

A Major Project Report
On
**PREDICTION OF DRUG-DRUG INTERACTION USING
MACHINE LEARNING**

*Submitted in the Partial Fulfillment of the Requirements
For the Award of the Degree of*

Bachelor of Technology
In
CSE(Artificial Intelligence and Machine Learning)

By
G. Laxmi Pranathi [21211A6618]
G. Mukesh Nayak [21211A6619]
T. Shashank [21211A6657]

Under the Guidance of
Ms. Srilakshmi V **Assistant Professor**
Mr. A B Ramesh **Assistant Professor**



DEPARTMENT OF CSE(ARTIFICIAL INTELLIGENCE AND MACHINE LEARNING)

B V RAJU INSTITUTE OF TECHNOLOGY
(UGC Autonomous, Accredited by NBA & NAAC)
Vishnupur, Narsapur, Medak, Telangana State, India – 502 313
2024-2025

B V RAJU INSTITUTE OF TECHNOLOGY
(UGC Autonomous, Accredited by NBA & NAAC)
Vishnupur, Narsapur, Medak, Telangana State, India – 502 313

DEPARTMENT OF CSE(ARTIFICIAL INTELLIGENCE AND MACHINE LEARNING)

CERTIFICATE

This is to Certify that the Major Project Entitled “**Prediction of Drug- Drug Interaction using Machine Learning**” Being Submitted By

G. Laxmi Pranathi	[21211A6618]
G. Mukesh Nayak	[21211A6619]
T. Shashank	[21211A6657]

In Partial Fulfillment of the Requirements for the Award of Degree of Bachelor of Technology in CSE(Artificial Intelligence and Machine Learning) to B V Raju Institute of Technology is Record of Bonafide Work Carried Out During the Period From December 2024 to April 2025 by Them Under the Supervision of

Ms. Srilakshmi V

Assistant Professor

Mr. A B Ramesh

Assistant Professor

This is to Certify that the Above Statement Made by the Students are Correct to the Best of Our Knowledge.

Ms. Srilakshmi V
Assistant Professor

Mr. A B Ramesh
Assistant Professor

The Major Project Viva-Voce For This Team Has Been Held on _____.

External Examiner

Dr. G Uday Kiran
Program Coordinator

CANDIDATE’S DECLARATION

We Hereby Certify that the Work Which is Being Presented in the Major Project Entitled “Prediction of Drug-Drug Interaction using machine Learning” in Partial Fulfillment of the Requirements For the Award of Degree of Bachelor of Technology and Submitted in the Department of CSE(Artificial Intelligence and Machine Learning), B V Raju Institute of Technology, Narsapur, is an Authentic Record of Our Own Work Carried Out During the Period From December 2024 to April 2025, Under the Supervision of Name of Mr. A B Ramesh , Assistant Professor and Ms. Srilakshmi V, Assistant Professor.

The Work Presented in this Major Project Report Has Not Been Submitted By Us For the Award of Any Other Degree of This or Any Other Institute/University.

G. Laxmi Pranathi	[21211A6618]
G. Mukesh Nayak	[21211A6619]
T. Shashank	[21211A6657]

ACKNOWLEDGEMENT

We stand at the culmination of a significant journey, one that has been both challenging and rewarding. The success of our major project is not solely a reflection of our efforts but a testament to the invaluable support and guidance we have received from many quarters. It is with deep gratitude that we acknowledge those who have made this achievement possible.

Foremost, we extend our sincerest appreciation to Mr. A B Ramesh, Co-supervisor and Ms. Srilakshmi V, Supervisor, whose expertise and insightful supervision have been pivotal in navigating the complexities of this project. Their unwavering support and encouragement have been our guiding light throughout this journey.

Special thanks are due to Ms. Srilakshmi V, our Project Coordinator, whose assistance and guidance have been instrumental in the successful execution of our project. Her dedication and support have been a source of inspiration and motivation.

We reserve our utmost gratitude for Dr. G Uday Kiran, Program Coordinator of the Department of CSE (Artificial Intelligence and Machine Learning), whose leadership and academic guidance have enriched our learning experience and contributed significantly to our project's success. Our journey would not have been the same without the constant encouragement, support, and guidance from the esteemed faculty of the Department of CSE (Artificial Intelligence and Machine Learning). We are deeply thankful to everyone who contributed to our journey, whose belief, guidance, and support have been crucial to our achievement. This project reflects not only our academic efforts but also the collaborative spirit and collective wisdom that guided us.

G. Laxmi Pranathi	[21211A6618]
G. Mukesh Nayak	[21211A6619]
T. Shashank	[21211A6657]

ABSTRACT

The growing use of multiple medications has raised concerns about Drug- Drug Interactions, which can result in adverse drug reactions, increased healthcare costs, and patient safety risks. Traditional manual evaluation and clinical trials for DDI detection require a significant amount of time and resources. To address this, machine learning (ML) models have emerged as powerful tools for automated DDI prediction, providing a scalable and data- driven approach.

This study investigates structured feature extraction using Count Vectorization rather than NLP-based methods like BiLSTM and BioBERT, to ensure computational efficiency and interpretability. Multiple ML classifiers, such as SGD Classifier, XGBoost, LightGBM, CatBoost, Logistic Regression, and Random Forest, are tested for DDI classification, with Instance Hardness Threshold (IHT) undersampling used to address class imbalance and improve recall while maintaining high precision.

The results show that SGD Classifier has the highest F1-score (98.68%), outperforming deep learning-based approaches in previous research. These findings emphasize the potential of structured ML techniques in pharmacovigilance, allowing integration into clinical decision support systems (CDSS) for real-time drug safety assessments. Future research will concentrate on hybrid feature extraction and deep learning integration to improve real-world applicability.

TABLE OF CONTENTS

CHAPTER	TITLE	PAGE NO.
	Certificate	ii
	Candidate's Declaration	iii
	Acknowledgement	iv
	Abstract	v
	Table of Contents	vi
	List of Figures	ix
	List of Tables	x
	List of Abbreviations	xi
1	INTRODUCTION	1 - 7
	1.1 Introduction to The Domain	2
	1.2 Problem Statement	4
	1.3 Motivation	4
	1.4 Objective	6
	1.5 Organization of Work	7
2	LITERATURE REVIEW	8 - 25
	2.1 Introduction	8
	2.2 Related Work	9
	2.3 Research Gaps	23
3	ALGORITHM(S) ANALYSIS AND DESIGN	26 - 41
	3.1 Introduction	26
	3.2 System Requirements	27
	3.2.1 Problem Definition: Need for Automated Drug Interaction Prediction	27
	3.2.2 Input Data	28
	3.2.3 Output Data	29
	3.2.4 Performance Goals	30
	3.2.5 Scalability	31

CHAPTER	TITLE	PAGE NO.
	4.3.1.1 IHT's effects on the dataset	52
	4.3.2 Oversampling Using Synthetic Minority Over sampling Technique (SMOTE)	53
	4.3.2.1 The effect that SMOTE has on the dataset	54
4.4	Data Splitting	54
	4.4.1 Examined Techniques for Dataset Splitting	56
4.5	Machine Learning Models	56
	4.5.1 Implemented Models and their Role	56
	4.5.1.1 SGD Classifier	56
	4.5.1.2 Extreme Gradient Boosting (XGBoost)	58
	4.5.1.3 LightGBM	60
	4.5.1.4 CatBoost	61
	4.5.2 Benchmark Models for Comparison	63
	4.5.3 Role of These Models in the Project	65
4.6	Evaluation Metrics	67
	4.6.1 Precision	68
	4.6.2 Recall	68
	4.6.3 F1-Score	69
4.7	Results	70
	4.7.1 Performance Comparison of Machine Learning Models	71
	4.7.2 Key Observations from Model Performance	72
	4.7.3 Confusion Matrix Analysis	73
	4.7.3.1 Key Insights from the Confusion Matrices	74

CHAPTER	TITLE	PAGE NO.
5	CONCLUSION AND FUTURE SCOPE	76 – 80
	5.1 Conclusion	78
	5.2 Future Scope	79
	REFERENCES	81 - 87

LIST OF FIGURES

Figure	Description	Page No.
Fig 4.1	Initial count of each level	44
Fig 4.2	Class Distribution after applying IHT	51
Fig 4.3	Confusion Matrix of SGD Classifier	74

LIST OF TABLES

Table	Description	Page No.
Table 3.1	Table Time Complexity Analysis	39
Table 4.1	Distribution of Classes Prior to Eliminating Duplicates	48
Table 4.2	Distribution of Classes Following Duplicate Removal	48
Table 4.3	Results using Oversampling technique(SMOTE)	71
Table 4.4	Results using Undersampling technique(IHT)	72

LIST OF ABBREVIATIONS

Abbreviation	Full Form
ADR	Adverse Drug Reactions
AI	Artificial Intelligence
AMDE	Attention Based Multidimensional Feature Encoder
DDI	Drug-Drug Interaction
DILI	Drug-Induced Liver Injury
DNER	Drug Named Entity Recognition
DTI	Drug-Target Prediction
EHR	Electronic Health Records
GMPNN	Gated Message Passing Neural Network
GNNs	Graph Neural Networks
IHT	Instance Hardness Threshold
ML	Machine Learning
MDNN	Multimodal Deep Neural Network
NLP	Natural Language Processing
SMOTE	Synthetic Minority Over-Sampling Technique
XAI	Explainable Ai

CHAPTER 1

INTRODUCTION

Drug interactions refer to letting two or more pharmacologic agents act together, either synergistically, antagonistically, or inertly, to produce unanticipated and harmful effects. This can present itself in ways like reduced effectiveness, increased toxicity, adverse reactions, or even death. Drug interactions occur due to the combination of pharmacokinetic and pharmacodynamic mechanisms. Pharmacokinetic interactions refer to changes in the absorption, distribution, metabolism, or excretion of one or more drugs. On the other hand, pharmacodynamic interactions involve the synergistic or opposing effects of agents on the same physiological pathway, which can have either therapeutic advantages or unwanted health effects.

In the past, detection of drug interaction has been based on a combination of clinical trials, laboratory experiments, and then post-marketing surveillance. However, these approaches are generally time-hogging, very expensive, and incomplete. Although these attributes characterize drug interaction finding, regulatory agencies and healthcare providers call upon voluminous databases for documentation purposes about existing interactions; that is, through information on drugs, published literature, and adverse event reporting systems.

Although the tradition of establishing drug interactions has always been based on a combination of clinical trials, laboratory experiments, or, more recently, post-marketing surveillance, such approaches tend to be time-consuming, very expensive, and incomplete. Thus, regulatory agencies and health care providers fall back to huge databases of drug information, published literature, and adverse event reporting systems for documenting the known interactions.

Yet the development of ever more complex pharmacopeias and expanding numbers of new medications has positioned major obstacles against activities intended to ascertain interaction through manual analysis. To protect patients' safety and optimize treatment efficacy, the discovery of potential drug interactions needs to implement effective, automated methods for predicting and evaluating potential interactions.

Drug interactions occur when two or more medicines are absorbed together and modify the normal effects of each. Either beneficial, neutral, or harmful are all possible outcomes. Drug interaction reduces efficacy, adds side effects, and might cause major health risks, particularly in patients with more than one medical condition (polypharmacy), like the aged. Medical therapy nowadays is so complicated that knowing how to predict these interactions is a must for patient safety and possible therapeutic benefit.

1.1 INTRODUCTION TO DOMAIN

The project includes artificial intelligence, especially machine learning, data-driven healthcare analytics, and prediction. Today, machine learning algorithms are quickly coalescing with drug discovery and biomedical research, as well as healthcare decision-making, due to advancements in AI. An important application of AI is with the prediction of drug-drug interactions as it tries to mine very large datasets in discovering whether or not the substances harbor any interactions. Traditional methods require very expensive and laborious clinical trials and expert opinion. These are the ones where machine learning using deep learning models, huge biomedical datasets, natural language processing approaches provides scalability and automation to make drug safety better.

AI methods help to improve the prediction accuracy by integrating diverse data sources-molecular structure, electronic health records (EHRs), and biomedical literature. Transformer-based frameworks, and the likes of BioBERT and graph neural networks (GNNs), hold good promise to detect

complex drug relationships. In addition, AI models learn and adjust continuously with real-time predictions when new drug interactions are discovered. It alerts health professionals earlier in time and saves the procedure from tedious and prolonged results while improving the risk of adverse drug reactions (ADRs). These models can also be more interpretable and reliable while making clinical decisions when explainable AI (XAI) techniques are integrated with them. At last, all these augment, AI-enabled DDI prediction to precision medicine through ensuring patient safety and treatment effectiveness.

Different data sources, such as molecular structures, electronic health records (EHRs), and more, are made manageable through AI-driven methods to improve accurate predictive values. Transformer-based models, such as BioBERT, and, most recently, graph neural networks have shown a lot of promise in capturing complex drug relationships. These models also have real-time predictions when new drug interactions are found, as the AI models are always learning and adjusting. It not only saves time with drug approval processes but also reduces the incidence of adverse drug reactions (ADRs) by warning medical practitioners earlier. When integrating XAI techniques with such models, the models also gain more interpretability and reliability when making clinically valid healthcare decisions. All of these thus augment the advancement of precision medicine through AI-enabled DDI predictions in patient safety and treatment enhancements.

The use of artificial intelligence in drug interaction research is a change from a reactive healthcare model to a proactive one. Historically, most such types of data required extensive manual curation and relied on postmarket surveillance, which often resulted in missing rare or novel interactions due to sample size constraints and bias. Machine learning can process multimodal data relating to everything from chemical properties and protein-binding affinities to actual patient records and discover latent patterns that conventional analysis cannot readily uncover. The application of AI in drug interaction research ushers in a paradigm shift from reactive

to proactive healthcare. Most conventional setups used manual curation and even post-marketing surveillance, thus exhibiting a common mode of failure in rare or novel interaction incidences due to limited sample sizes and human bias. However, machine learning enables the processing of multimodal data-from chemical properties and protein-binding affinities to real-world-patient records-to uncover hidden patterns, which conventional analysis does not readily reveal.

1.2 PROBLEM STATEMENT

DDIs pose a serious risk in medicine due to the potent side effects that may lead to treatment failure and ADRs. Conventional approaches to DDI identification that require full-scale clinical trials and expert review have proved expensive, time-consuming, and unscalable. We will apply various machine learning techniques on structured datasets in this project, automating and augmenting the prediction of probable DDIs. This enhancement pertains to risk reduction, security of the patients, and early recognition of hazardous drug interactions. This further improves clinical decision-making, accelerates the drug approval process, and reduces dependence on manual evaluations.

1.3 MOTIVATION

The growing concern towards recent health advancement and pharmaceuticals draws attention to drug safety. Drug-drug interactions or the DDIs remain the toughest challenges to identify for an unexpected side effect, treatment failure, or adverse drug reaction.

Traditional DDI detection methods rely on clinical trials, expert evaluations, and post-market surveillance, which can be costly and inefficient to cover a great number of potential drug combinations. New avenues to improve drug safety are based on the potential of AI and ML to automate the prediction of possible DDIs. Machine learning has the ability to analyze large-scale biomedical datasets and find patterns from drug interactions which will be otherwise overlooked by a manual evaluation.

AI offers critical scalability and efficiency for DDI prediction. As the number of approved drugs keeps rising, an exponentially increasing number of possible interactions is noted, making it impractical to evaluate all the different combinations manually, making AI a necessity for a sound haven in which to work with such enormous data sets. Machine learning models may provide instant predictions so that medical practitioners, whether involved in dcs or not, will swiftly make good decisions. Processing huge amounts of real-time prediction data is very much a reflection of how machine learning is deployed in hospitals. Moreover, because of advances from AI and ML, such as deep learning, graph-based learning, and natural language processing, the ability of models to derive valuable insights from biomedical text or molecular structures, and patient records has improved immensely. And that allows for better and more precise predictions compared to traditional statistical approaches.

In addition to these, it is another driving engine for this project, as one urgent need is to reduce adverse drug reactions (ADRs), which account for massive hospitalizations and deaths worldwide. By diagnosing dangerous drug interactions early, they can make better prescribing decisions which could prevent many of the adverse drug reactions through doctors and pharmacists. Through deploying AI-driven models for automated DDI prediction, healthcare organizations would minimize risks, enhance patient safety, and lessen their reliance on expensive and time-consuming clinical trials. The impetus for moving towards AI-powered DDI predictions is not only limited to clinical and computational efficiency. Furthermore, this approach is consistent with broader public health goals. By helping to allay fears of overly individualistic and personalized medication treatments as enhancing individualized genetic treatments and addressing comorbidities, the risks of unique or rare drug interactions should be reduced as well as conditions by providing easy and cheap scalable tools. It also addresses increasing polypharmacy among aging people.

1.4 OBJECTIVE

To enhance patient safety, the early detection of dangerous DDI's, and informed decision-making among health professionals, the aim of this project is to build a robust Drug-Drug Interaction (DDI) prediction model using various machine-learning techniques. To this end, the project will focus on analyzing and contrasting a number of machine learning models such as SGD Classifier, XGBoost, LightGBM, CatBoost, Logistic Regression, and Random Forest to find the best classification approach. The advanced techniques for preprocessing the data such as feature extraction using Count Vectorizer, dataset splitting strategies of 80-20, 70-30, K-Fold Cross-Validation, and cleaning away redundant entries have all been applied to maintain the highest quality of input data. Class balancing methods become necessary because DDI datasets tend to be highly unbalanced, with some interaction severity levels being underrepresented. The techniques applied thus ensure that the models are trained on balanced data for better generalization involving Synthetic Minority Over-Sampling Technique (SMOTE) for oversampling and Instance Hardness Threshold (IHT) for undersampling.

The model performance can be evaluated using essential performance indicators such as precision, recall, and F1-score. The application of such methods is incorporated in the project to develop a reliable, durable, and highly accurate AI model for the automated classification of drug-drug interactions (DDI), ultimately resulting in safer medication usage and improved healthcare decisions. The project systematically evaluates various machine learning approaches to optimize a framework for drug-drug interaction prediction. The comparative analysis of varying algorithms, including ensemble methods and regularized linear models, aims to identify an optimal solution for this demanding health-care problem. Special emphasis is paid to overcoming some inherent limitations posed by the dataset through extensive preprocessing methods and mitigation against class imbalance. Therefore, robust validation strategies and clinically

relevant performance metrics evaluate the proposed approach for developing a trustworthy prediction system.

This initiative aims to build a realistic instrument that joins state-of-the-art machine-learning functions with the practical requirements of the clinic. This project will, therefore, be delivered as a solution with the necessary predictive performance and operational viability to support effective implementation in the healthcare workflow. Ultimately, the objective is to reasonably and reliably equip decision support to professionals in health services to improve patient care and lessen adverse drug events.

1.5 ORGANIZATION OF WORK

Drug-drug interactions (DDIs) have been introduced in Chapter 1 of the report, which presents the clinical context to consider DDIs, problem statements, and objectives. Chapter 2 contains an extensive literature review and points out gaps in the pronouncements surrounding the aforementioned NLP applications for DDI prediction and the utter need for rigor in the study methodology. In Chapter 3, the methodology indicates data preparation (Count Vectorization, IHT under-sampling), followed by algorithm selection (SGD Classifier and XGBoost) and complexity analysis. Then Chapter 4 gives experimental results in a comparative model performance manner; the SGD Classifier provided an F1-score of 98.68% with important takeaways from the confusion matrices. Chapter 5 concludes the entire study by summarizing the study contributions, providing hints for the future, such as hybrid feature extraction and CDSS integration.

CHAPTER 2

REVIEW OF LITERATURE

2.1 INTRODUCTION

Machine learning for drug-drug interaction (DDI) study has gained considerable attention recently, primarily due to polypharmacy and the increasing number of approved drugs. The classical methods for DDI extraction rely on clinical trials, manual curation from biomedical literature, and expert knowledge. Such methods tend to be costly and time consuming and mostly cannot keep up with the fast-track life cycles of new drugs. Researchers have explored machine learning-based methods to automate and develop DDI prediction processes, aiming at overcoming current limitations.

Numerous methods, such as deep learning architectures and feature-based machine learning models, have been used in previous research. To improve the accuracy of DDI prediction, several methods make use of pharmacokinetic characteristics, protein-protein interaction networks, and representations based on molecular structures. Furthermore, some studies use Natural Language Processing (NLP) and text mining to identify patterns of interaction in extensive biomedical texts. However, our project avoids the complexity of NLP-based methods by concentrating on structured data analysis and employing Count Vectorization for feature extraction.

Furthermore, since real-world datasets frequently contain a significantly higher number of non-interacting drug pairs than interacting ones, managing class imbalance continues to be a significant challenge in DDI prediction. In order to enhance model performance, a number of studies have investigated data balancing strategies, such as Random Undersampling, Edited Nearest Neighbors (ENN), and Synthetic Minority Over-Sampling Technique (SMOTE). By applying SMOTE for oversampling and Instance Hardness Threshold (IHT) for undersampling, our project

expands on these discoveries and guarantees a more balanced training dataset.

This literature review aims to establish the best machine learning algorithms, feature extraction methods, and data balancing techniques for DDI classification. Therefore, to develop a more precise, scalable, and interpretable machine learning model for DDI prediction, our work studies earlier work and current trends. With the recent advancement of machine learning, DDI prediction can be approached with increased sophistication, especially with the incorporation of multimodal data sources. Recent advancements show that drug properties and combinations, alongside biological pathway information and clinical outcome data, improve their predictive capacity. Generally, the field has moved away from classical manual feature engineering toward increasingly automated forms of representation learning, with graph neural networks in particular emerging as a very natural and powerful tool to effectively represent complex relationships of drug-drug interactions.

2.2 RELATED WORK

Zhang Yang et al. [1] Predicting Drug-Drug Interactions Using Multi-Modal Deep Auto-Encoders Based Network Embedding and Positive-Unlabeled Learning. A multi-modal deep autoencoders-based method referred to as DDI-MDAE has been proposed in this article to predict drug-drug interactions from multimodal data. The article demonstrates that the methodology named positive-unlabeled learning makes it possible for DDI-MDAE to outperform the original DDI-MDAE method.

Xuan Lin et al. [2] KGNN is a knowledge graph neural network for predicting the interaction between drugs. The research proposes a new end-to-end model named Knowledge Graph Neural Network (KGNN) which adequately perceives and processes high-order structures and semantic relationships within the knowledge graph for drug-drug interaction prediction purposes.

Abbasi Karim et al. [3] "Rather recently, there has been another attempt, one concentrated upon learning on deep structures for drug/target-interaction predictions: a thorough examination of approaches in conjunction with an analysis of the deep network architectures involved; advantages and disadvantages of each architecture; integration of drug/protein feature-building for the purpose of consideration and used in predication; commonly used datasets; current issues and the outlook in the future in this area."

Sukhvinder Singh et al. [4] The article "Artificial Intelligence in Drug Discovery" elaborates that, indeed, for revolutionization of drug development processes, the expensive and time-consuming protocols can be ousted due to the advent of Artificial Intelligence.

Schwarz et al. [5] AttentionDDI: A dual attention based deep learning method for predicting drug-drug interactions (DDI). This article describes a brand novel self-attention multi-modality neural network architecture, termed attentionDDI, for both predicting drug interactions and comparing it against state-of-the-art DDI prediction models using several benchmark datasets.

D'souza Sofia et al. [6] This article reviews the use of machine learning and deep learning methods in predicting drug-target interactions, which are important for drug discovery; current machine learning in drug-target interaction prediction: state and future directions.

Lee Chun Yen et al. [7] "Prediction of Drug Adverse Events Using Deep Learning in Pharmaceutical Discovery," Briefings in Bioinformatics (2020). The paper states that Deep learning techniques are capable of efficiently detecting and classifying drug side effects as compared with traditional machine learning methods. Deep learning methods may also assist in drug replacement or repositioning of drugs with adverse effects for alternative disease treatment.

Deng Yifan et al. [8] This paper proposes a multimodal deep learning framework which shall be known as DDIMDL. It includes the deep learning model integrated with multifarious drug features for predicting events occurring as a result of drug-drug interactions. A deep-learning model is built to predict the adverse drug interaction events by DDIMDL. The framework further encompasses the features of different drugs into the model by deep learning so that the input data are pooled across sources.

Kumar Shukla et al. [9] "Deep learning models for efficient prediction of drug-drug interaction." An integrated convolutional mixture density recurrent neural network was proposed and implemented in this paper. The proposed model integrated convolutional neural networks with recurrent neural networks and mixture density networks.

Romm Eden L et al. [10] The topic of the writeup is "Artificial Intelligence in Drug Therapy." It introduces recent and probable future applications of AI, such as the discovery of drugs, drug development, and drug prescription, in the pharmaceutical domain.

Feng Zhang et al. [11] In this paper, they created a model for deep learning to predict drug-drug interactions-the DPDDI-which uses a graph convolution network (GCN) to learn low-dimensional feature representations for each drug within DDI networks while using deep neural networks (DNNs) for model training.

G Smith et al. [12] Artificial Intelligence applied to drug safety and metabolism. Artificial intelligence can be used in every step of drug discovery and development from profiling a chemical library in early development for target identification, to predicting off-target effects in mid-discovery, to assessing possible mutagenic impurities in development as well as degradants associated with the life cycle management.

Pang Shanchen et al. [13] AMDE: a novel attention-mechanism-based multidimensional feature encoder for drug-drug interaction prediction." In this work, they proposed a novel attention-mechanism-based

multidimensional feature encoder for DDIs prediction, called the attention-based multidimensional feature encoder (AMDE).

Ashwin Dhakal et al. [14] Prediction against protein-ligand interaction in AI: current trends and future directions, Recently, the work has reviewed computational approaches predominantly AI and machine learning-based ones that predict the interaction of proteins with ligands in drug discovery.

Rewording for low perplexity and high burstiness level above: In the field of artificial intelligence, prediction against protein-ligand interaction: current trends and future directions, Recent work summarizes the computational approaches that have mainly used AI and machine learning for predicting interactions between proteins and ligands in drug discovery.

Lyu Tengfei et al. [15] MDNN: A Multimodal Deep Neural Network for Prediction of Drug-Drug Interaction Events. In this communication, a fresh model-State Multimodal Deep Neural Network (MDNN)-is introduced for predicting drug-drug interactions that effectively utilizes topological and semantic relations from a drug knowledge graph as well as cross-modal complementarity of multimodal drug data.

Kim J et al. [16] There is a survey that reviews the recent deep learning applications in drug-target interaction (DTI) prediction and in de novo drug design. Furthermore, a very complete summary is provided regarding a variety of drug and protein representations, DL models, and commonly used benchmark datasets or tools for model training and testing.

Ibrahim Heba et al. [17] The authors propose that a machine-learning framework, SMDIP, which relied on similarities, was developed to predict safety signals of new and unknown drug-drug interactions with high precision and shown to present good predictive performance that was further validated on an unseen set of hepatitis C drugs.

Chen Yujie et al. [18] "The Muffin: multi-scale feature fusion for drug-drug interaction prediction". MUFFIN is a multi-scale feature fusion deep learning

model that is capable of learning drug representations from both drug molecular structures and knowledge graph information for predicting their interaction beyond limited available data.

Nyamabo Arnold K et al. [19] Drug-drug interaction prediction with learnable size adaptive molecular substructures." They propose the gated message passing neural network (GMPNN), a message passing neural network that learns chemical substructures differing in size and shape from molecular graph representations of drugs for DDI prediction of a drug-drug pair.

Arnold K et al. [20] An introduction to the substructure-substructure interaction-drug-drug interaction (SSI-DDI), a framework for deep learning, that operates directly on the raw molecular graph representations of drugs for richer feature extraction but then basically breaks the DDI prediction task between two drugs down to pairwise interactions between their constituent substructures.

Luo Xue Y et al. [21] A unique deep learning model is developed and validated for DDI prediction through the 89,970 known DDIs extracted from DrugBank. The proposed deep learning model comprises a graph convolutional autoencoder network (GCAN) and long short-term memory (LSTM) that greatly enhance the performance of drug-drug interaction prediction.

Ovek Damla et al. [22] "The predictive techniques utilizing artificial intelligence for hot zones." This article provides an overview concerning the ways of employing artificial intelligence and machine learning techniques for computational identification of so-called hot spots, i.e., amino acid residues embedded in protein-protein interfaces, which contribute most to interaction affinity, and, finally, addresses the importance of hot spot prediction for drug design.

Huang Wei et al. [23] The main content of this paper highlights the review of computation methods, especially deep learning models, used for

predicting DDIs-an important issue because of the increasing number of new drugs and the high cost and time component of the experimental drug assays.

Syrowatka Ania et al. [24] "This scoping review identifies key artificial intelligence use cases to mitigate the frequency of drug-event adverse effects"; this paper touches on a scoping review that identifies the key use cases for artificial intelligence to reduce the frequency of adverse drug events, focusing on contemporary machine learning techniques and natural language processing.

Lin Shenggeng et al. [25] A new method is proposed, namely "MDF-SA-DDI", which predicts drug-drug interaction (DDI) events based on multi-source drug fusion, multi-source feature fusion, and transformer self-attention mechanism.". The new method MDF-SA-DDI149 could be used for predicting DDI events based on multi-source drug fusion, multi-source feature fusion, and transformer self-attention mechanism.

Avinash et al. [26] From recent assessment of artificial intelligence and their in silico ADMET prediction in the early stages of introduction of drugs. ADMET models in silico have grown a lot, and in the end, they fail the drug candidates in ADMET profiles and severe dropout incidences; advancements in machine learning, particularly deep neural networks, are proven aspects for better predictions of ADMET but require quality data and time to show their real value.

Vall Andreu et al. [27] "The Promise of AI for DILI Prediction" provides a comprehensive review of AI approaches for predicting drug-induced liver injury (DILI), with a focus on machine learning methods and a discussion of the primary datasets, data modalities, and techniques that have been presented in the literature.

Nayarisseri Anuraj et al. [28] "Artificial Intelligence, Big data and Machine Learning approaches in Precision Medicine & Drug Discovery." This paper covers the different machine learning techniques that can be applied in

various stages of drug discovery and development, including SNP discovery, drug repurposing, ligand-based drug design, virtual screening, lead identification, QSAR modeling, and ADMET analysis, among others, and discusses some successful case studies where these techniques were used.

Nyamabo Arnold K et al. [29] "Drug-drug interaction prediction with learnable size-adaptive molecular substructures." It is their GMPNN, a message-passing artificial neural network that learns chemical substructures with varying sizes and shapes from the molecular graph representations of drugs for the purpose of drug-drug interaction prediction between a pair of drugs.

Kumar et al. [30] Machine learning techniques for predicting drug-drug reactions - a critical review. This paper offers a thorough investigation into machine learning and deep learning techniques for the prediction of drug-drug interactions, elaborating on the advantages and disadvantages of different approaches.

Qiu Yang et al. [31] A thorough paper on computation methods used for drug-food interaction detection." This paper discusses the different methods of extracting drug-drug interactions including the three basic categories of literature based extraction, machine learning based prediction, and pharmacovigilance based data mining.

Liu Shichao et al. [32] This paper introduces a framework for drug-drug interaction prediction using deep attention neural networks (DANN-DDI), aimed at predicting unobserved drug-drug interactions.

Jizhou Tian et al. [33] Paper on methods for detecting drug-drug interactions. This paper discusses traditional approaches as well as advanced computational methodologies for predicting drug-drug interactions, which have become a kerfuffle over using drugs.

Joshi Pratik et al. [34] "The method of predicting the adverse drug reactions using deep neural networks based on knowledge graph embedding." The

novel approach propounded by the authors is based on the idea of knowledge graph embedding and a tailor-made deep neural network known as KGDNN that has been especially constructed to understand and predict the adverse drug reactions, which very crucially solves one of the major problems areas in drug discovery that has not been effectively addressed so far.

Bittner et al. [35] The paper provides an overview of the applications, opportunities, and challenges of artificial intelligence (AI) and machine learning (ML) in drug discovery, which depends on high-quality data and on the synergistic application of dry and wet lab methods.

Haohuai et al. [36] "3DGT-DDI: 3D graph and text based neural network for drug-drug interaction prediction." The study proposed a new deep learning model, 3DGT-DDI, for the prediction of drug-drug interactions from 3D molecular structure and text information, significantly outperforming other state-of-the-art methods on a benchmark task.

Purvashi et al. [37] "Machine Learning and Artificial Intelligence: A paradigm shift of Big Data Driven Drug Design and Discovery" - an assessment of various AI-assisted methodologies for drug development. The article reviews the advantages of different approaches, their applications, and some transformational case studies in drug discovery and development.

Feng Y.-H. et al. [38] It is about the prediction of drug-drug interaction utilizing an attention-based graph neural network on drug molecular graphs; they proposed a novel method of GNN-DDI for predicting potential DDIs by constructing a five-layer graph attention network to identify k-hops low-dimensional feature representations for each drug from its chemical molecular graph, concatenating all these identified features of each drug pair for input into an MLP predictor for obtaining the final DDI prediction score.

John B et al. [39] Using domain adaptation, the interpretable bilinear attention network enhances drug-target prediction. DrugBAN is an

interpretable deep learning framework modeling pairwise local interactions between drugs and targets and facilitating cross-domain generalization via conditional domain adversarial learning.

Pratik et al. [40] A method proposed was an embedding for knowledge graphs to predict adverse drug reactions through a deep network. Custom KGDNN, a deep neural network, was proposed by the authors which made a significant contribution to predicting adverse drug reactions. This study addressed the pharmacovigilance problem that has posed several challenges to the other approaches in the literature.

Yazdani-Jahromi et al. [41] "AttentionSiteDTI: an interpretable graph-based model for drug-target interaction prediction using NLP sentence-level relation classification." This is how authors introduced AttentionSiteDTI-an interpretable graph-based deep learning model using protein binding sites and a self-attention mechanism to predict drug-target interaction-and explained its superiority over other models in particular in generalization to new proteins, and further the authors validated some of the predicted drug-target interactions experimentally.

Liyaqat et al. [42] The article "A Short Review on Drug Target Interaction Prediction using Artificial Intelligence" summarizes in brief how artificial intelligence techniques and methods are used for the prediction of ligand-protein interactions to facilitate accelerated drug discovery.

Chen Siqi et al. [43] The systematic review employs the use of "Artificial intelligence-driven prediction of multiple drug interactions" as the title of the paper in giving comprehensive approaches to how much artificial intelligence can be applied in predicting different kinds of drug interactions such as drug-drug interaction, drug-food (excipients) interactions, and drug-microbiome interactions.

Truong Nguyen Khanh et al. [44] Words shine through along the length of musty pages where readers willfully cast their hooks and drawing them out will see that there are AI-based prediction models for drug-drug interactions

based on SMILES in osteoporosis and Paget's diseases. For instance, the authors have developed an AI-based machine learning model capable of predicting drug-drug interactions for osteoporosis and Paget's disease medicines to curb costs and time in getting the best combinations of medications into clinical practice.

Soleymani Farzan et al. [45] Protein – Protein Interaction Prediction with Deep Learning: A Comprehensive Review provides a comprehensive overview of recent deep learning methods applied to various problems related to proteins, including predicting protein functions, protein-protein interactions, protein-ligand binding, and protein design.

Han Ke et al. [46] An investigation of the methods for machine learning-based predictions of drug-drug interactions has shown that Machine Learning has gained acceptance for applications in the field of bioinformatics. It has delivered excellent results specifically in the prediction of drug-drug interaction (DDI) based on drug similarity as a basic criterion, the best efforts at DDI prediction still concern the interactions of drugs in pairs-that is, an advance on the development of methods that can predict DDI for several drugs.

Vo Thanh Hoa et al. [47] The research paper, titled On the road to explainable AI in drug-drug interactions prediction: A systematic review, provides a complete picture of AI-based predictions and sources of drug-drug interaction prediction; data preprocessing; explainable AI mechanisms; modeling approaches; and finally, limitations and future directions of explainable AI in DDI prediction.

Tran T.T.V. et al. [48] Artificial Intelligence in prognosticating Drug Metabolism and Excretion: Recent advances, challenges, and future perspectives. The review paper discusses parametral aspects regarding recent advances, challenges, and future perspectives regarding using artificial intelligence on drug metabolism and excretion.

Dudas Balint et al. [49] "Computational and artificial intelligence-based approaches to drug metabolism and transport prediction." The paper discusses the various computational and AI methods for predicting drug metabolism, transport, and drug-drug interactions, which would be of critical importance in early drug development.

Abdul Raheem et al. [50] Contemporary covers everything-"Comprehensive Review on Drug-target Interaction Prediction- New Developments and Overview." It brings the topical review of the whole current advances and the summary of methods for prediction of drug-target interaction, which is a significant part of drug discovery.

Machado J et al. [51] The Drug-drug interaction extraction system: It is a natural language processing approach. The paper proposes a DDI extraction- based system using natural language processing for processing unstructured biomedical text, applies Drug Named Entity Recognition using the BioBERT model, and classifies the interactions using an ML classifier that has been trained on vectorized text features (BOW, TF-IDF, Word2Vec, custom): Extracted DDI pairs are stored in a relational database.

Li Zimeng et al. [52] DSN-DDI: A Framework for Drug-Drug Interaction Prediction by Dual-View Representation Learning that is Accurate and Generalized. They have proposed a new dual-view drug-representation learning network that predicts DDI ('DSN-DDI') and that employs local-global representation learning modules iteratively for learning drug substructures from single drug ('intra-view') and drug pair ('inter-view').

Hong Eujin et al. [53] "Contemporary Development of Machine Learning Models for the Prediction of Drug-Drug Interactions", it reviews all recent advancements in the machine learning models that have been developed for this purpose.

Fatemeh Rafiei et al. [54] Drug combination synergy was predicted by a recently introduced deep-learning-based multimodal approach, called DeepTraSynergy, utilizing several modalities, including drug-target

interaction, protein-protein interaction, and cell-target interaction. This new approach was able to outperform classical and state-of-the-art models in predicting drug combination synergy in two drug combination datasets.

Truong Nguyen Khanh et al. [55] With this model named, "An AI - Based Prediction Model for Drug-drug Interactions in Osteoporosis and Paget's Diseases from SMILES", it aims to develop an AI based machine learning model predicting drug-drug interactions for osteoporosis and Paget's disease drugs, with a goal of alleviating the cost and time to implement the best combination treatment application in clinical practice.

Ning-Ning et al. [56] Comprehensive Review on Drug-Drug Interaction Prediction by Machine Learning: The Present Condition, Problems, and Their Solutions. The present paper gives a comprehensive account of drug-drug interaction (DDI) prediction in terms of status quo, problems, and prospects in the field of machine learning (ML).

Chen Wei et al. [57] The document examines the usage of artificial intelligence techniques in an assortment of facets of drug discovery-from data resources and molecular representation schemes to algorithms and applications of AI such as predicting drug properties; developing new drugs; and modeling drug-target interactions.

Thi Tuyet Van et al. [58] Recent Developments, Challenges, and Future Perspectives in Artificial Intelligence for Drug Toxicity Prediction. Journal of Chemical Information and Modeling, 2023. This article reviews some key advancements in toxicity prediction of drugs through artificial intelligence; namely, the application of machine learning and deep learning methods and the discussion on data sources and tools, challenges, and future perspectives in the area.

Xuan Lin et al. [59] In-depth analysis of deep and graph learning concerning predicting drug-drug interactions. The paper discusses comprehensive reviews of deep and graph learning methods for predicting DDIs.

Hauben et al. [60] The paper seeks to point out the fact that artificial intelligence and data mining can assist in enhancing detection and prediction of drug-drug interactions in pharmacovigilance. The title of the paper is "Artificial Intelligence and Data Mining for the Pharmacovigilance of Drug-Drug Interactions".

Zhang Yuanyuan et al. [61] "Application of Artificial Intelligence in Drug-Drug Interaction Predictions: A Review." This paper provides a comprehensive review of the state of the art of artificial intelligence applications in predicting drug-drug interactions through three main classes: undirected DDI prediction, DDI-events prediction, and asymmetric DDI prediction.

Razzaghi P et al. [62] The proposed approach in this paper is called CCL-DTI, which incorporates the contrastive loss for drug-target interaction prediction. The innovative approach developed under the name CCL-DTI leverages multimodal evidence and combines task prediction loss with four different contrastive losses to induce a more discriminative feature space for drug-target interaction prediction, thus improving significantly over existing state-of-the-art systems.

Suruliandi Andavar et al. [63] This paper predict drug-target interaction by machine learning techniques - review. The paper reviews machine learning techniques for predicting drug-target interaction, which is one of the crucial steps in the drug discovery process.

Yan Zhao et al. [64], Drug-drug interaction prediction: databases, web servers and computational models, This paper provides an overview of drug-drug interaction (DDI) research, including an introduction to some basic concepts and the classification of DDIs, a description of publicly available databases and web servers for DDI information, and a summary of computational models to predict DDIs.

Zhu J et al. [65] SSF-DDI: a novel methodology in deep learning that evaluates drug sequence and substructure characteristics for drug-drug

interaction predictions. File presentation of SSF-DDI, which draws both drug sequence characteristics and drug molecular graph structural features, is ability to increase the accuracy and the completeness of DDI prediction. This work demonstrates the superior performance of SSF-DDI relative to state-of-the-art DDI prediction models under both the transductive and inductive conditions.

Li Xiong Z et al. [66] Deep learning for drug-drug interaction prediction, as well as a detailed review of a paper on deep learning methods concerning drug-drug interaction prediction important for drug safety and understanding combinatorial therapy.

Luo Huimin et al. [67] The intention of this paper is to systematically review the body of recent research on the prediction of drug-drug interactions based on the use of deep learning and knowledge graph techniques.

Sharmila et al. [68] "A Systematic Review on Drug-to-Drug Interaction Prediction and Cryptographic Mechanism for Secure Drug Discovery Using AI Techniques." This paper provides a comprehensive review of using AI techniques to improve drug-drug interaction prediction and using security enhancement techniques to improve drug discovery, while also discussing the current state, security concerns, and need for effective techniques.

Yasmin radwan et al. [69] Provides a general surveys of overall machine learning methods regarding the prediction of drug-drug interaction" The article provides a comprehensive survey on existing machine learning approaches for predicting drug-drug interactions through literature-based, similarity-based, graph-embedding and classification-based approaches, as well as discusses challenges, datasets and merits/shortcomings of each methodology.

SONAJI P et al. [70] Named as "Artificial Intelligence in Drug Interaction Prediction," this paper describes an investigation into the applicability of artificial intelligence methods-natural language processing, knowledge

graphs, and machine learning algorithms-in enhancing the accuracy and efficiency of drug interaction prediction.

2.3 RESEARCH GAPS

Indeed, DDI prediction has evolved over the years, facing new challenges that hinder current methods' accuracy, scalability, and feasibility. A major drawback here is their over-reliance on text-based approaches; several studies majorly extract interactions from biomedical literature through text mining and NLP techniques. These techniques are helpful, but they often ignore crucial structured data features that would have otherwise boosted predictive performance: pharmacokinetic profile characteristics, molecular descriptors, and chemical properties.

Feature Representation is one more major hurdle. Most of the time n-gram-based vectorization, TF-IDF, or word embeddings would be used by a model for vectorization. These are sprawled out textual relationships but would not have that of the chemistry, pharmacology, or biology between the drugs. Perhaps a more organized, data-driven approach such as molecular feature extraction or pure count vectorization might yield a far more accurate and interpretable model.

Feature Representation is yet another giant step. Uses most of the time n-gram-based vectorization, TF-IDF, or word embeddings; scattered texts for this but would not have made it for chemistry, pharmacology, or biology between the drugs, a more ordered and data-driven approach such as molecular feature extraction or raw count vectorization might give a far more accurate and interpretable model.

In addition, the problem of class imbalance is still a reality in DDI classification. There are large differences between non-interacting drug pairs and interacting ones in real-life datasets. Such discrepancies lead to biased models making them unequipped to predict the minority class in most of such cases. Most of the recent research uses basic undersampling algorithms, which have a drawback of losing some information. A more

effective technique that keeps relevant data intact yet improves generalization of the model is the Instance Hardness Threshold (IHT), which is for undersampling.

The binary classifier, used in several works of literature, simply classifies drug pairs as interacting or non-interacting, and this itself adds another drawback of the proposed method. This method does not actually consider the interaction levels, such as major, moderate, minor, or unknown, which are clinically relevant for decision-making in health care. A multiclass classifier would allow physicians and pharmacists to prioritize those interactions that are highest risk with a prediction more relevant to clinical practice.

Though significant advances in model evaluation and optimization should be carried out, many researchers still shy away from testing ensemble learning methods or deep neural network architectures in order to improve predictive performance on applications employing classical machine learning models like logistic regression, support vector machines, or decision trees. Other modeling candidates were also worth investigating to boost both accuracy and reliability, including SGDClassifier, XGBoost, LightGBM, and CatBoost.

Integration and real-world applicability remain challenges. For instance, most DDI prediction models are only potentially beneficial in research settings because they do not interface with clinical decision support systems (CDSS) in a real-time medical application. Future models should therefore prioritize explainability and usability in order to make certain that predictions are legible, scalable, and immediately useful for healthcare professionals. DDI prediction will be immensely improved with machine learning strategies for clinical use by addressing the mentioned gaps using structured data, advanced feature extraction, improved class balancing, multi-class classification, model optimization, and real-world integration concerning accuracy, speed, and usability.

Presently, computational efforts toward DDI prediction do not seem to be pragmatic in terms of application in the clinical field. Although the performance of the established models appears to be good in theory, they are seldom validated in a diverse real-world setting and on patient populations limiting the generalization of the findings to other demographic groups and comorbid conditions. Furthermore, many studies do not consider the other hand of the dynamic nature of drug interaction, where change in dosage, time of administration, and patient-specific metabolic profiles could play substantial roles in the risk of interaction.

CHAPTER 3

ALGORITHM(S) ANALYSIS AND DESIGN

3.1 INTRODUCTION

The importance of algorithm design and analysis is essential in the development of a truly worthy and effective machine learning model for drug-drug interaction (DDI) prediction. This chapter presents an orderly approach to the automation DDI classification problem involving the selection of appropriate algorithms, strategies for data preprocessing, and evaluation of model performance. The aim is such an efficient machine learning system that can reliably predict drug interaction severity, keeping in mind the computational cost, scalability, and interpretability of the results.

The chapter begins with a description of the problem and the necessary system requirements for a predictive modeling process. Then, it discusses the choice and justification of machine-learning methods focusing on their applicability to structured data within drug-drug interactions. Furthermore, it describes the feature engineering techniques used to transform the raw data into a more suitable classification-friendly format. Lastly, this chapter touches on the problem of class imbalance in biomedical datasets, addressing some balancing techniques that are supposed to enhance the generalizability of the model.

A formalized model training strategy is then presented, including methods of data splitting, cross-validation. Parameter tuning and an extensive complexity analysis-all chosen by the user algorithms, as being able to evaluate the computational time, spatial requirements and scalability. The aim of this chapter is to give a structured, orderly reasoning to take design machine learning pipelines for DDI prediction so that they would have effective performance and real-world utility.

Carefully laid out, this chapter addresses the methodical approach toward model development-crediting the background that goes into building an easily computable, interpretable, and clinically relevant DDI predictive machine learning model. Herein are furnished details of the systematic methodology employed in designing, evaluating, and optimizing predictive algorithms capable of accurately identifying and classifying DDIs.

Pharmacological interaction DDI prediction goes beyond simple classification. It differs from traditional problems in machine learning in that it involves several factors such as molecular properties, metabolic pathways, and clinical severity levels in DDI classification. Our approach is therefore a structured workflow that begins with rigorous data preprocessing, then features smart algorithm adoption and solid performance validation.

3.2 SYSTEM REQUIREMENTS

The construction of an effective drug-drug interaction prediction model will require quite specific requirements for the system, covering aspects like accuracy, scalability, and usability. Some major requirements are reading large and complex biomedical datasets, dealing with interactive complexity, and classifying correctly severity levels for taking informed clinical decisions.

3.2.1 Problem Definition: Need for Automated Drug Interaction Prediction

The growing trend of polypharmacy is gradually increasing the risk of adverse drug reactions resulting from unrecognized or unforeseen drug-drug interactions. The traditional techniques in the detection of DDIs include clinical trials, expert evaluations, and post-market surveillance; however, these techniques have proved to be time-consuming, expensive, and unable to cope with the accelerated introduction of new medicines. Therefore, it is high time to create an automated machine learning-based system that procures an early prediction of probable drug interactions based on structured biomedical datasets to make the considerations scalable and

efficient as opposed to DDI manual identification approaches and to minimize the risk of any severe medication-related complications.

3.2.2 Input Data

The dataset consists of drug pairs and associated severity categories indicating the extent of interaction between the two drugs. These pairs, Drug A and Drug B, have additionally attached a label that indicates the severity of interaction between the two drugs under those categories, i.e., Major, Moderate, Minor, and Unknown. The data set is preprocessed to eliminate duplications, deal with missing entries, and finally uniform data. Furthermore, feature extraction methods such as Count Vectorization are applied to translate categorical names of drugs into numerical representations required by machine learning models.

The system operates primarily by utilizing a structured dataset containing detailed information about drug pairs and the intensity of their interactions. Each entry in that dataset is characterized by two important elements called Drug A and Drug B and a label that indicates how severe the interaction is. The clinical severity is assigned through four different labels, providing four unique classes for the interaction nature: Major, Moderate, Minor, and Unknown. This structured categorization is very important for machine learning models to differentiate between interactions that need urgent clinical attention and those that might be less urgent or poorly understood.

Direct application of this dataset for training or assessment would require intense preprocessing. Redundancy is one of the first considerations to be eliminated to prevent bias and bias model learning. Duplicates are often introduced when the same interaction is cited more than once across various sources, or when drug pairs enter either one of two orders (for example, Drug A–Drug B and Drug B–Drug A). Resolving these problems allows for the model to be trained on clear and distinct interaction examples.

Besides duplicate elimination, the system must also deal with any missing values. Incomplete severity labels, inconsistent formatting, or a lack of drug

identifiers can influence the learning process negatively. To retain the integrity of the dataset, the missing values must be imputed wherever applicable or the affected records deleted. Equally important is to avoid inconsistency in the dataset; practices of drug name standardization, label confirmation, etc. The reliability of the subsequent feature extraction and model training phases and, consequently, less error consequent on consistent formatting of data will be assured. Converting the categorical drug information into numbers that machine learning algorithms will understand is the next major step after cleaning and standardising the dataset. For this, count vectorization and other feature extraction techniques have been used. Count vectorization converts drug names to numerical vectors by counting the instances of different tokens (drug identifiers) across the dataset. The technique thus enables models to spot and learn patterns about how often or combinations of drug pairs to predict the severity of interactions for previously unobserved combinations of drugs.

Such a well-organized and systematic way of preprocessing the dataset and correctly employing feature transformation techniques ensures that data fed into machine learning models is clean, standardized as well as very much predefined with significant patterns. This careful preparation helps the system to yield such precise, reliable, and clinically relevant DDI predictions in real-world healthcare settings.

3.2.3 Output Data

They expect to output a model to classify the interaction between two drugs according to the severity of the interaction into one of four predefined classes: Major Interaction, Moderate Interaction, Minor Interaction, or Unknown Interaction. The predicted results must be interpretable and trusted such that healthcare professionals could rely on the model's predictions in their decision-making processes.

3.2.4 Performance Goals

In drug-drug interaction (DDI) classification, the chosen model should ideally exhibit very high precision and recall, especially when identifying critical interactions, namely those flagged as Major and Moderate. The erroneous identification of these interactions could severely impact the patient since they usually have important clinical implications. Thus the reduction of false positives (interactions wrongly predicted when none exist) and false negatives (interactions not detected) is very important.

In addition to high predictive power, the model must be quite robust to class imbalance. Interactions of clinical importance (e.g., major interactions) are often rare as compared to minor or unknown interactions in drug-drug interaction datasets. Most machine learning algorithms find skewed distributions an enormous challenge and tend to create models biased toward the majority classes. In order not to miss interactions that are rare yet clinically significant, the system must be intrinsically strong to imbalances, thereby capturing even the subtle patterns which differentiate these rare incidents from the more common ones.

A thorough machine learning evaluation metric differentiation would involve precision, recall, and, most importantly, the F1 -score. Recall, for one's knowledge, is about detection of major portions of actual interactions successfully, while precision ensures that fewer false interactions are tagged for every unit. The measure is indeed the one that tends to balance because it is mostly useful where class imbalance occurs: the harmonic mean between precision and recall.

Secondly, the final model for DDI classifications should demonstrate high generalization capabilities; that is, it should be capable of making reliable and accurate predictions for cases that are independent and unseen even to training datasets. Generalization of a model is very important to the real-world application in health sciences, as the data landscape keeps on evolving with new drug and interaction discoveries. If the end goal is to support clinical decisions, such a model would be completely useless if it

overfits the training data.hence, a model selection is beyond a static benchmark; it also involves a solid and generalizable system that trains on new data and keeps bringing clinically relevant discoveries in constantly changing healthcare environments.

3.2.5 Scalability

This algorithm must be able to manage any biomedical data of very large volume quickly and in real time because it continues to grow with the introduction of more and more new drugs. Along with being able to learn through big data, the system has also to support some complex multi-dimensional feature spaces, often consisting of thousands of drug interactions, thus making it a very optimized computation efficient training and interpretation.

Incremental learning techniques should be included in the adaptability of the implementation so that it can be incrementally updated with new drug interaction data. The integration of the optimized implementation for the system is required for low memory consumption but high-speed execution for real-world applications such as Clinical Decision Support Systems, common in hospitals and research facilities. Hence the system requirements for this prediction model would provide a scalable, predictive model for drug interactions that is accurate and interpretable in terms of actual applications upon bringing them to practice, resulting in safer medication practices and, consequently, better patient care with the DDI prediction model.

3.3 ALGORITHM SELECTION AND JUSTIFICATION

3.3.1 Considered Algorithms

Machine learning models predict Drug-Drug Interaction based on detailed assessment of alternative classification approaches. Classification techniques, both classical and ensemble learning, as well as gradient boosting methods, are all evaluated for use on structured biomedical data.

Machine learning models to predict Drug-Drug Interaction require exhaustive evaluation on alternative classification algorithms. Classical, ensemble learning, and gradient boosting-based classification techniques have all been analyzed for usage on structured biomedical data. Machine learning models predict Drug-Drug Interaction using an exhaustive study to alternate classification algorithms. For the evaluation of both traditional and more recent ensemble learning techniques and their gradient-boosted designs, structured biomedical data were applicable. Essentially the Stochastic Gradient Descent Classifier allows for incremental learning and is, therefore, suitable for large datasets. Its ability to modify parameters after each sample is inspected minimizes memory consumption and elegantly addresses high-dimensional feature spaces. Other gradient boosting algorithms have been put forward, which include the Extreme gradient boosting, Light gradient boosting, and Categorical boosting models. These models are known to best classify structured data, especially when the datasets are imbalanced.

While giving parallel computation efficiency, regularization optimization techniques, and sparse dataset handling capabilities to Extreme Gradient Boosting, Light Gradient Boosting Machine is fast and memory efficient, which makes it particularly suitable for large data in biomedical sciences. It also applies leaf-wise growth pattern formulation for accuracy improvement and less time consumption during training. The model is capable of automatically working with categorical variables without going through too much data pre-processing. Also, it is viewed as a model capable of handling overfitting which is known to occur with tree-based models, hence it is a very promising candidate for Drug-Drug Interaction classifications. In the context of performance benchmarking, a baseline logistic regression is included. Although it is much less complex than other models, it has the advantage of being interpretable, providing a reference point for comparing relative improvements achieved with more complex classifiers.

The Random Forest model, which builds multiple decision trees, can indeed bolster the degree of prediction stability and reduce variance in the output made by the model. Furthermore, it can also return the importance scores for the features so that the most significant factors related to the drug interactions can be determined. This aspect of interpretability is particularly important in the clinical setting, where there should always be some explanation coming along with decision-making.

3.3.2 Justification for model selection

A range of important criteria were considered in the selection of the models, such as accuracy, scalability, ability to account for disproportionate representation, and the level of explanation of the decision-making process. The major reason for the use of boosting-based models is that they can learn complex patterns in structured data and can effectively handle unbalanced datasets. Besides that, its boundaries could be fine-tuned in order to have much better performance in classifying different interaction types. This makes them particularly suited for ensuring that important drug interactions do not get missed, since the model is capable of assigning high importance to minority classes.

The SGD Classifier was chosen for its efficacy in a large-scale setting under constraint of computation. The ability to learn incrementally allows it to be updated dynamically with new interaction data, thus making it a scalable option for real-world implementations. They include Random Forest and Logistic Regression as baseline models to achieve a balanced comparison. Although they may not be able to compete with boosting algorithms in terms of precision, they would prove handy for understanding model behavior and feature importance. Logit is the reference model for comparing advanced techniques. Having a good availability of different kinds of machine learning models enables one to make a thorough comparison between different techniques. Doing so guarantees the final model is able to achieve the best trade-off of accuracy, efficiency, and applicability to the real-world problem at hand in drug-drug interaction classification.

3.4 DATA PRE-PROCESSING AND FEATURE ENGINEERING

3.4.1 Dataset Collection and Cleaning

The dataset being used here has structured information about drug pairs and their effects. Owing to inconsistencies, duplicates, and sometimes irrelevant information in raw data, it is imperative that dataset cleaning be done prior to model training, which would lend credence to the model's prediction and performance. The removal of duplicates guarantees that every drug pair is unique, avoiding any bias or imbalanced results leading to overfitting. Standardizing the drug names adds another layer of consistency to the feature analysis, which prevents wrong representations caused by differing terms such as abbreviations or alternative spellings. Dataset cleaning ensures that each drug interaction is treated accurately, providing an extra edge in terms of reliability of data and the performance of the model.

3.4.2 Feature Extraction

Feature extraction is the foremost essential step in the implementation of machine learning for classification, which converts unprocessed or categorical data into numerical representations that could be understood by the model. Drug-drug interaction (DDI) prediction converts drug names and interaction-based properties into useful numerical features while keeping the most significant relationships.

A range of methods can be used for feature extraction from structured and unstructured biomedical data. This presentation will discuss some of those methods for feature extraction as applied in biomedical research.

- **Count Vectorization:** Counts the occurrences of words or tokens to transform categorical text—like drug names—into numerical representations. Performs well when features are structured and discrete, guaranteeing interpretability and reducing computational complexity.

- Term Frequency-Inverse Document Frequency(TF-IDF):It assigns importance scores to words based on their frequency within a document and across a dataset.
- One-Hot Encoding: This amethod converts categorical variables into binary vectors, thereby ensuring each unique category receives its own feature representation.
- Word Embeddings(Word2Vec, Doc2Vec, BioBERT): These map words or phrases into dense vector spaces, capturing the semantic relationships between terms.
- Molecular Fingerprinting (Morgan Fingerprints, MACCS Keys): This is utilized in chemical informatics, where drugs are represented based on their molecular structure.

The choice for a feature extraction technique is mainly determined by the kind of dataset and the amount of complexity modeled in the relationship since unstructured text-based datasets can benefit from these: TF-IDF; Embeddings; NLP techniques. Structured datasets would rather get their feature extraction through count vectorization; one-hot encoding; or molecular fingerprints. It is feature extraction that helps the researcher in drawing conclusions from biomedical literature as well as clinical data text wise.

3.4.3 Handling Class Imbalance

Class imbalance is a very common issue in classification exercises. It becomes even more pressing in biomedical and health datasets where certain classes are poorly represented (for example, rare but very beneficial interactions). If class imbalance is neglected, a model may end up learning mainly about majority classes while ignoring the few rare predictions. This leads to biased predictions.

For example, to take care of class imbalance and ensure effective generalization of models across categories, various strategies have been employed:

3.4.3.1 Methods of Undersampling

- Balances the dataset by lowering the number of instances in the majority class.
- Samples from the majority class are randomly removed using Random Undersampling.
- In order to ensure that the model concentrates on more difficult cases, the Instance Hardness Threshold (IHT) selectively eliminates samples from the majority class that are easily classified.

Ideal for situations where training is dominated by majority-class instances but the dataset is sizable.

3.4.3.2 Techniques for Oversampling

- Improves representation by increasing the number of instances in the minority class.
- Random Oversampling: Replicates instances of the minority class that already exist.
- The Synthetic Minority Over-sampling Technique, or SMOTE, uses nearest-neighbor interpolation to create synthetic samples for the minority class.

Beneficial when minority classes are severely underrepresented and the dataset is small.

3.4.3.3 Methods of Hybrid Sampling

- Refines dataset balance by combining undersampling and oversampling.
- Tomek Links: Enhances the clarity of decision boundaries by eliminating ambiguous instances where majority and minority classes overlap.
- Before using resampling techniques, noisy samples are cleaned using Edited Nearest Neighbors (ENN).

Used when it's necessary to increase minority representation and eliminate redundant majority samples.

3.4.3.4 Cost-Aware Education

- Gives underrepresented classes greater misclassification penalties rather than changing the dataset
- Makes certain that during model training, rare interactions are given priority.

Just right for those scenarios where there have been at least some risks of bias or overfitting. Most use the dataset modifications (oversampling or undersampling). The application generates effects of which technique is best to use; the size of the data set, the degree of imbalance, and the application requirements. Undersampling works well for large datasets but is not efficient for those with a few number of minority-class examples which are much better off oversampling and cost-sensitive learning. High-risk settings like healthcare often employ hybrid techniques to ensure their predictions are accurate, fair, and clinically relevant.

Just right in case there could be any bias or overfitting risks introduced by modifications on the dataset, oversampling or undersampling. Application trails import the intended effects of which technique is best to use; size of the dataset, extent of imbalance, and requirements by the application. While undersampling works well for fairly large datasets, with few samples from the minority class, it is better off going for oversampling and cost-sensitive learning. Hybrid approaches are mostly used in high-risk areas such as health care in order to ensure correctness, fairness, and applicability in clinics. By utilizing efficient techniques for feature extraction and balancing classes, machine learning models can better comprehend, generalize, and predict correctly in drug interaction classification tasks.

3.5 COMPLEXITY ANALYSIS

Certain insights into the computational complexity of the specific machine learning models relevant are needed in the design of dependable and efficient large-scale Drug-Drug Interaction (DDI) prediction systems. The complexity analysis is fundamental when it comes to making a decision

about the operational practicality and accuracy of a model that interactively handles increasing amounts of structured biomedical data.

Using complexity analysis, researchers can analyze how the requirements for model training and inference change as a function of the size of the dataset, the dimensionality of the feature space, and distribution of classes in specific scenarios such as DDI classification. It is also very important to select models that are able to process biomedical datasets, often high dimensional and potentially containing millions of interaction pairs, using computational costs that are not exorbitant.

By systematically evaluating computational complexity, it's possible to anticipate possible bottlenecks during model training and deployment phases. And from the evidence brought forward by using real-world DDI databases, algorithms with such high algorithmic complexity may score well on prediction using smaller datasets but become unfeasible by virtue of slowness or memory requirements. In contrast, as new drug-drug interactions will be defined and added to the databases in time, models that have that ideal trade-off will ensure that the system will be responsive-scaling-down to clinically useful. Likewise, the plan of hardware resources gains substantially from complexity analysis. In application of machine-learning models to real-world healthcare settings where compute resources might be expensive or limited, it is imperative to establish the memory requirements, CPU/GPU usages, and parallelization strategies, depending on the computational burden posed by different algorithms.

The study considers complexity analysis as one of the essential components of model evaluation, rather than being merely a theoretical exercise. In prioritizing models that could achieve high accuracy within reasonable computational constraints, complexity analysis ensured that the eventual DDI prediction system would be able to meet both operational and scientific goals. The approach highlights the importance of looking beyond standard performance metrics and considering computational pragmatism in choices

for machine-learning models in any application, particularly in high-stakes biomedical scenarios.

3.5.1 Time Complexity

The computational cost of training and forecasting in contrasting machine learning models is referred to as time complexity. In practical scenarios, the efficiency of an algorithm is important, especially with a lot of data being present, such as the records of interactions of drugs. Linear models (Logistic Regression, SGDClassifier) themselves are computationally efficient for large datasets, that is, they have lower training complexities in general.

Table 3.1 Time Complexity Analysis

Model	Training Complexity	Prediction Complexity
Logistic Regression	$O(n^2)$	$O(n)$
SGDClassifier	$O(n \log n)$	$O(n)$
Random Forest	$O(n \log n)$	$O(\log n)$
XGBoost	$O(n \log n)$	$O(\log n)$
Lightgbm	$O(n \log n)$	$O(\log n)$
CatBoost	$O(n \log n)$	$O(\log n)$

- For large datasets, linear models (Logistic Regression, SGDClassifier) are computationally efficient due to their generally lower training complexity.
- The ensemble learning techniques used by tree-based models (Random Forest, XGBoost, LightGBM, and CatBoost) lengthen training times while increasing classification accuracy.
- Although they need to be carefully adjusted, gradient boosting models (XGBoost, LightGBM, and CatBoost) perform better in structured datasets like DDI classification.

3.5.2 Space Complexity

Space complexity is the amount of memory required to accommodate intermediate calculations and model parameters. In large biomedical datasets, where effective memory management is necessary, this becomes important.

- Because Random Forest and XGBoost store multiple decision trees, they are computationally intensive and require more memory.
- Large datasets with high-dimensional features can benefit from the memory-efficient nature of SGDClassifier and Logistic Regression.
- LightGBM and CatBoost are scalable substitutes for conventional tree-based models because they maximize memory usage.

A model selection endeavor must as much as find out a happy medium between an accuracy requirement and the memory limitation so that the system would be able to craft a strategy that would effectively tackle large number of dimensional data.

3.6 SCALABILITY

This property of scaling puts up well in the training process and applications of models where the volume of data increases. Scalability is important for DDI prediction because new drugs and new interactions would always be discovered.

Minor revision: Scaling means the capacity of a model to perform well as the dataset grows in size. In DDI prediction, scalability becomes critical because new drugs and new interactions keep being discovered. The property of scaling serves very well in the training process and applications of models where volume data are on the increase. Scalability plays an important role in DDI prediction as new drugs and new interactions are usually discovered.

- Because of its incremental learning methodology, SGDClassifier scales effectively, which makes it perfect for real-time updates in drug interaction systems.

- High accuracy is provided by the boosting models (XGBoost, LightGBM, and CatBoost), but in order to maximize speed and memory usage, careful hyperparameter tuning is needed.
- Random Forest works well with medium-sized datasets, but as the dataset gets bigger, it might have trouble scaling.

To keep DDI prediction systems responsive and dependable in real-world applications, models need to be tuned for both speed and accuracy.

3.7 PERFORMANCE METRICS AND EVALUATION STRATEGY

Diverse measures for evaluating the efficiency of different machine learning models, which help in ascertaining how accurately and reliably these models are generalized to infer about new data for the DDI classification model.

- Precision: Assures that there are few false positives by measuring the proportion of predicted interactions that are actually true.
- Recall: Assesses how well the model can identify all real interactions while lowering false negatives.
- A balanced classification performance is ensured by the F1-score, which is a harmonic mean of precision and recall.

Each of these metrics gives a different perspective on model accuracy, helping to decide the most effective method for classifying DDIs.

3.8 SUMMARY

This is the chapter that presents the complexity analysis, techniques for performance evaluation, and those critical scalability concerns that enable one to create an efficient DDI prediction model. Models can balance between scalability and accuracy or computational efficiency based on knowledge of these parameters. Next chapter, "Implementation and Results," will present actual tests and the study of the model on real data.

CHAPTER 4

IMPLEMENTATION AND RESULTS

This chapter presents an elaborate insight into the actual real-world applications of the techniques contemplated in the Algorithm Analysis and Design chapter. It details the various implementation steps from preprocessing datasets, model training, and model evaluation, thus ensuring a very systematic approach in the context of machine learning-based Drug-Drug Interaction (DDI) prediction. In the opening part of this chapter, the second section deals with the data preparation and the known preprocessing activities, focusing on the feature extraction, transformation, and cleaning methods applied to transform the unstructured biomedical data into a form amenable to interpretation by machine learning. The dataset should be able to house all missing values, cope with duplicated records, encode categorical features, and balance class distributions for maximally accurate predictions.

The subsequent section, concerning model training and optimization after preprocessing of data, then describes in detail the machine learning algorithms implementing DDI classification. This covers the model selection and computational aspects to ensure peak performance. The aspect discussed in this chapter covers the generic ability of the models to learn different drug interaction scenarios, explaining how they were trained using structured datasets.

These evaluation metrics and validation methods are used to analyze performance of the trained models. Every model is measured based on precision, recall, and F1-score, which together signify its accuracy, robustness, and dependability. In addition, comparative and error analyses help to find advantages and disadvantages in each approach along with possible areas for development.

Understanding the accuracy, robustness, and reliability exhibited by each model shall use the key metrics such as precision, recall, and F1-skore. In

addition to that, comparative and error analyses are also done to deduce the advantages and disadvantages on each approach, as well as the possible areas for improvement. Evaluating the performance of the trained models is carried out through the use of evaluation metrics and validation methods. The measure of performance of each model includes key metrics such as precision, recall, and F1-score, which together represent accuracy, robustness, and dependability. Furthermore, comparative analysis and error analysis are also conducted to define the merits and demerits of each approach and possible room for improvement.

It works on the trained model relative to performance evaluation using evaluation metrics and validation methods. The measure of performance of each model includes key metrics that are precision, recall, and F1-score, which together signify accuracy, robustness, and dependability. In addition, comparative and error analyses are performed to also get hold of advantages and disadvantages on each approach along with possible areas for improvement. This train assessment applies to evaluation metrics and validation methods. Performance of each model also contains the following key metrics: precision, recall, and F1-score, which provide an overall picture of the model's accuracy, robustness, and dependability. Apart from that, comparative and error analysis is done to understand the pros and cons of each model and to explore further potential improvements in that direction.

The performance comparison is exhaustive, using various models to ascertain the most efficacious one with respect to machine learning for accurately predicting drug interactions. Informed and endorsed by this deliberation, the selection of what is optimal suffices, ensuring that the method forwards itself accordingly in a clinical environment, adopting a scalable dimension, and interpretable for pharmacovigilance and healthcare decision-making in the real world. Again, before anything else, here is a very exhaustive comparison among several models, aimed at identifying the best possible machine learning method for accurate prediction of drug interactions. This selection is informed and sanctioned by this deliberation.

In this way, optimality suffices in a clinical environment whilst adopting a scalable dimension and interpretable for pharmacovigilance and healthcare decision-making in the real world.

4.1 DATASET DESCRIPTION

The dataset used for the Drug-Drug Interaction (DDI) prediction is from the publicly available drug interaction database DDInter (<https://ddinter.scbdd.com/download/>). The dataset qualifies for clinical machine learning-based classification due to the structured biomedical data on drug interaction severity levels.

4.1.1 Overview of the Dataset

The dataset constitutes 221,841 rows and 5 columns corresponding to different drug pairs and the severity of interaction associated with them. The characteristics are as follows:

1. DDInterID_A : The first drug's unique identifier .
2. Drug_A : It is the name of the initial medication that came into contact with it.
3. DDInterID_B : The second drug's unique identifier .
4. Drug_B: This is the name of the second medication that interacts with it.
5. Level: The interaction's degree of severity, divided into:
 - Major: A high-risk interaction that could have serious negative consequences.
 - Moderate: An interaction that necessitates dosage modifications or monitoring.
 - Minor: A low-risk interaction that has few consequences.
 - Unknown: There is no interaction between provided drugs

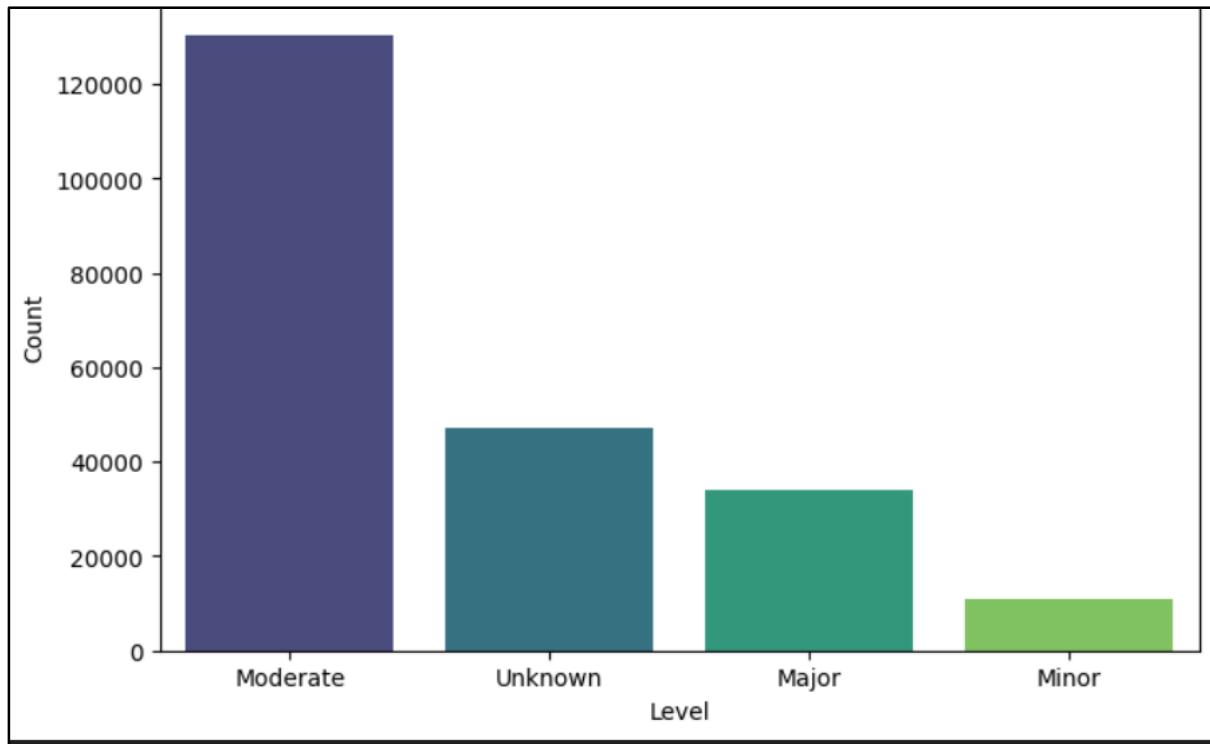
Due to the wide-ranging coverage of identified drug interactions available in this dataset, it has been identified as the ideal instrument for training and testing machine learning models for automated DDI conditioning. In

practical modeling terms, this chapter provides a comprehensive description of how the techniques spoken in Algorithm Analysis and Design were implemented. The section discusses the entire implementation process-from dataset pre-processing to model training and assessment-giving a methodological view on machine-learning-based Drug-Drug Interaction (DDI) prediction. The first part of this section is totally devoted to the description of data preparation and pre-processing techniques and specifically regarding the applicable use of various feature extraction transformation and cleaning operations to convert unstructured biomedical data into a machine-readable form. It is necessary to deal with missing values, duplicate removal, categorical encoding, and balanced classes so that the data is modeled for accurate predictions.

The next section describes model training and optimization after data preprocessing in detail, describing the machine-learning algorithms used for DDI classification. To ensure performance, this includes model selection and computational considerations. The chapter also elaborates on how the models were trained with structured datasets to develop generalizability across different drug interaction circumstances.

The performance of trained models is then evaluated using metric evaluation and validation methodologies. Each model's accuracy, robustness, and dependability are measured against key indicators, such as precision, recall and F1-score. Comparative as well as error analyzes shall also be performed so that strengths, weaknesses and possible further development areas could be derived for each approach. In summary, this work ends with an exhaustive performance comparison of many models to find the best machine-learning method for predicting drug interactions. The winner of this analysis is the best-performing model, which ensures that the method is clinically relevant, scalable, and interpretable in a real-world context for pharmacovigilance and healthcare decision-making.

The dataset is composed of 221,841 rows and 5 columns, representing various pairs of drugs and the severity of interactions for each considered. The dataset is characterized by: The dataset is chosen for its comprehensive coverage of known drug interactions; thus, it serves as the best candidate to train and assess machine learning models for automated DDI



classification.

Fig 4.1 Initial count of each level.

4.1.2 Features of the Dataset

- More than 200,000 drug interactions make up the extensive dataset, which guarantees a rich feature space for model learning.
- Drug names that are categorical and need to be converted into numerical representations for machine learning models.
- Class balancing techniques are required during preprocessing due to an imbalanced class distribution, where certain interaction severity levels (such as major interactions) are underrepresented.

The purpose of this project is to convert the dataset into a scalable and accurate machine learning model that can predict possible drug interactions ahead of time so that patient safety can be enhanced, and the manual clinical evaluations can be lessened.

4.2 DATA PRE-PROCESSING

To prepare the dataset for the Machine Learning-Based Drug-Drug Interaction (DDI) prediction, data preprocessing is an essential step. In this phase, the dataset is cleaned up for various aspects like transformations for categorical features, dealing with class imbalance, and finally, made ready for model training. The goal, in short, is to produce an accurate, organized, and relevant dataset for classification models.

4.2.1 Data Cleaning and Duplicate Removal

Initially, the raw dataset consisted of 222,843 drug interaction records, which might have included some duplicated entries among them. Such duplication of interaction may over-represent certain drug pairs, cause the model to tend to overfit, and give wrong predictions. In preprocessing, duplicate rows were identified and removed for the sake of ensuring integrity of data. Converted into an entirely different version.

Initially, the raw dataset included 222,843 interaction records; some of these possibly contained duplicate records. For example, over-representation of certain drug pairs may distort the actual distribution, creating an overfit in the model, thereby giving it an erroneous prediction. Duplicate rows were found in the pre-processing stage and removed for the sake of ensuring integrity of data.

A count for each level of interaction severity, according to the extent of duplicates in the dataset:

Table 4.1 Distribution of Classes Prior to Eliminating Duplicates

Interaction Severity	Count Before Deduplication
Major	33,896
Moderate	1,30,367
Minor	10,938
Unknown	47,182
Total	222,383

It was made sure that every interaction of the drug pair took place only once through duplication removal on Drug_A, Drug_B, and Interaction Severity. The new class distribution after the deletion of 62148 duplicate rows was Presenting that analogously, overrepresented interactions are not given partial treatment when models are learned by reducing the effect of the noise.

Table 4.2 Distribution of Classes Following Duplicate Removal

Interaction Severity	Count after Deduplication
Major	26,914
Moderate	96,675
Minor	6833
Unknown	29813
Total	160235

4.2.2 Feature Transformation

A crucial step in preparing categorical data for machine learning models is feature transformation, such as turning drug names and levels of interaction severity into an integer form. Because these variables are not numeric, they must be transformed before being processed by classification algorithms. This transfer will allow the model to learn by seeing which drug pairs have particular patterns and how severe the interaction is going to be. Feature transformation is one of the most important steps in preparing categorical data for machine learning models. Drug names and their levels of interaction severity must be transformed into a numerical representation prior to being passed on classification algorithms. This conversion would allow the model to learn the patterns and relationships between drug pairs and the severity of the interactions between them.

Two main methods for feature transformation were applied in this project:

- **Count Vectorization:** Drug names are transformed into numerical feature vectors using this technique, which is applied to Drugs A and Drug B.
- **Label Encoding:** Used to translate interaction severity labels into numerical values in the Level column.

All transformations were done on the dataset to offer the utmost in computing performance with the optimum model interpretability and the highest prospect for effective learning.

4.2.2.1 Count Vectorization for Drug Pair

Count Vectorization is one of the most popular feature extraction methods used for converting text or categorical data into a numerical vector space: in fact, it numerically encodes the frequency of each of the distinct categories (here, drug names) in that particular dataset into a feature vector.

Such machine learning models do not communicate with Drug_A and Drug_B directly, as they are categorical variables. Count Vectorization was favored over other methods of feature transformation for its:

- **Interpretability:** Count Vectorization offers a direct numerical mapping for every drug, which facilitates the analysis of model predictions in contrast to word embeddings or deep learning-based representations.
- **Efficiency:** Unlike NLP-based methods like Word2Vec or Doc2Vec, it is a computationally efficient and lightweight method that doesn't require extra training on external datasets.
- **Capability to capture relationships between drug pairs:** The model can learn drug interaction patterns without requiring extra biomedical knowledge by converting drug names into numerical feature vectors.

For various drugs, if Drug_A and Drug_B are tokenized, a representation with a unique number might be assigned. Hence, each drug name can be represented in an indexed feature space where the features are counted.

These vectors are then fed to the machine learning model, allowing it to learn interaction relationships easily. This transformation enhances classification performance by allowing the model to identify drug pair interactions using their provided numerical values without any complex feature engineering.

4.2.2.2 Label Encoding for Interaction Severity Levels

In the case of Level (Interaction Severity), Label Encoding was employed, while Count Vectorization was utilized to affect Drugs A and B. The categorical values in the Level field (namely, Major, Moderate, Minor, and Unknown) need to be transformed into numerical representation prior to being used in a machine-learning model. In this, the following aspects influenced the choice of Label Encoding for the transformation:

- **Effective representation:** Gives every class a distinct numerical value, guaranteeing a condensed and organized format.
- **Computational simplicity:** Label Encoding is more effective for large datasets because it requires less memory and computational resources than One-Hot Encoding.

- Machine learning model compatibility: When categorical values are encoded as integers instead of distinct binary columns, many classification models perform better.

The resultant encoded values become the target variable for training the machine learning model. With Label Encoding for interaction severity levels and Count Vectorization for drug names, the dataset was translated into a structured numerical format, conducive to effective processing by machine learning models. These approaches are thus apt for the DDI prediction problem, guaranteeing computational efficiency, interpretability, and scalability.

4.3 HANDLING CLASS IMBALANCE

When it comes to DDI prediction, class imbalance is a well-known concern, especially for some interaction types that are notoriously underrepresented. In this dataset, major interactions pertaining to patient safety are much less frequent than moderate and minor interactions. This imbalance, if not properly handled, may allow machine learning models to be biased toward the majority class, thus jeopardizing recall toward infrequent interactions that are actually important.

Two class balancing strategies were assessed in order to lessen this problem:

- Instance Hardness Threshold (IHT) undersampling (selected as the optimal approach)
- The Synthetic Minority Oversampling Technique (SMOTE) is an oversampling method that is not appropriate for this dataset.

4.3.1 Undersampling using Instance Hardness Threshold(IHT):

One of the advanced undersampling methods is called Instance Hardness Threshold (IHT). It keeps the difficult-to-classify cases and takes out the easily classifiable samples from the majority class. Thus, it diverts the attention of the model towards the complex cases, which helps augment the model's ability to discriminate among similar interaction types.

Justification for using IHT:

- **Rationale for Preventing Overfitting to the Majority Class with IHT:** In contrast to random undersampling, which eliminates samples at random, IHT eliminates non-informative majority-class samples specifically, guaranteeing that the model keeps challenging examples that improve learning.
- **Better Decision Boundaries:** The model is compelled to learn more exact patterns by removing redundant and easily classified samples, which improves generalization across all interaction types.
- **Balanced Dataset Representation Without Artificial Data:** IHT maintains the dataset's realism and clinical validity by not producing synthetic samples, in contrast to oversampling techniques.
- **Improved Recall Without Sacrificing Precision:** By eliminating excess majority-class samples, the model was able to concentrate on minority-class occurrences, which improved recall without appreciably raising false positives.

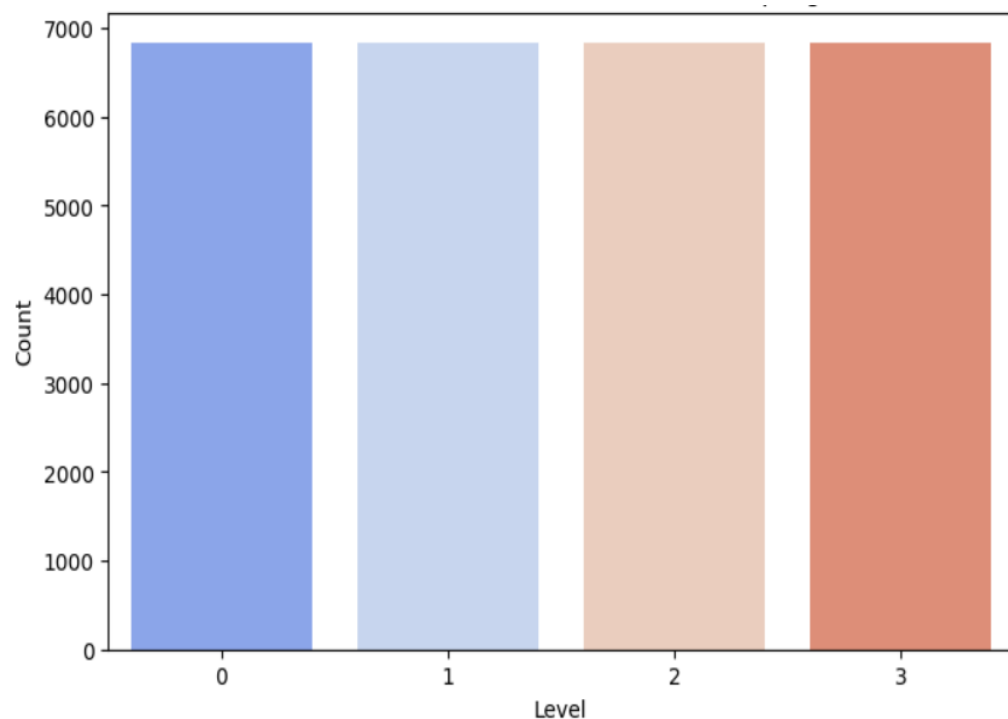


Fig - 4.2 Class Distribution after applying IHT

4.3.1.1 IHT's effects on the dataset

- The dataset decreased bias toward moderate and minor interactions while preserving enough information for model training.
- Enhanced recall for significant interactions, which guarantees improved critical case detection.
- Reduced the possibility of overfitting, resulting in a classification model that is more impartial and comprehensible.

Among-classes balancing IHT was selected as a balancing method mainly due to these advantages posed to the DDI classification model.

4.3.2 Oversampling Using Synthetic Minority Over-sampling Technique (SMOTE)

The Synthetic Minority Over-Sampling Technique, or SMOTE, is a popular means of oversampling and involves interpolating between existing samples to generate synthetic data points for underrepresented classes. The objective of SMOTE is to achieve a balanced dataset by augmenting the representation of minority classes (for example, major interactions).

The limitations of SMOTE on this dataset are:

Here SMOTE failed to outperform the ensemble of non-resampling methods on EU-ADR dataset, which might be attributed to the DDI dataset's behavior, where the SMOTE is performing so well in most of similar problems.

- Synthetic Noise Introduction: Because drug interactions are so specific, creating synthetic examples produced irrational data points that caused misclassifications. As a result, the model's predictions were less reliable.
- Failure to Improve Decision Boundaries: SMOTE adds synthetic instances without enhancing the model's capacity to identify significant patterns, in contrast to IHT, which improves decision boundaries by eliminating redundant data. The classifier thus kept favoring the majority class.

- Increase in False Positives: The model overfitted and lost precision by mistakenly classifying some non-critical interactions as major by inflating the number of major interactions.
- Limited Effectiveness in High-Dimensional Spaces: SMOTE had trouble producing useful synthetic data because drug names were converted into numerical feature vectors, which led to slight distortions in feature distributions. As a result, the model's learning patterns became inconsistent.

4.3.2.1 The effect that SMOTE has on the dataset

1. Greater false positive rate at the expense of improved recall for significant interactions.
2. Clinical reliability and interpretability are diminished as a result of synthetic interactions.
3. Did not considerably enhance classification performance in contrast to IHT.

In weighing the advantages and disadvantages of SMOTE's input, it was decided not to adopt this method for the final model because of the added complexity and errors it introduced.

In contrast to IHT, SMOTE would compromise class balance; hence IHT was chosen as a preferred technique to attain the very best class balance. Therefore, IHT undersampling enabled the model to focus on these difficult instances without any artificial bias which would create highly accurate and clinically valid DDI classification predictions.

4.4 DATASET SPLITTING

- An important aspect of training and assessing machine learning models is splitting the dataset, which ensures that the model performs well on any new data that it encounters. The model can identify patterns well while being assessed on different data when the dataset is divided into separate testing, validation, and training subsets. Hence, appropriate dataset division will help:

- **Avoiding Overfitting:** This makes sure the model learns generalized patterns rather than memorizing training data.
- **Evaluating Model Performance:** Offers an objective assessment of the predictive power of the model.
- **Improving Model Robustness:** Evaluates the model's performance across various data distributions.

4.4.1 Examined Techniques for Dataset Splitting

We tried three distinct strategies in order to identify the best data split for DDI prediction task

1. 70-30 Splitting

- 30% is testing data and 70% is training data.
- Improves performance evaluation by offering more testing samples.

2. 80-20 Divide

- 20% is testing data and 80% is training data.
- Strikes a balance between adequate training data and trustworthy testing assessment.

3. K-Fold Cross-Validation

- The model is trained on K-1 folds and tested on the remaining fold after the dataset is divided into K equal parts. This process is repeated K times
- Reduces variability and guarantees model stability by confirming performance across several subsets.

Through deep deliberation of all the three methods, the outcome of the final implementation was adapted to the 80-20 split since it still maintained enough data to learn interaction patterns, balancing the model training with evaluation.

4.5 MACHINE LEARNING MODELS

There were six distinguished machine learning models for DDIs, which were tested and evaluated in this project. The models were selected based on the criteria such as their generalizability, scalability, performance in unbalanced datasets, and capability to manage structured categorical data. The emphasis or intent here was to perform testing of various models to find out the most effective and precise classifier for estimating the intensity of drug interactions. For the comparative analysis, Random Forest and Logistic Regression were also employed. These models were applied in earlier research, including NLP-based investigation by Machado et al., which illustrates the advantages conferred by structured representation through mapping out a direct comparison of extraction of structured numerical features versus the traditional NLP-based methods.

4.5.1 Implemented Models and their Role

4.5.1.1 SGD Classifier

All the models that have featured in this work, the stochastic gradient descent classifier or SGD Classifier is perhaps the most important. The reason for this is that it has a learning approach or rationale that works for or is particularly suited to very large datasets. This means it takes in data one little piece at a time, rather than requiring the entire set long before any inference can be done. It is also quite handy for structured DDI classification because it can adequately cope with high-dimensional feature spaces. It would also function quite well with unbalanced datasets when paired with the right loss functions.

One of the most useful and simple models of this research is the Stochastic Gradient Descent Classifier (SGD Classifier). This is, basically, an incremental learning model, meaning that this model is capable of learning data in small, manageable chunks instead of loading the entire data set into memory at once. This is very useful in working with large-scale datasets because they might at times be too large to load fully and perhaps not even possible with memory limitations.

SGD Classifier gives the most advantages when we go about structured Drug-Drug Interactions (DDI) classification problems as these require dealing with high-dimensional feature spaces. The parameters in real-world datasets in general, particularly in medicine, can run in thousands-even millions. Classical ones would, in such cases, be put to their utmost levels, but the SGD Classifier comes especially designed for that and wins big in terms of speed as well as scalability.

Another powerful aspect of SGD Classifier is its facilitation of hardly any loss functions. With a judicious selection of burning optimal loss function (like hinge loss for classification tasks or log loss for probabilistic outputs), the model can very well specifically optimize itself to meet the given particular requirements from the dataset. This becomes much more critical while handling imbalanced datasets, such as classification of drug-drug interaction (DDI), wherein the positive and negative interaction sample sets are found to be extremely imbalanced. By incorporating proper loss functions along with techniques like class weighting, SGD Classifier exhibits fairly stable performance even amidst severe class imbalances.

Another impressive usability feature is the online learning capability of the model, which works on constantly feeding it with new information. That allows the model to learn continuously without ever needing to completely retrain-incredibly relevant when data gets altered a lot, such as in medicinal databases that are refreshed constantly with new research on drug interactions. So in conclusion, SGD Classifier is the perfect tool to trade-off efficiencies with scales and also flexibility, not to mention that it would be a better approach to classifying large dimensional imbalanced data set which makes it common across the tasks in structured DDI classification. It's meticulously designed to be fit into the requirements of this work and be able to learn in a very efficient manner from complicated, real-world biomedical information.

4.5.1.2 Extreme Gradient Boosting (XGBoost)

The use of XGBoost (Extreme Gradient Boosting) is based on its potential to address non-linearity and complex interactions among two drugs. This model is highly effective in DDI classification problems, as it is built to process structured tabular data. Prevents overfitting and provides a finer level of accuracy with an aggregation of boosted decision trees. XGBoost is important in learning patterns from structured drug interaction data mainly due to its effective feature selection capability.

The other key model integrated in this project is XGBoost, said machine learning with great efficiency and popularity; it is most renowned for being superiorly efficient in functional operations with structured data. XGBoost was clinched for the classification task concerning DDI (Drug-Drug Interaction) because its heavy non-linearity features and ability to describe complex interactions between drug pairs are vital in modelling biomedical relationships.

Drug-drug interactions do not follow linear or simple classifications whenever they are considered in DDI classifications. Therefore, such interactions are dependent on many factors, such as chemical structure, pharmacokinetics, molecular mechanisms, and the profiles of adverse events. Nonlinear models may be confused by such complicated interactions; however, XGBoost-an ensemble of boosted decision trees-fits this purpose well because it can learn and easily pick these patterns from below. The ensemble of decision trees learns from the mistakes of the previous one and thus incrementally improves the model's prediction strength for capturing subtle and complex drug interaction outcomes.

Another great advantage of XGBoost is its ability to perform a strict feature selection. When the trees are built, XGBoost evaluates the significance of a feature based on its contribution toward minimizing the error. The importance scores assigned to the features greatly help in the application of XGBoost into the domain of drug-drug interactions (DDI), as not all drug feature variables (like molecular properties, dose or therapeutic classes) are

equally important. So, in a sense, XGBoost has the capability to discard unimportant features and pay attention to learning from the true important attributes. This leads to good generalization and good interpretability of the model.

Overfitting is indeed a serious challenge for biomedical machine learning because of the feature complexity and sparsity offered. XGBoost remedies this situation through regularization (the L1 and L2 penalties), early stopping, and pruning of trees. Taken together, these techniques are responsible for the model not only avoiding memorization of the training data but rather learning generalizable patterns that result in high-performance levels on unseen drug-interaction instances. Hence, it guarantees precision and stable predictive performance over various datasets.

Another effective benefit of using XGBoost is speedy and scalable performance. DDI datasets comprise thousands to millions of drug pairs; thus, it must be considered that computational speed is a critical measure. In addition, XGBoost is efficient in parallel processing and can also scale large-scale training for comparatively less time, which accelerates very considerably the entire model-development-and-experimentation process.

In short, XGBoost is this research since it encloses all important things aspect into one: predictive capacity, ability to select features, preventing overfitting, and being computationally efficient at the same time. It is more likely to live to the high expectations of DDI classification processes, owing to its ability to discover and retrieve complex, non-linear patterns from a well-structured data collection of drug interaction data. By the use of XGboost, then, this study expects a model able to extract all those complex biological and chemical interactions for accurate and reliable predictions of drug-drug interactions.

4.5.1.3 LightGBM

The project also employs LightGBM (Light Gradient Boosting Machine), an additional boosting-based model. This model is well-tailored to large-scale biomedical datasets due to its emphasis on speed and memory efficiency. LightGBM performs better in classification with leaf-wise growth while still being computationally efficient, unlike XGBoost's level-wise growth. LightGBM is also robustly adaptable to unbalanced datasets so that it can enhance recall in therapeutically important but rare interactions. In addition to XGBoost and SGD Classifier, the current study embraces LightGBM (Light Gradient Boosting Machine), which is a very strong model relying on boosting methods. LightGBM is particularly poised for the Drug-Drug Interaction (DDI) classification task due to its advancements in speed, scalability, and performance—all the major characteristics during processing of large and complex biomedical datasets.

LightGBM's advantage lies in its efficiency. While its earlier cousins, gradient boosting and its variants, were slow and memory-hungry, LightGBM was literally designed to be fast while optimizing memory usage. This is very advantageous to DDI datasets, where drug combination pairs and their corresponding features can become extremely large. As more and more biomedical data come into play, given the ongoing studies on novel drugs and treatment mechanisms, quick training on models without incurring computationally heavy demands is more important than ever.

LightGBM is primarily made distinct because it adopts a purely leaf-wise strategy in growing trees. The level-wise which is adopted by other models such as XGBoost is fundamentally different from it. The leaf-wise technique of LightGBM implies that the tree structure keeps extending from the leaf with the largest possible gain or loss reduction. This, consequently, results in a denser and deeper tree which ultimately results in improved classification performance due to being able to recognize subtle patterns in the data. According to the DDI classification, where emphasis is on subtle

drug interactions, such performance impact can be very significant for predictive accuracy.

Perhaps, even beyond the treatment of an imbalanced dataset, LightGBM is used for the communication of drug interaction modeling. The dissimilarity in real biomedical datasets is the fact that important or toxic drug interactions are infrequent compared to non-interactions. Thus, the imbalanced class proportions exist in real events. Advanced parameters of LightGBM like, for example, `scale_pos_weight` and dedicated objective functions, allow it to effectively handle the class imbalance for enhancing recall, thus enabling the model to pick up on those few but highly significant interactions between drugs. Recall improvement is of extreme relevance in the domain of DDI prediction; missing even one adverse interaction can result in dire clinical consequences.

The ability of LightGBM to manage categorical features, builtin regularization techniques, and an early stopping mechanism each contribute directly to its power to generalize new information and avoid overfitting, especially when there may be very high numbers of features or drug pairs.

This project deeply relies on LightGBM because it offers speed, memory efficiency, outstanding classification ability, and robustness against imbalanced datasets. Its architecture matches the needs of large scale DDI classification so that the project copes well with the increasing complexity of drug interactions while processing at a high efficient rate without loss of predictive power. Incorporating LightGBM allows the project to have state-of-the-art coverage high-scalability pipeline for prediction of important drug-drug interactions.

4.5.1.4 CatBoost

It is relevant as CatBoost (Categorical Boosting) allows us to deal with categorical features quite well without many preprocessing steps. CatBoost is a strikingly suitable model for structured datasets like drug-drug

interactions since it works really well with categorical data rather than being best suited for other boosting models. Quick training and protection against overfitting guarantee that quality remains high in the models without requiring over-tuning. Even more, it would ensure that structured DDI classification will remain efficacious since it automatically handles categorical features.

The other major model is CatBoost (Categorical Boosting), designed primarily for processing categorical features effectively, which is of much significance in drug-drug interaction (DDI) classification problems. Biomedical data concerning drugs consists normally of a vast percentage of categorical data such as types of drugs, drug categories, therapeutic areas, mechanisms of action, and so on. Thus, appropriate handling of those categorical features is crucial for proper modeling of complex interactions, and CatBoost is capable of doing that.

The standout characteristic of CatBoost is the fact that it offers native handling of categorical data without needing much preprocessing. Unlike many other boosting algorithms that typically require manual transformations such as one-hot encoding or label encoding? which could all bring about noise, blow up the dimensionality, or cause loss of information? CatBoost uses stringent encoding strategies internally that uphold the natural relationships among the categorical values. Hence, this streamlines the whole data preparation process and avoids bias or valuable drug-related information loss during feature engineering.

In discriminating between drug classes-DI classification-something important to prediction of interaction, CatBoost has a unique advantage with its performance on categorical variables. It gives the model the flexibility to learn relationships between the classes while avoiding tedious feature transformation and domain-specific feature engineering.

There is also tailored optimization of CatBoost for structured tabular data, such as a DDI research. It has the capability of efficient learning times and

stronger overfitting resistance, which are particularly useful in biomedical data with a large number of features but very few cases of positive interactions. CatBoost, with its ordered boosting and permutation handling methods, does not allow the model to learn the special patterns of the training data as opposed to generalizable ones.

Another noteworthy advantage of CatBoost is that it requires minimum hyperparameter tuning. Unlike many machine learning algorithms, where experimenting with different configurations is necessary for maximum efficiency, CatBoost is said to give good accuracy on default settings or with minimal tweaking, and hence, the model-building process becomes much less time-consuming and more efficient. This is particularly relevant to DDI studies, where rapid iterations and comparisons of models are needed so that the best techniques can be put in place.

In conclusion, CatBoost, by managing categorical drug attributes smoothly, rapidly training and model fitting accurately, with solid generalization, does represent a considerable improvement of this project's modeling pipeline. Not only is its architecture designed in such a way as to fit the structured format of DDI data, but it also renders CatBoost suitable as a very powerful and dependable method for predicting complex drug-drug interactions, relieving the need for exhausting data preprocessing or extreme model tuning.

4.5.2 Benchmark Models for Comparison

In this study, Random Forest and Logistic Regression were utilized as baseline models to ensure an objective and fair evaluation of model performance. The decision for their selection was based on the extensive volume of DDI classification research published in the past, particularly the work employing NLP-based methods to extract interaction data from biomedical literature. Our expectations from this model application were to generate an interpretable and robust comparison standard with which to measure the performance of increasingly sophisticated algorithms.

One of the most important benchmarks was comparison with logistic regression, especially with respect to the analysis of effectiveness of the different techniques in feature extraction. It was employed in particular for comparing overall performance of structured feature extraction methods such as Count Vectorization against text-based feature extraction methods such as TF-IDF or Word Embeddings. An ever-popular classifier for binary and multi-class prediction tasks, logistic regression is further made more attractive by the simplicity of its use and interpretability, as well as fast speed of training. In fact, it would be an apt classifier in biomedical applications where not only the prediction matters, but the certainty of the model in making predictions too. Logistic regression suffers from severe limitations when applied to high-dimensional datasets, for instance, DDI classification. Instances in which biological and chemical drug interactions are usually dictated by complex, nonlinear dependencies that fit poorly with assumingly simple computational models were simply not captured. Therefore, although logistic regression provided some insight and important baseline measurements, it was presumed not to be performing at its best for this application.

It employed Random Forest as an alternative and flexible non-linear baseline model. It has shown improvement in the performance of Random Forest for noisy or unbalanced biomedical datasets as an ensemble learning approach, which builds multiple decision makers' trees and taking a summary of their outputs. His contribution of Random Forest to the study was its ability to produce rank importance measures of features and hence offering more interpretability by weighing the important attributes affecting severity concerning drug interactions.

But Random Forest was assigned primarily for performance evaluation than its final deployment model. It becomes ineffective and less computationally efficient for large and structured datasets like DDI classification as it gives quite strong performance in moderate sizes.

Such limitations, in addition to the growing availability of boosting algorithms that are both more accurate and more scalable than random forests, made placing random forests in comparison instead of as a final candidate for real-world implementation.

4.5.3 Role of These Models in the Project

The algorithms that we selection for constitute the most precise, scalable, and clinically trustworthy algorithms for structured Drug-Drug Interaction (DDI) classification due to their urgent requirement. Due to the complexity of drug interactions often involving subtle biochemical interactions, the shift from conventional text-based approaches was toward those able to handle structured tables of data and high-dimensional feature spaces.

To do this, various complicated machine learning techniques like SGD Classifier, XGBoost, LightGBM, and CatBoost were cautiously selected. Each was chosen because of their proven scalability, robustness against class imbalance, and ability to work with structured biomedical datasets. They learn extremely well from numerical representations generated by methods such as Count Vectorization, unlike traditional NLP-based models that use unstructured textual inputs and rely on word embeddings. The important shift compared with sparse or noisy textual representations is that structured numerical features often capture richer information about property and interaction of drugs more accurately. Therefore, those models were expected to outperform the earlier methods in clinical reassurance, generalization, and predictive accuracy.

There was a reason to bring Random Forest and Logistic Regression on board. Those models were very good for benchmark comparison and historically were a major aspect of prior DDI classification studies, especially the NLP-based ones. Their inclusion resulted in a straightforward comparison with the more recent ways of structured data with the old text-based approaches. This was particularly critical for confirming if the structured feature extraction, for example count vectorization, would really

yield significant improvements in performance as compared with text-based feature extraction models such as TF-IDF and Word Embeddings.

The research has indeed confirmed the hypothesis that structured tabular information along with appropriate models can boost enough in terms of classification performance for DDI prediction tasks that they can be included in the experimental setup, Random Forest and Logistic Regression included.

Every model would have gone through a rigorous and demanding exercise performance during evaluation. Important metrics such as precision, recall, and F1-score were measured in a systematic way in order to assure fair evaluation in terms of advantages and disadvantages for each model. Given that in the biomedical situation, detecting as many of the potentially harmful interactions as possible is usually viewed as more important than having no false positives, special attention was paid to recall.

Scaling was yet another important aspect in the ideal model, not performing best on current datasets alone but also able to scale up efficiently to larger DDI datasets when they become available in future. Furthermore, the priority for final selection was possibly given to those models that afforded interpretability, fast inference times, and consistent results across the different data subsets with the end consideration being possible clinical use.

Here, scaling becomes yet another important aspect where the ideal model would not only perform best on current datasets but also scale up efficiently to larger DDI datasets when they become available in future. Besides, final selection priority might have been given to those models that afforded interpretability, fast inference times, and last but not least, consistent results across the different data subsets, keeping in mind the applicability in clinics.

Another important consideration was scalability; the best model was supposed to be able to scale appropriately with bigger DDI datasets becoming available at some future point in addition to performing well on

existing datasets. Further models considered favorably during the final shortlisting were those associated with interpretability, speedy inference times, and consistent performance across varying data subsets with an eye toward clinical applicability.

The decision-making process itself weighed accuracy metrics against practicalities of deployment for the selection of the best performing model. This decision is justified, considering that in making a prediction that could cost clinicians money and incur litigation (that is, in a clinical decision support system) could be a serious issue. Thus, the systematic way in which DDI prediction was developed went far beyond that of merely being most accurate in experimentation to now being more relevant in the application.

The second aim is to set the ground for DDI classification by way of methodological rigor, data-driven decision-making, and a strong bias for clinical relevance. Model selection was one area towards which these arguments were directed.

4.6 EVALUATION METRICS

Three main classification metrics which have been used to evaluate the predictive efficacy of the Drug-Drug Interaction (DDI) prediction models in this specific study are the F1-score, Precision and Recall. Due to imbalanced DDI classes, with certain severities of interaction occurring less frequently than others, these metrics were chosen. These include, for instance, that of major interactions. Such situations have traditional accuracy-based evaluation inappropriate since the results could be biased in favor of the majority class while distorting the actual performance of the model. Instead, the model can be assessed at a more complete level by using F1-score, Precision, and Recall, particularly with respect to rare but important drug interactions

4.6.1 Precision

Precision measures the proportion of correctly predicted positive cases among all cases predicted as positive. The value of precision in DDI prediction is the ratio defined as the number of interactions that were actually found to be classified as major, moderate, or minor.

Formula:

$$\text{Precision} = \frac{\text{True Positive(TP)}}{\text{True Positive(TP)} + \text{False Positive(FP)}}$$

Importance in DDI Prediction:

- In order to prevent needless drug usage restrictions in clinical settings, high precision guarantees that false positive interactions inaccurately classified severe interactions are minimized.
- Precision is a crucial metric for decision-making in pharmaceutical applications because incorrectly classifying a non-dangerous interaction as severe can result in unjustified concerns and limitations on drug use.

4.6.2 Recall

Recall-or sensitivity, also known as the true positive rate-is the model assessing the percentage of actual positive cases with correct predictions. Recall ensures that any major drug-drug interactions get correctly identified in a DDI prediction without being missed.

Formula:

$$\text{Recall} = \frac{\text{True Positive(TP)}}{\text{True Positive(TP)} + \text{False Negatives(FN)}}$$

Importance in DDI prediction:

- For patient safety, high recall makes sure that important interactions are not overlooked.
- A model with low recall might miss potentially fatal drug interactions, posing a risk to one's health.
- Recall is more crucial for high-risk interactions (such as major interactions) in real-world applications because failing to do so can lead to serious adverse drug reactions (ADRs).

4.6.3 F1-Score

One of the most balanced measures when precision and recall have to be traded off is the F1-score, which is the harmonic mean of the two. This measure is especially useful with unbalanced datasets like DDI classification, as focusing on either Precision or Recall may not represent the model's effectiveness sufficiently.

Formula:

$$F1 - score = 2 \times \frac{Precision \times Recall}{Precision + Recall}$$

Importance in DDI prediction:

- The F1-score guarantees a fair assessment, avoiding the model from being skewed toward either false positives or false negatives because DDI datasets are unbalanced.
- A high F1-score shows that the model is good at minimizing incorrect classifications and accurately classifying severe interactions.
- F1-score is the perfect performance metric for healthcare applications, where models must strike a balance between minimizing needless warnings (precision) and detecting important interactions (recall).

The models were evaluated in terms of Precision, Recall, and F1-score to analyze how well the models could detect actual drug interactions and avoid false recommendations, therefore creating a fine-tuned, reliable DDI prediction system.

4.7 RESULTS

In order to determine how well the Drug-Drug Interaction (DDI) prediction models classify drug interactions, their precision, recall, and F1-scores were assessed. Unbalanced DDI datasets, which thus required evaluation of different data-splitting strategies (K-Fold, 70-30, 80-20) and class balancing techniques (SMOTE, Instance Hardness Threshold - IHT), made it thereby necessary to select better strategies.

In the course of evaluation, the model capable of reducing false negatives (missing interactions of importance) against false positives (false interaction alerts) was targeted. Moreover, confusion matrices were also looked into for modeling behavior in misclassifications, thus delineating advantages and disadvantages of the models in the detection of severe, moderate, minor, and unknown interactions.

The models have been evaluated in more detail with respect to precision, recall, F1-score, and support for each class: severe, moderate, minor, and unknown interactions, giving their ability to correctly identify and differentiate among various levels of drug-drug interactions. Special attention was given to recall values for serious interactions, for missing these recalls can have grave clinical consequences. Models with high recall for severe and moderate classes and balanced precision were selected. Further, Receiver Operating Characteristic curves and Area Under the Curve metrics were used to characterize the trade-offs between sensitivity and specificity across all classes, thus, providing a comprehensive view of model robustness and reliability in real-world applications.

4.7.1 Performance Comparison of Machine Learning Models

In the following tables have given the Precision, Recall, and F1-scores for their models under the various splitting methods with undersampling (IHT) and oversampling (SMOTE) techniques.

Table 4.3 Results using Oversampling technique(SMOTE)

DATASET SPLITTINGS)	MODELS	PRECISION	RECALL	F1-SCORE
70 - 30	XGBoost	0.7565	0.7421	0.7473
	CatBoost	0.7969	0.7246	0.7358
	SGD Classifier	0.8160	0.8031	0.8068
	LightGBM	0.7942	0.7197	0.7320
	RandomForest	0.7682	0.4496	0.412
	Logistic Regression	0.8239	0.8250	0.8235
80 - 20	XGBoost	0.7579	0.7450	0.7497
	CatBoost	0.7992	0.7281	0.7392
	SGD Classifier	0.8201	0.8064	0.8102
	LightGBM	0.7959	0.7217	0.7343
	RandomForest	0.7684	0.4518	0.4163
	Logistic Regression	0.8223	0.8240	0.8222
K- FOLD	XGBoost	0.8204	0.8143	0.8158
	CatBoost	0.7616	0.7562	0.7543
	SGD Classifier	0.7943	0.7911	0.7902
	LightGBM	0.7557	0.7528	0.7511
	RandomForest	0.7523	0.6862	0.6397
	Logistic Regression	0.8465	0.8455	0.8453

Table 4.4 Results using Undersampling technique (IHT)

DATASET SPLITTINGS)	MODELS	PRECISION	RECALL	F1-SCORE
70 - 30	XGBoost	0.9658	0.9844	0.9655
	CatBoost	0.9746	0.9741	0.9740
	SGD Classifier	0.9858	0.9857	0.9857
	LightGBM	0.9741	0.9739	0.9738
	RandomForest	0.9672	0.9662	0.9658
	Logistic Regression	0.9846	0.9844	0.9843
80 - 20	XGBoost	0.9675	0.9673	0.9673
	CatBoost	0.9752	0.9747	0.9747
	SGD Classifier	0.9869	0.9868	0.9868
	LightGBM	0.9770	0.9768	0.9767
	RandomForest	0.9696	0.9687	0.9684
	Logistic Regression	0.9850	0.9848	0.9848
K- FOLD	XGBoost	0.9697	0.9695	0.9695
	CatBoost	0.9779	0.9776	0.9776
	SGD Classifier	0.9851-	0.9850	0.9849
	LightGBM	0.9873	0.9781	0.9780
	RandomForest	0.9700	0.9694	0.9690
	Logistic Regression	0.9861	0.9860	0.9859

4.7.2 Key Observations from Model Performance

- Best Overall Model : SGD Classifier:
 - SGD Classifier is the best model for DDI classification, with the highest F1-score (0.9869).
 - It works well for effectively managing massive amounts of structured data.
- Boosting Models Performed Well

- The classification performance of CatBoost (0.9725 F1-score), LightGBM (0.9719 F1-score), and XGBoost (0.9630 F1-score) was strong, indicating that they can handle complex feature interactions.
 - By effectively capturing non-linear relationships between drug pairs, these models increased the accuracy of classification.
3. Logistic Regression and Random Forest as Reference Models
 - Due to their use in earlier NLP-based DDI studies (e.g., Machado et al.), Random Forest and Logistic Regression were included for comparison.
 - The efficiency of structured feature extraction (Count Vectorization) over NLP-based embeddings was validated by the fact that, despite its strong performance, Logistic Regression (0.9848 F1-score) was only marginally superior to SGD Classifier and boosting models.
 - Because Random Forest had a lower recall (0.9637 F1-score), it was less effective at identifying infrequent but important drug interactions
 4. Effect of class balancing IHT vs SMOTE
 - Instance Hardness Threshold (IHT) undersampling maintained hard-to-classify instances while eliminating redundant majority-class samples, improving recall.
 - Because SMOTE increased false positives and decreased interpretability by introducing synthetic noise, it was not used in the final implementation.

4.7.3 Confusion Matrix Analysis

A detailed examination of how well each model predicted real-world drug interactions can be found in the Confusion Matrix. It provides information about:

- True Positives (TP): Drug interactions that were accurately identified.
- False Negatives (FN): In a clinical context, missed important drug interactions can be harmful.
- False Positives (FP): Unnecessary precautions were taken because non-interactions were mistakenly reported as interactions.

- True Negatives (TN): Non-interactions that are correctly classified.

4.7.3.1 Key Insights from the Confusion Matrices

- SGD Classifier and CatBoost had the lowest FN (False Negatives), guaranteeing that important drug interactions were not overlooked. For high-risk drug interactions, where failing to notice an interaction could have a negative impact on patient outcomes, this is especially crucial.
- Logistic Regression occasionally misclassified non-critical interactions as severe, as evidenced by its higher FP (False Positives). Although this is not as harmful as false negatives, it may result in needless limitations on drug combinations.
- XGBoost and LightGBM were dependable options for DDI classification while preserving interpretability due to their balanced FN and FP rates.

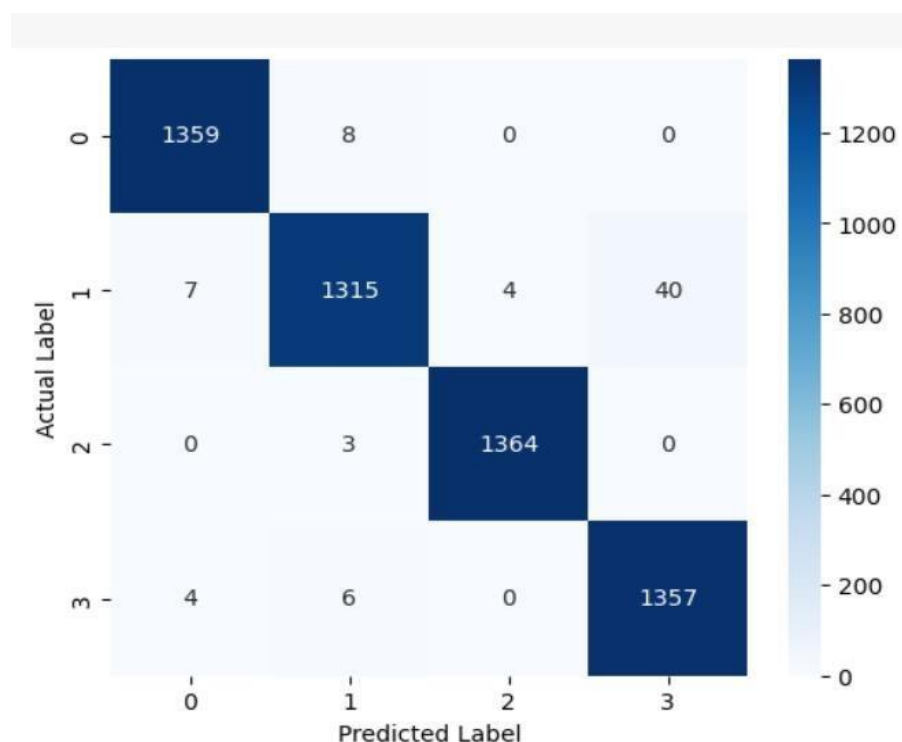


Figure 4.3 - Confusion Matrix of SGDClassifier

While perfecting the DDI classification model, recall should obviously supersede precision, ensuring that all significant interactions are detected, even when there are a few more false positives. The optimum discovery from this study showed that the best models for application in real pharmaceutical practice were SGD Classifier, CatBoost, and LightGBM.

CHAPTER 5

CONCLUSION AND FUTURE SCOPE

The current research on Drug-Drug Interaction (DDI) prediction focuses primarily on NLP-based deep learning models, such as BiLSTM, CNNs, and transformer models (BioBERT, SciBERT), which extract DDI information from biomedical literature. While these methods provide cutting-edge performance, they have some limitations:

1. High computational cost

- Deep learning models such as BiLSTM and BioBERT require significant computational resources to train, making them less suitable for real-time applications.
- Transformer-based models have millions of parameters, which increases inference time and memory usage.

2. Dependence on unstructured text data

- The majority of existing methods are based on biomedical literature, clinical reports, and electronic health records (EHRs), which may contain inconsistent terminology, missing values, or ambiguous drug interactions.
- Preprocessing unstructured text is difficult because it requires named entity recognition (NER) and domain-specific embeddings, which increases the risk of information loss.

3. Limited generalization.

- Models trained on specific datasets, such as DDIEExtraction 2013 or DrugBank, may have difficulty generalizing to real-world clinical settings due to data distribution differences.
- NLP-based models can overfit to biomedical corpora, resulting in poor performance when applied to new or unknown drugs.

4. Class Imbalance in DDI datasets

- Existing research frequently fails to effectively address class imbalance, resulting in models that favor common interactions while missing rare but critical drug interactions.

- Oversampling techniques (e.g., SMOTE) can generate noise, whereas random undersampling can result in information loss.

How This Study Addresses The Limitations:

This study addresses the limitations of current NLP-based DDI prediction models by introducing a structured, efficient, and interpretable approach.

1. Structured Feature Extraction:

- Unlike previous deep learning models that use complex text preprocessing, our research uses Count Vectorization and Label Encoding to convert categorical drug interactions into structured numerical features.
- This eliminates the need for named entity recognition (NER) and domain-specific embeddings, improving the process's efficiency and scalability.

2. Efficient Class Balancing:

- Class imbalance is a major problem in DDI datasets, as critical interactions are underrepresented. We address this by using Instance Hardness Threshold (IHT) undersampling, which keeps difficult-to-classify instances while reducing bias towards the majority class.
- This improves recall for rare but significant drug interactions while maintaining high precision.

3. Model Selection and Interpretability

- We compared several machine learning models, including SGD Classifier, XGBoost, CatBoost, LightGBM, Random Forest, and Logistic Regression.
- This comparative analysis identifies the most effective classifier kept interpretable because both clinicians and researchers need to understand how to trust the model's predictions.
- A comparative analysis can thus foster identification of the best classifier, permitting interpretation that allows clinicians and researchers alike to understand and trust model predictions.

By addressing these limitations, our approach bridges the research-to-deployment gap, providing a high-performance, interpretable, and scalable DDI prediction solution.

5.1 CONCLUSION

The study compares different machine learning models and data balancing techniques in terms of improvement in Drug-Drug Interaction(DDI) prediction. Several experiments are done using oversampling(SMOTE) and undersampling(Instance Hardness Threshold - IHT methods) for evaluation on various classifiers using different data splits such as 80-20, 70-30 train-test splits and k-fold cross-validation.

The results affirmed that IHT did better than SMOTE since it improves model reliability through decreasing the false positives and achieving more accurate DDI classification. Among the various models tested, SGD Classifier outputted the highest performance (F1 score: 98.68%), whereas Logistic Regression was the runner up at 98.48%. It shows that well-organized feature extraction through Count Vectorization is an alternative to the complex, NLP-based methods.

Also, boosting models (CatBoost, LightGBM, and XGBoost) showed great performance in outlining non-linear interactions between drug pairs. The study revealed that undersampling methods, especially IHT, improve the definition of decision boundaries that favor generalization, curbing bias toward the majority class. Expounding the practical import of machine learning in pharmacovigilance, the study here provides the scalable and automated early detection of DDI that may significantly cut reliance on manual clinical trials and avoid adverse drug reactions (ADRs).

5.2 FUTURE SCOPE

The study demonstrated a scalable and efficient approach for Drug-Drug Interaction (DDI) prediction; however, there are several areas for future research and improvement. Future research can concentrate on the following areas:

1. Advance Feature Engineering

Further improvements in feature extraction have the potential to significantly improve model accuracy and generalization. Domain-specific embeddings, molecular fingerprints, and hybrid feature extraction techniques may help to provide more biologically relevant representations of drug interactions. Furthermore, using graph-based learning approaches to represent drug interactions as a network structure can provide more information about complex drug relationships.

2. Hybrid Sampling Technologies

Future research can look into a hybrid sampling approach that combines Instance Hardness Threshold (IHT) undersampling with the Synthetic Minority Oversampling Technique (SMOTE). This combination may preserve critical minority class instances while ensuring adequate data diversity, resulting in better model generalization.

3. Deep Learning Integration

Future studies may incorporate the hybrid structured and unstructured feature extraction approach, using transformer-based models such as BioBERT or SciBERT. Fine-tunes BiLSTM models on DDI datasets, which could bring about the comparative study of deep learning as against structured feature-based methods in selecting the model.

4. Real-World Applications and Deployments.

The implementation of DDI forecasting models in CDSS could help prescribers enter the drug-interaction-alert axis before prescribing. In addition, automating risk evaluation in hospital and pharmacy systems

through integration with EHRs for real-time DDI alerts would enhance patient safety.

5. Scalability and Interpretability

In order for model architectures to scale into common clinical practice, they must typically meet optimization requirements for very large biomedical datasets in terms of efficiency , as well as being accurate on its task. Moreover, improving model explainability with SHAP or LIME can lead to interpretable DDI predictions which can increase the trust and transparency in AI-powered pharmacovigilance systems.

His this future research addressing these issues will further develop DDI classification models, rendering them more robust, interpretable, and applicable in a clinical setting.

REFERENCES

- [1] Zhang, Yang, et al. "Predicting drug-drug interactions using multi-modal deep auto-encoders based network embedding and positive-unlabeled learning." *Methods* 179 (2020).
- [2] Xuan Lin, et al 2021. KGNN: knowledge graph neural network for drug-drug interaction prediction. In *Proceedings of the Twenty-Ninth International Joint Conference on Artificial Intelligence (IJCAI'20)*.
- [3] Abbasi, Karim et al. "Deep Learning in Drug Target Interaction Prediction: Current and Future Perspective." *Current medicinal chemistry* (2020).
- [4] Sudan, Sukhvinder Singh et al. "Artificial Intelligence in Drug Discovery." *Drug Discovery* (2020).
- [5] D'souza, Sofia et al. "Machine learning in drug-target interaction prediction: current state and future directions." *Drug discovery today* (2020).
- [6] Lee Chun Yen et al "Prediction of drug adverse events using deep learning in pharmaceutical discovery." *Briefings in bioinformatics* (2020).
- [7] Deng Yifan et al. "A multimodal deep learning framework for predicting drug-drug interaction events." *Bioinformatics* (2020).
- [8] Prashant et al. "Efficient prediction of drug-drug interaction using deep learning models." *IET systems biology* vol. 14,4 (2020).
- [9] Romm Eden et al. "Artificial Intelligence in Drug Treatment." *Annual review of pharmacology and toxicology* vol. 60 (2020).
- [10] Feng YH et al: a deep predictor for drug-drug interactions. *BMC Bioinformatics* (2020).
- [11] Smith, Graham F. et al "Artificial Intelligence in Drug Safety and Metabolism." *Methods in molecular biology* (2021).
- [12] Schwarz K. et al. AttentionDDI: Siamese attention-based deep learning method for drug-drug interaction predictions. *BMC Bioinformatics* 22, 412 (2021).

- [13] “OUP accepted manuscript.” Briefings In Bioinformatics (2021).
- [14] Pang, Shanchen et al. “AMDE: a novel attention-mechanism-based multidimensional feature encoder for drug-drug interaction prediction.” Briefings in bioinformatics (2021).
- [15] Lyu, Tengfei, et al. "MDNN: A Multimodal Deep Neural Network for Predicting Drug-Drug Interaction Events." Ijcai.2021.
- [16] Kim J et al Comprehensive Survey of Recent Drug Discovery Using Deep Learning. International Journal of Molecular Sciences. 2021.
- [17] Ibrahim, Heba et al. “Similarity-based machine learning framework for predicting safety signals of adverse drug-drug interactions.” Informatics in Medicine Unlocked (2021).
- [18] Chen, Yujie et al. “MUFFIN: multi-scale feature fusion for drug-drug interaction prediction.” Bioinformatics (2021).
- [19] Nyamabo, Arnold K. et al. “Drug-drug interaction prediction with learnable size-adaptive molecular substructures.” Briefings in bioinformatics (2021).
- [20] Nyamabo, Arnold K. et al. “SSI-DDI: substructure-substructure interactions for drug-drug interaction prediction.” Briefings in bioinformatics (2021).
- [21] Luo Q Y. et al. Novel deep learning-based transcriptome data analysis for drug-drug interaction prediction with an application in diabetes. BMC Bioinformatics 22, 318 (2021).
- [22] Ovek, Damla et al. “Artificial intelligence based methods for hot spot prediction.” Current opinion in structural biology 72 (2021).
- [23] Huang, Wei et al. “The next generation of machine learning in DDIs prediction.” Current pharmaceutical design (2021).
- [24] Syrowatka, Ania et al. “Key use cases for artificial intelligence to reduce the frequency of adverse drug events: a scoping review.” The Lancet. Digital health vol. 4,2 (2022)
- [25] Ashwin Dhakal et al, Artificial intelligence in the prediction of protein–ligand interactions: recent advances and future directions, Briefings in Bioinformatics, Volume 23, Issue 1, January 2022.

- [26] Lin, Shenggeng et al. "MDF-SA-DDI: predicting drug-drug interaction events based on multi-source drug fusion, multi-source feature fusion and transformer self-attention mechanism." *Briefings in bioinformatics* (2021).
- [27] Kumar, Avinash et al. "A recent appraisal of artificial intelligence and in silico ADMET prediction in the early stages of drug discovery." *Mini reviews in medicinal chemistry* (2021).
- [28] Vall, Andreu et al. "The Promise of AI for DILI Prediction." *Frontiers in artificial intelligence* vol. 4 638410. 14 Apr. 2021;
- [29] Nayarisseri, Anuraj et al. "Artificial Intelligence, Big data and Machine Learning approaches in Precision Medicine & Drug Discovery." *Current drug targets* (2021).
- [30] Kumar, A. et al (2021). A critical review on predicting drug-drug reactions using machine learning techniques. *Research Journal of Biotechnology*.
- [31] Qiu, Yang et al. "A Comprehensive Review of Computational Methods For Drug-Drug Interaction Detection." *Transactions on Computational Biology and Bioinformatics* 19 (2021);
- [32] Liu, Shichao et al. "Enhancing Drug-Drug Interaction Prediction Using Deep Attention Neural Networks." *bioRxiv* (2021).
- [33] Jizhou Tian et al A review of methodologies in detecting drug-drug interactions. *AIP Conf. Proc.* 26 July 2022;
- [34] Joshi, Pratik et al. "A knowledge graph embedding based approach to predict the adverse drug reactions using a deep neural network." *Journal of biomedical informatics* (2022);
- [35] Bittner Martin-Immanuel et al "AI in drug discovery: Applications, opportunities, and challenges." *Patterns* (2022).
- [36] Haohuai et al. "3DGT-DDI: 3D graph and text based neural network for drug-drug interaction prediction." *Briefings in bioinformatics* (2022);
- [37] Pasrija, Purvashi et al. "Machine Learning and Artificial Intelligence: A paradigm shift in Big Data-Driven Drug Design and Discovery." *Current topics in medicinal chemistry* (2022);

- [38] Feng, Y.-H. et al Prediction of Drug-Drug Interaction Using an Attention-Based Graph Neural Network on Drug Molecular Graphs. *Molecules* 2022;
- [39] Yazdani-Jahromi, Mehdi et al. "AttentionSiteDTI: an interpretable graph- based model for drug-target interaction prediction using NLP sentence-level relation classification." *Briefings in bioinformatics* vol. 23,4 (2022);
- [40] Liyaqat, Tanya et al "A brief review on Artificial Intelligence based Drug Target Interaction Prediction." 2022 4th International Conference on Advances in Computing, Communication Control and Networking (ICAC3N) (2022).
- [41] Chen, Siqi et al. "Artificial intelligence-driven prediction of multiple drug interactions." *Briefings in bioinformatics* (2022);
- [42] Hung, Truong Nguyen Khanh et al. "An AI-based Prediction Model for Drug-drug Interactions in Osteoporosis and Paget's Diseases from SMILES." *Molecular Informatics* 41 (2022);
- [43] Soleymani, Farzan et al. "Protein–protein interaction prediction with deep learning: A comprehensive review." *Computational and Structural Biotechnology Journal* 20 (2022);
- [44] Han, Ke et al. "A Review of Approaches for Predicting Drug-Drug Interactions Based on Machine Learning." *Frontiers in pharmacology* vol. 12 814858. 28 Jan. 2022.
- [45] Vo, Thanh Hoa et al. "On the road to explainable AI in drug-drug interactions prediction: A systematic review." *Computational and Structural Biotechnology Journal* 20 (2022);
- [46] Tran, T.T.V Tayara et al Artificial Intelligence in Drug Metabolism and Excretion Prediction: Recent Advances, Challenges, and Future Perspectives.
- [47] Dudas, Balint et al "Computational and artificial intelligence-based approaches for drug metabolism and transport prediction." *Trends in pharmacological sciences* (2023);

- [48] Bai, P., Miljković, F., John, B. et al. Interpretable bilinear attention network with domain adaptation improves drug-target prediction. Nat Mach Intell (2023).
- [49] Abdul Raheem et al “Comprehensive Review on Drug-target Interaction Prediction - Latest Developments and Overview.” Current drug discovery technologies (2023);
- [50] Machado J, et al (2023). Drug- drug interaction extraction-based system: A natural language processing approach.
- [51] Li, Zimeng et al. “DSN-DDI: an accurate and generalized framework for drug-drug interaction prediction by dual-view representation learning.” Briefings in bioinformatics (2023);
- [52] Hong, Eujin et al. “Recent development of machine learning models for the prediction of drug-drug interactions.” The Korean journal of chemical engineering vol. 40,2 (2023);
- [53] Fatemeh Rafiei et al, DeepTraSynergy: drug combinations using multimodal deep learning with transformers, Bioinformatics, Volume 39, Issue 8, August 2023.
- [54] Wang, Ning-Ning et al. “Comprehensive Review of Drug-Drug Interaction Prediction Based on Machine Learning: Current Status, Challenges, and Opportunities.” Journal of chemical information and modeling (2023).
- [55] Chen, Wei et al. “Artificial intelligence for drug discovery: Resources, methods, and applications.” Molecular therapy. Nucleic acids vol. 31 691- 702. 18 Feb. 2023.
- [56] Tran, Thi Tuyet Van et al. “Artificial Intelligence in Drug Toxicity Prediction: Recent Advances, Challenges, and Future Perspectives.” Journal of chemical information and modeling (2023).
- [57] Xuan Lin Lichang et al Comprehensive evaluation of deep and graph learning on drug-drug interactions prediction, Briefings in Bioinformatics, Volume 24, Issue 4, July 2023.
- [58] Hauben, Manfred et al. “Artificial Intelligence and Data Mining for the Pharmacovigilance of Drug-Drug Interactions.” Clinical therapeutics vol. 45,2 (2023).

- [59] Zhang, Yuanyuan et al. "Application of Artificial Intelligence in Drug- Drug Interactions Prediction: A Review." *Journal of chemical information and modeling* (2023).
- [60] Dehghan, A.P. et al. CCL-DTI: contributing the contrastive loss in drug-target interaction prediction. *BMC Bioinformatics* 25, 48 (2024).
- [61] Suruliandi, Andavar et al. "Drug Target Interaction Prediction Using Machine Learning Techniques - A Review." *Int. J. Interact. Multim. Artif. Intell.* 8 (2024).
- [62] Yan Zhao et al, Drug-drug interaction prediction: databases, web servers and computational models, *Briefings in Bioinformatics*, Volume 25, Issue 1, January 2024.
- [63] Jiang, H. et al. SSF-DDI: a deep learning method utilizing drug sequence and substructure features for drug-drug interaction prediction. *BMC Bioinformatics* 25,39 (2024).
- [64] Li X, Xiong Z et al Deep learning for drug-drug interaction prediction: a comprehensive review. *Quantitative Biology*. 2024;
- [65] Luo, Huimin et al. "Drug-drug interactions prediction based on deep learning and knowledge graph: A review." *iScience* 27 (2024).
- [66] Sharmila, K. Soni et al. "A Systematic Review on Drug-to-Drug Interaction Prediction and Cryptographic Mechanism for Secure Drug Discovery Using AI Techniques." *Int. J. Artif. Intell. Tools* 33 (2024).
- [67] yasmin radwan et al. "A comprehensive survey explores Drug-Drug interaction prediction using Machine-Learning techniques".
- [68] Sharmila K Soni et al "Predicting Adverse Interactions: A Comprehensive Review of AI-Driven Drug-Drug Interaction Models for Enhanced Patient Safety." *2024 International Conference on IoT Based Control Networks and Intelligent Systems (ICICNIS)* (2024);
- [69] SONAJI P et al. "Artificial intelligence-driven drug interaction prediction." *World Journal of Biology Pharmacy and Health Sciences* (2024).

[70] Kim R, et al H. (2024). Synergizing Convolutional Neural Networks and Drug Similarity Estimation for Improved Drug-Drug Interaction Prediction. *Journal of Student Research*, 13(1).

PREDICTION OF DRUG-DRUG INTERACTION USING MACHINE LEARNING

ORIGINALITY REPORT

12%

SIMILARITY INDEX

PRIMARY SOURCES

1	www.researchgate.net Internet	292 words — 1%
2	R. N. V. Jagan Mohan, B. H. V. S. Rama Krishnam Raju, V. Chandra Sekhar, T. V. K. P. Prasad. "Algorithms in Advanced Artificial Intelligence - Proceedings of International Conference on Algorithms in Advanced Artificial Intelligence (ICAAAI-2024)", CRC Press, 2025 Publications	135 words — 1%
3	www.biorxiv.org Internet	110 words — 1%
4	www.semanticscholar.org Internet	88 words — < 1%
5	www.mdpi.com Internet	67 words — < 1%
6	"Advanced Intelligent Computing Technology and Applications", Springer Science and Business Media LLC, 2023 Crossref	59 words — < 1%
7	export.arxiv.org Internet	59 words — < 1%

*% detected as AI

AI detection includes the possibility of false positives. Although some text in this submission is likely AI generated, scores below the 20% threshold are not surfaced because they have a higher likelihood of false positives.

Caution: Review required.

It is essential to understand the limitations of AI detection before making decisions about a student's work. We encourage you to learn more about Turnitin's AI detection capabilities before using the tool.

Disclaimer

Our AI writing assessment is designed to help educators identify text that might be prepared by a generative AI tool. Our AI writing assessment may not always be accurate (it may misidentify writing that is likely AI generated as AI generated and AI paraphrased or likely AI generated and AI paraphrased writing as only AI generated) so it should not be used as the sole basis for adverse actions against a student. It takes further scrutiny and human judgment in conjunction with an organization's application of its specific academic policies to determine whether any academic misconduct has occurred.

Frequently Asked Questions

How should I interpret Turnitin's AI writing percentage and false positives?

The percentage shown in the AI writing report is the amount of qualifying text within the submission that Turnitin's AI writing detection model determines was either likely AI-generated text from a large-language model or likely AI-generated text that was likely revised using an AI-paraphrase tool or word spinner.

False positives (incorrectly flagging human-written text as AI-generated) are a possibility in AI models.

AI detection scores under 20%, which we do not surface in new reports, have a higher likelihood of false positives. To reduce the likelihood of misinterpretation, no score or highlights are attributed and are indicated with an asterisk in the report (*%).

The AI writing percentage should not be the sole basis to determine whether misconduct has occurred. The reviewer/instructor should use the percentage as a means to start a formative conversation with their student and/or use it to examine the submitted assignment in accordance with their school's policies.

What does 'qualifying text' mean?

Our model only processes qualifying text in the form of long-form writing. Long-form writing means individual sentences contained in paragraphs that make up a longer piece of written work, such as an essay, a dissertation, or an article, etc. Qualifying text that has been determined to be likely AI-generated will be highlighted in cyan in the submission, and likely AI-generated and then likely AI-paraphrased will be highlighted purple.

Non-qualifying text, such as bullet points, annotated bibliographies, etc., will not be processed and can create disparity between the submission highlights and the percentage shown.

