S. Shan, "Support vector machine." *Machine learning models and algorithms for big data classification: thinking with examples for effective learning*, 207-235, 2016.
- [16] K. Nong, H. Zhang and Z. Liu, "Comparative Study of Different Machine Learning Models for Heat Transfer Performance Prediction of Evaporators in Modular Refrigerated Display Cabinets", *Energies*, 17, 6189, 2024. <https://doi.org/10.3390/en17236189>
- [17] A. T. Azar, H. I. Elshazly, A. E. Hassanien and A. M. Elkorany, "A random forest classifier for lymph diseases." *Computer methods and programs in biomedicine*, 113(2), 465-473, 2014.
- [18] W. Wang and S. Dongchu, "The improved AdaBoost algorithms for imbalanced data classification." *Information Sciences* 563, 358-374, 2021.
- [19] S. Li and Z. Xiaojing, "Research on orthopedic auxiliary classification and prediction model based on XGBoost algorithm." *Neural Computing and Applications* 32.7 ,1971-1979, 2020.
- [20] P. Iacobescu, V. Marina, C. Anghel, and A-D. Anghel, "Evaluating Binary Classifiers for Cardiovascular Disease Prediction: Enhancing Early Diagnostic Capabilities. *Journal of Cardiovascular Development and Disease*", 11(12):396, 2024. <https://doi.org/10.3390/jcdd11120396>.
- [21] P. Pranjali, S. Mallick, A. Das, A. Negi and M.R. Panda, "Alzheimer's Disease Prediction Using Modern Machine Learning Techniques." 2024 International Conference on Advancements in Smart, Secure and Intelligent Computing (ASSIC). IEEE, 2024.