



Dissertation on

**“Importance of Drug Features in Drug-Drug Interaction:
A Comparative Study”**

Submitted in partial fulfillment of the requirements for the award of degree of

**Bachelor of Technology
in
Computer Science & Engineering**

UE20CS461A – Capstone Project Phase - 2

Submitted by:

**Chetana N Patil
Naveen B N
Mekala Sanjana
Kumari Shivangi**

**PES2UG20CS504
PES2UG20CS217
PES2UG20CS194
PES2UG20CS173**

Under the guidance of

Dr.Prajwala.T .R
Associate Professor
PES University

June - Nov 2023

**DEPARTMENT OF COMPUTER SCIENCE AND ENGINEERING
FACULTY OF ENGINEERING
PES UNIVERSITY**

(Established under Karnataka Act No. 16 of 2013)
Electronic City, Hosur Road, Bengaluru – 560 100, Karnataka, India



PES UNIVERSITY

(Established under Karnataka Act No. 16 of 2013)
Electronic City, Hosur Road, Bengaluru – 560 100, Karnataka, India

FACULTY OF ENGINEERING

CERTIFICATE

This is to certify that the dissertation entitled

‘Importance of Drug Features in Drug-Drug Interaction: A Comparative Study’

is a bonafide work carried out by

**Chetana N Patil
Naveen B N
Mekala Sanjana
Kumari Shivangi**

**PES2UG20CS504
PES2UG20CS217
PES2UG20CS194
PES2UG20CS173**

In partial fulfilment for the completion of seventh semester Capstone Project Phase - 2 (UE20CS461A) in the Program of Study -Bachelor of Technology in Computer Science and Engineering under rules and regulations of PES University, Bengaluru during the period June 2023 – Nov. 2023. It is certified that all corrections / suggestions indicated for internal assessment have been incorporated in the report. The dissertation has been approved as it satisfies the 7th semester academic requirements in respect of project work.

Signature
Dr.Prajwala.T.R
Associate Professor

Signature
Dr. Sandesh B J
Chairperson
External Viva

Signature
Dr. B K Keshavan
Dean of Faculty

Name of the Examiners

Signature with Date

1. _____

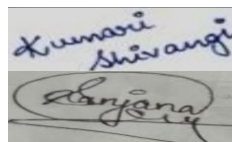
2. _____

DECLARATION

We hereby declare that the Capstone Project Phase - 2 entitled “**Importance of Drug features in Drug-Drug Interaction: A Comparative Study**” has been carried out by us under the guidance of Dr.Prajwala T R , Associate Professor and submitted in partial fulfillment of the course requirements for the award of degree of **Bachelor of Technology** in **Computer Science and Engineering** of **PES University, Bengaluru** during the academic semester June – Nov. 2023. The matter embodied in this report has not been submitted to any other university or institution for the award of any degree.

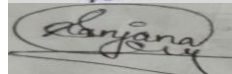
PES2UG20CS173

Kumari Shivangi



PES2UG20CS194

Mekala Sanjana



PES2UG20CS217

Naveen B N



PES2UG20CS504

Chetana N Patil



ACKNOWLEDGEMENT

I would like to express my gratitude to Dr.Prajwala T R, Associate Professor, Department of Computer Science and Engineering, PES University, for her continuous guidance, assistance, and encouragement throughout the development of this UE20CS461A - Capstone Project Phase – 2.

I am grateful to the Capstone Project Coordinator, Dr.Sarasvathi V, Professor and Dr. Sudeepa Roy Dey, Associate Professor, for organizing, managing, and helping with the entire process.

I take this opportunity to thank Dr. Sandesh B J, Chairperson, Department of Computer Science and Engineering, PES University, for all the knowledge and support I have received from the department. I would like to thank Dr. B.K. Keshavan, Dean of Faculty, PES University for his help.

I am deeply grateful to Dr. M. R. Doreswamy, Chancellor, PES University, Prof.Jawahar Doreswamy, Pro Chancellor – PES University, Dr.Suryaprasad J, Vice-Chancellor, PES University and Prof.Nagarjuna Sadineni, Pro-Vice Chancellor - PES University, for providing to me various opportunities and enlightenment every step of the way. Finally, this project could not have been completed without the continual support and encouragement I have received from my family and friends.

ABSTRACT

Drug target interaction is used to represent how a drug molecule interacts with a specific biological target, such as a protein or another drug, in order to develop a therapeutic effect. Drug repurposing is the process of identifying new therapeutic uses for existing drugs that were originally developed for a different purpose. It works as an attractive strategy in drug discovery, because of the several advantageous offers it has over traditional drug development, including reduced development time and costs, as well as increased safety and efficacy compared to newly developed drugs. The drug-target interaction is an important event in Drug repurposing as it defines whether the interaction is positive or negative. If an effect is positive then, mostly drugs can be repurposed. If it's negative then, repurposing will be mostly a failure.

In this project, Deep Learning techniques will be used for Drug Repurposing as it can train a model on a large dataset of known drug-target interactions. The model then uses the knowledge of all these interactions of drugs to predict potential drug-target interactions for existing drugs that have not been previously tested. These predictions can then be validated experimentally, and the most promising drug-target interactions can be further evaluated for their potential as new therapies. Additionally, deep learning can be used to predict the potential efficacy and toxicity of drugs, which can help to guide the selection of promising drug candidates for further development.

Thus, Using Deep Learning techniques as an approach for Drug Repurposing can accelerate the drug discovery process and help to bring new therapies to patients more quickly and efficiently.

TABLE OF CONTENTS

Chapter No.	Title	Page No.
1.	INTRODUCTION	01
	1.1. Overview	01
	1.2. Scope of the Project	01
2.	PROBLEM DEFINITION	02
	2.1. Problem Statement	02
3.	LITERATURE SURVEY	03
4.	DATA	15
5.	SYSTEM REQUIREMENTS SPECIFICATION	16
	5.1 Product Perspective	16
	5.1.1 Product Features	16
	5.1.2 User Classes and Characteristics	17
	5.1.3 Operating Environment	17
	5.1.4 General Assumptions and Dependencies	17
	5.1.5 Risks	18
	5.2 Functional Requirements	18
	5.3 External Interfaces Requirement	18
	5.3.1 User Interface	18
	5.3.2 Hardware Requirements	19
	5.3.3 Software Requirements	19
	5.4 Non Functional Requirements	20
	5.5 Other Requirements	21

6.	SYSTEM DESIGN (detailed)	22
6.1	Introduction	22
6.2	Current System	22
6.2.1	Decision Tree	22
6.2.2	Random Forest	23
6.3	General Low level Design Diagram	23
6.3.1	Exploratory data analysis(data visualization)	24
6.3.2	Data Preprocessing	25
6.3.3	Feature Engineering	25
6.4	ML Model Development	25
6.4.1	Train and Test Split	25
6.4.2	Evaluate Performance	26
6.5	ML GUI Deployment	26
6.6	Architecture Choice	27
6.7	Constraints, Assumptions And Dependencies	27
7.	PROPOSED METHODOLOGY	29
8.	IMPLEMENTATION AND PSEUDOCODE	33
8.1	Random Forest	33
8.2	SVM (Support Vector Machine)	34
8.3	XgBoost	35
8.4	XGBoost + AutoEncoders	36
8.5	Capsule Network	37
8.6	Convolutional Neural Network	38

9.	RESULTS AND DISCUSSION	39
9.1	Random forest	39
9.2	SVM	39
9.3	XGBoost	40
9.4	XGBoost + Autoencoders	41
9.5	Capsule Networks	41
9.6	CNN	42
9.7	Table for Comparision	43
	REFERENCES	47
	APPENDIX A ABBREVIATIONS	49

LIST OF FIGURES

Figure No.	Title	Page No.
4.1	The content of dataset	15
6.1	Graphical representation of Random forest	23
6.2	An overview of the low level design of the system	24
6.3	An overview of the design of the system	27
7.1	An overview of the high level design of the system	29
8.1	Code for Implementation of Random forest	33
8.2	Code for implementation of SVM	34
8.3	Code for implementation of XgBoost	35
8.4	Code for implementation of XgBoost + Autoencoder	36
8.5	Code for implementation of Capsule Network	37
8.5	Code for implementation of Convolution Neural Network	38
9.1	ROC Curve for Random Forest	39
9.2	ROC Curve for SVM	39
9.3	Precision-Recall Curve for SVM	40
9.4	ROC Curve for XGBoost	40
9.5	ROC Curve for XGBoost + Autoencoders	41
9.6	ROC Curve for Capsule Networks	41
9.7	Learning Curve for Capsule Networks	42

9.8	ROC Curve for CNN	42
9.9	Model accuracy and loss Curve for CNN	43
9.11	Importance of features	45
9.12	GUI of our system when interaction is present	45
9.13	GUI of our system when there is no interaction	46

LIST OF TABLES

Figure No.	Title	Page No.
9.10	Metrics analyzed for different algorithm	43

CHAPTER 1

INTRODUCTION

1.1 Overview:

Drug-Drug interaction is a process in which a drug molecule interacts with another drug molecule. People all across the world consume multiple drugs all day long. Different combination of drugs are administered-prevalent to underlying diseases or treatments. Combination of few drugs when administered in small doses daily may accumulate side-effects. We therefore need to learn the essence of such interaction and prevent future consequences.

As the world is progressing, there has been a huge upsurge in the number of different diseases and healthcare issues like vitamin deficiency, flues, and other life threatening diseases. Recently, we have faced a huge pandemic caused due to COVID-19 which has shown the need for improvement in the drug discovery process which can be possible using ML models.

Also, according to the stats generated in recent years, at least more than 60% of old age populations and some amount of the working population have health issues such that they have to consume multiple drugs simultaneously. The intake of different drugs all together can lead to side effects which could prove to be life threatening.

1.2 Scope of the Project

Our project has future in the Pharmaceutical industries as it can

- Reduce the time required for new drug discovery.
- Gives an essence for Drug Repurposing.

CHAPTER 2

PROBLEM DEFINITION

2.1 Problem Statement

The project depicts the importance of drug features in the interaction of two drugs when ingested. Drugs interact with each other and can produce synergistic effects (enhance the effect of medication), antagonistic effects (reduce the effect of medication) or neutral.

Each Drug has its own characteristics i.e., features. Some features might relate to its chemical structure whilst others depict its chemical properties. In any Drug - Drug Interaction, properties of metabolism, play an important role in deciding Pharmacokinetics ie, one drug might affect absorption or excretion of another drug ultimately leading to loss of efficacy of ingested drugs or causing adverse side effects. Likewise Drugs with similar structure might take similar metabolic pathways ultimately leading to interactions. Drugs with equivalent binding affinity might interfere with the receptor, increasing or decreasing their concentration in the body .Therefore, it becomes a major source of task to identify underlying causes and provide insights on such problems.

We tend to analyze the importance of these features in the Drug - Drug Interaction. For attaining such we aim to build multiple ML models and make a comparative study for Chemical properties.

CHAPTER 3

LITERATURE SURVEY

3.1 Inferring drug-target interactions based on random walk and convolutional neural network [1]

DTI-Pred is a framework implemented using random walk and CNNs to predict DTIs. The model works by calculating the drug similarity using the data on the drug-drug side effect association. Then, the existing information about various drug-protein interactions and similarities as well as other interactions and similarities are used to create a heterogeneous network. Since, DTI-Pred has CNNs so it provides good feature extraction for the representation of drug and protein.

The advantages are: Novel drug-protein interactions can be discovered. It gives best prediction performance due to the large amount of knowledge on the multiple drug-target interactions. First k candidates are correctly classified as the recall rate is high.

The disadvantages are: Due to imbalance in class (i.e; the positive example and the negative example) of the training data, the ROC cannot be used.

3.2 DDI-Pred: Graph Convolutional Network-based Drug-Drug Interactions Prediction using Drug Chemical Structure Embedding. [2]

SMILES representation is converted to SELFIES. (Syntactic, semantically correct and valid molecules).SELFIES are converted to Vectors by embedding Techniques. Link in graph says that interaction is present. Take the feature vectors and Existing DDI prediction to give new DDI styles. 4 layers GCN with 16d hidden node feature, Activation Function -RELU, Adam Optimizer and dropout at 10%.Link generator takes 2 nodes and checks whether a link is possible between nodes(drug)and outputs the same.

The advantages are: The values for AUPRC and the AUROC were closer to one which means prediction is good. SELFIES outperform SMILES representation with Accuracy, AUROC and Loss values.t-SNE visualization of the vectors gave an outright synergistic combination of DDI's and model beat other models in DDI predictions; an example would be mol2vec.

The drawbacks are: There is a requirement of more number of hidden layers would have been a bonus for further increase in accuracy. It requires registered DDI's for prediction and if input of DDI network is very sparse, it provides poor results.

3.3 Graph Convolutional Auto-encoder and Generative Adversarial Network-Based Method for Predicting Drug-Target Interactions [3]

A heterogeneous drug-target network is used to concatenate different drug-target connections, drug-target or drug-target similarities and interactions, and drug-target interactions. A graph convolutional auto-encoder was used to learn network embedding of drug and target nodes in a low-dimensional feature space, and the auto-encoder deeply combined different kinds of connections in the network. A GAN was used to regularize the feature vectors of the nodes into a Gaussian distribution. Severe imbalance between known and unknown DTI. Therefore, we developed an ensemble learning model-based classifier, Light GBM, to estimate drug-target interaction propensities. This classifier made full use of all unknown DTIs and balanced the negative effect of class imbalance. Experimental results have shown that GAN DTI outperforms several state-of-the-art methods for DTI prediction.

The advantages are: By building several decision trees, light GBM may effectively alleviate the effect of class imbalance. It trains distinct trees using a different dataset of negative samples, ensuring that the negative samples are fully utilized.

The limitations are: Complex architecture. Difficulty in handling noisy.

3.4 HS-DTI: Drug-target Interaction Prediction Based on Hierarchical Networks and Multi-order Sequence Effect [4]

The drug encoder obtains functional group information from hierarchical learning on the molecule graph constructed from the SMILES sequence. The protein encoder has a convolution block and a multi-order sequence learning block in which the 1st and 2nd order sequence information is assembled to extend the protein encoder and extract different levels of features from the convolution process. Drug functional group information and multi-order protein sequence information are obtained efficiently. Moreover, we found that the HS-DTI model containing these two modules provides the best performance. Both the drug property and the protein property are concatenated and fed to a multilayer perceptron to predict the interaction.

The advantages are: Reduced experimental costs and time. Ability to handle heterogeneous data. HS-DTI seeks to extract numerous properties of drugs and proteins through functional group structures and incorporating multi-line sequence information.

The limitations are: The structure and function of certain functional groups in pharmaceutical preparations and the multiple sequence effect of proteins were not taken into account in this model. Data quality issues: The accuracy of the HS-DTI model is strongly dependent on the quality of the input data. Inaccuracies or biases in the data can lead to incorrect predictions. Limited data availability: The HS-DTI model requires large amounts of data for training, which may be limited in some cases, leading to over-fitting or poor performance.

3.5 A Graph Convolution-Transformer Neural Network for Drug-Target Interaction Prediction [5]

Three different Datasets from Drug Bank. Drugs-SMILES were converted to graphs using an adjacency list (containing atoms and bond) and also maintained feature representation of 9 attributes.

(Mol2Vec).Proteins-FASTA sequence was converted into a vector using Word2Vec.Transformer is a Deep learning model encompassing the nature of “self-attention”, meaning predicting linkages between nodes (query, key and value vectors).Output of Mol2Vec is passed onto the GCN layer and extracts information, whilst protein sequences undergo 1d convolution and attention mechanism .The outputs namely molecules and protein are transformed/converted to similar dimension with ending a fully connected layer. Lastly, inter-sequence self-attention is provided by the decoder (3 layers) and predicts binary class (present/not).

The advantages includes -Mapping attention weights to atoms displayed the model decision mechanism. The model also showed by paying attention to the atoms, which enormously contributed to the interaction’s chemical bond formation. Self-attention mechanism along with GCN transformer was a novel approach.

The disadvantages that follow using the model would be that it did not outperform the Deep DTA model in terms of 2 datasets. Partially extract the molecular information in an interaction. The performance decay on the Kinome imbalanced dataset shows that the dataset can dominantly influence performance. The model focuses on H-bond donors (oxygen, nitrogen).

3.6 Using Supervised Machine Learning Algorithms For Drug Target Interaction Prediction [6]

This study obtains a low-dimensional vector for graph nodes using the node2vec method. Node2vec framework is used to get low dimensionality for nodes. Different types of supervised learning algorithms have been used to estimate drug-target associations from low-level representations. Multiple machine learning methods have been used to predict drug-target interactions. To evaluate the proposed models, the calculated AUROC and AUPRC values and obtained results indicate that nonlinear SVM and logistic regression performed better than other models with AUROC and AUPRC values of 0.8317 and 0.8260, respectively. A bipartite network is used.

The advantages are: Down sampling was used to try to fix the imbalance in the dataset. The Node2vec framework is used to create small dimensions for a node.

The cons of the model are: It may lead to over-fitting. Logistic regression is not suitable for modeling complex relationships between variables, such as interactions or nonlinear effects. Logistic regression is limited to linear relationships.

3.7 Deep learning in drug design: protein-ligand binding affinity prediction [7]

Deep Atom is a framework to predict the protein-ligand binding affinity. A 3D grid box is used on the protein ligand 3D structure at their binding site so as to get information about it. PC Max algorithm is used to collect information from the binding site surroundings. Light weighted CNN is used as feature extractors to get the useful information regarding the interaction according to the binding affinity labels. The positive aspects that we inferred from the study include, Due to light weighted CNN, the feature extraction was more accurate than other models and in future gives the model ability to be more accurate as the amount of training data increases.

The advantages are: Reduced cost and computation even though 3D models are used. No restrictions on the color channels for protein and ligand representations unlike in computer vision (where we have to show RGB). Through training and Over fitting is avoided using a shared weights FC (Fully Connected) layer. Improved performance and reduced variance in prediction.

The cons of study include: Performance is not that good as a limited amount of training data was used. Due to 3D CNNs computation cost and trainable parameters increases. Increased resolution gives more clarity on the interaction but causes a huge computational cost. Therefore, the optimal resolution is used instead of the best solution.

3.8 Drug target interaction prediction: end-to-end deep learning approach [8]

A model using CNNs to get the 1D representations from protein sequences and compound SMILES strings. The output received can be used in FCNNs for binary classification as CNNs are good at feature detection which gives useful information about patterns or other dependencies. The deep learning model is better than any traditional descriptor for classification of positive as well as negative interactions between the drugs.

The advantages are: Good and efficient identification of meaningful regions for the study of drug-target interaction. Due to CNNs feature extraction is good.

The drawbacks are: Not all deep learning models are that efficient in comparison to traditional models. (eg.- CNN with random forest walk). Also, additional information may prove to be useful to correctly identify positive interactions. As a CNN, auto encoder and FCNN combined model has more information. Hence, more sensitivity than the used model.

3.9 Drug-Target Interaction Prediction Based on Multisource Information Weighted Fusion [9]

Most present research presume that there are no known interactions between medicines and targets. Un-labelled samples are defined in this paper as those for which the link between medicines and targets has not been determined. Three approaches are used to screen un-labelled samples: drug similarity, random walk with restart, and WNN-GIP. Finally, the training set's interaction matrix is changed based on the fusion results, and the BLM-NII model is used to forecast interactions. Enzymes (Es), ion channels (ICs), G-protein-coupled receptors (GPCRs), and nuclear receptors (NRs) are all represented in the databases. To demonstrate the validity of the multisource information fusion-revised dataset in the proposed technique is by comparing accuracy, sensitivity, specificity, and precision to the BLM-NII method.

The advantages are: Experiments show that the suggested method has significantly higher specificity, sensitivity, precision, and accuracy than the BLM-NII method. DTI prediction concerns, samples with uncertain labels are frequently viewed as negative samples, influencing the results and imposing certain constraints. The weighted fusion method improves the efficacy and dependability of screening outcomes.

The disadvantages are: It performs better in datasets with more samples, but its generalization ability deteriorates in datasets with fewer examples. The prediction accuracy needs to be increased, especially for datasets with fewer positive results.

3.10 Drug-drug interaction extraction from biomedical texts based on multi-attention mechanism [10]

In this model again we look at attention mechanisms to calculate import scores on the neural network. Layer called BI-LSTM performs conversion of each word into multidimensional vectors and then to learn meaningful and useful sentences or creating their representation is done through Bi directional LSTM layers. As we know in attention layer we calculate contribution of each particular word against all others. The use of [CLS] tags provokes reduction of training time and performance improvement.

If we look at key takeaways from the article we get to see that the model allows to calculate the attention weights in the neural network, which can be utilized to measure the contributions of different words while the model makes decisions. Key feature include classified tag word, named as [CLS] used to learn the global information for DDI classification, experimental results have shown that our model has a competitive advantage in extracting the relation between the two candidate drugs in one instance.

3.11 Drug-target interaction prediction using Multi Graph Regularized Nuclear Norm Minimization [11]

Multi-Graph Regularized Nuclear Norm Minimization, a new technology, predicts drug-target interaction networks utilizing three inputs: known drug-target interaction networks, similarities across drugs, and those across targets. To capture the proximities, numerous drug-drug and target-target similarities are used as multiple graph Laplacian (across drugs/targets) regularization terms. The four newly incorporated similarities computed from the interaction matrix are Cosine similarity, Pearson's linear correlation, Hamming similarity, and Jaccard similarity. They used Singular value shrinkage and Multi Graph regularized Nuclear Norm Minimization, both of which are efficient solvers for Nuclear Norm minimization. To compare our method to five other state-of-the-art methods (WGRMF, CMF, RLS_WNN, NRLMF, TMF), they employed three separate cross-validation settings: CVS1 (random drug-target combinations excluded), CVS2 (entire drug profile excluded), and CVS3 (entire target profile omitted).

The advantages are: The multi graph regularized version of Nuclear norm minimization tries to prevent over-fitting and considerably improve generalization capabilities. The suggested method focuses on constructing a low-rank interaction matrix based on the proximities of medicines and targets encoded by graphs.

The limitations are: The issue with traditional nuclear norm minimization (NNM) is that it does not support linked information such as Similarity matrices for Drugs and Targets.

3.12 Predicting Drug-Drug Interactions Using Deep Neural Networks. [12]

This paper includes a model which would predict about 80 different Drug-Drug Interaction patterns using deep neural network as its model. Information pertaining to the structure of drugs was taken from the Drug bank. Drug description (descriptors) of drugs is taken through Chemical databases. 3,126,838 drug-drug pair interaction was refined and added 2216 additional features to existing ones and predicted interaction (4432). 4 hidden layers, Activation function -RELU, Loss- Categorical Cross Entropy. Mini batch size of 256 to improve performance. Dropouts of 20% (randomly) to improve accuracy.

If we were to look at the positive aspects of the paper we find that the performance obtained from this model is relatively good with values of Area under the Curve (AUC) -93.2% and for test case (set) was 94%. The estimation was applied on IBD drugs which showed positive synergistic effect. Showed a significant 5% rise in accuracy compared to SVM implementation. Increased antipsychotic activities.

Paper also has few key points of failure which include that the ROC curve of certain DDI types did not show optimized results due to fewer samples (smaller dataset -5,134 drugs and more hidden layers would have been a bonus for further increase in accuracy).

CHAPTER 4

DATA

4.1 OVERVIEW:

Drug Bank is a comprehensive, publicly available online collection of medication and drug target information. It is a valuable resource for researchers, healthcare professionals, and students in the pharmaceutical and biotech industries. The Drug Bank dataset is huge and consists of information about numerous drug molecules.

Data regarding interaction of Drugs and SMILES were obtained from the link given below,

(<https://github.com/haichengyi/DDIPred>)[10].

Properties of Drugs were obtained from Kaggle for 8289 Drugs with their associated Drug Bank id given below

(<https://www.kaggle.com/code/abdelrahmankhalil/drug-bank-smiles>)

	A	B	C	D	E	F	G
1	drugbank_id	cas	logP ALOGPS	logP ChemAxon	solubility ALOGPS	pKa (strongest acidic)	pKa (strongest basic)
2	DB000006	128270-60-0	-0.76	-14	4.64e-02 g/l	2.79	11.88
3	DB000007	53714-56-0	1.04	-2.4	3.38e-02 g/l	9.49	11.92
4	DB000014	65807-02-5	0.3	-5.2	2.83e-02 g/l	9.27	10.82
5	DB000035	16679-58-6	-1	-6.1	1.10e-01 g/l	9.5	11.77
6	DB000050	120287-85-6	1.33	-1.7	6.94e-03 g/l	9.5	11.79
7	DB000067	11000-17-2	-1.4	-7.2	1.24e-01 g/l	7.65	11.5
8	DB000080	103060-53-3	-0.47	-9.4	1.73e-02 g/l	2.98	9.59
9	DB000093	56-59-7	-1.1	-5.8	4.53e-02 g/l	11.39	10.18
10	DB00104	83150-76-9	0.42	-1.4	1.22e-02 g/l	11.4	10.17
11	DB00106	183552-38-7	2.84	-0.46	3.71e-03 g/l	9.47	10.66
12	DB00114	54-47-7	-0.55	-2.1	5.70e+00 g/l	1.68	4.11
13	DB00115	68-19-9	2.66	-3.2	2.02e-02 g/l	1.82	8.68
14	DB00116	135-16-0	-0.96	-4.2	2.69e-01 g/l	3.51	3.58
15	DB00117	71-00-1	-2.7	-3.6	7.13e+01 g/l	1.85	9.44
16	DB00118	29908-03-0	0.05	-5.3	1.62e+00 g/l	1.7	9.41
17	DB00119	127-17-3	-0.38	0.066	1.34e+02 g/l	2.93	-9.6
18	DB00120	63-91-2	-1.4	-1.2	4.14e+00 g/l	2.47	9.45
19	DB00121	58-85-5	0.17	0.32	1.22e+00 g/l	4.4	-1.9

4.1 The content of dataset.

The properties include pKa acidic and basic values, logP ALOGPS, logP ChemAxon, solubility ALOGPS.

CHAPTER 5

SYSTEM REQUIREMENTS SPECIFICATIONS

5.1 Product Perspective

The following section deeply describes the essential product representatives that are essential to our model.

5.1.1 Product Features

- 1. Predictive Modeling:** The product can look through large amounts of data and levels of information and predict and model the required drug candidates effectively.
- 2. High Throughput Screening:** The obtained interaction prediction is subjected to testing and validation for further use.
- 3. Personalized Medicine:** One can predict based on one's health and genetic data and can be monitored.
- 4. Good Visualization:** Interaction can be visualized using 3d modeling techniques.
- 5. Downloadable:** The predicted interaction can be stored on local storage.
- 6. Patient-specific recommendations:** The ability to generate patient-specific recommendations for managing drug-drug interactions based on the patient's medical history and current medication regimen.
- 7. Integration with electronic health records:** The ability to integrate with electronic health records (EHRs) to access patient data and provide real-time alerts and recommendation.
- 8. Customizable alerts:** A feature that allows users to set alerts for specific drugs or drug combinations that may be of particular concern.

5.1.2 User Classes and Characteristics

- 1. Industrial User:** These are pharmaceutical companies that are used to synthesize new drugs or are involved in a quality driven process.
- 2. Pharmacists and Research Professionals:** Drug Discovery involves a lot of steps from drug design to toxicity prediction. They can use this as a source of information and process them effectively.
- 3. Diagnostic User:** He is responsible for diagnosing a disease or medical condition using drugs.
- 4. Regulatory agencies:** These are governmental organizations responsible for validation and approval of drugs.
- 5. Prophylactic User:** They are to prevent onset of pandemic/epidemic diseases.

5.1.3 Operating Environment

Operating Systems: Windows, Linux, Mac-OS.

Graphics: Nvidia. High Graphics processor /GPU preferred.

Operating Processors: 4.

RAM: 8GB

5.1.4 General Constraints, Assumptions and Dependencies

- 1. GPU:** Calculations and large datasets are used so the processing power should be good.
- 2. Dataset:** Lack of biological and structural data.
- 3. AI Model:** It is an AI model which cannot give perfect accuracy. If an AI model leads to incorrect prediction of target and it may harm patients.
- 4. Complex to interpret:** Being computer science Students, we lack the knowledge of molecules and Pharmacology.

We assume that, Drug reaches the depicted target and interacts with that. Drug combines with the target to provide desired effects. Drug -target interaction speeds up the drug discovery process.

5.1.5 .Risks

1. GPU's are required for faster processing.
2. Resources and Scientific Research pertaining the same are at the bottom.
3. Lack of visualization and protein structure information.
4. Chemical knowledge and libraries are less.

5.2 Functional Requirements

1. The system would be highly interactive and visualizable.
2. It shall take varied levels of drugs, an input (with a wide range of different properties), protein sequences and predict the possibility of interaction.
3. The interaction would be depicted in high dimensional resolution.
4. The interaction model would also depict labels on 3d models of different amino acids and gene sequences.
5. Chemical descriptions would also be provided.
6. The interaction depiction can be easily downloaded.
7. Drug screening would depict further use of drugs for either drug repurposing or new drug synthesis.

5.3. External Interface Requirements

5.3.1 User Interfaces:

Command-Line Interfaces (CLIs): CLIs are text-based interfaces that enable users to interact with the system using command-line input. They are often used by advanced users and can provide more flexibility and control than GUIs, but can be less user-friendly for novice users.

5.3.2 Hardware Requirements

1. GPU for computations.- Nvidia
2. Good Graphics for displaying interaction effectively.
3. Capable RAM size.-8GB
4. Operating Processors: 4 (like 2.20 GHz Intel i7-8750H)

5.3.3 Software Requirements:

5.3.3.1 Operating System

Windows, Linux, Mac-OS

5.3.3.2 Tools and libraries

- 1) **Tkinter:** A python library used to create graphical user interface.
- 2) **Numpy:** NumPy can be used to perform a wide variety of mathematical operations on arrays. It adds powerful data structures to Python that guarantee efficient calculations with arrays and matrices and it supplies an enormous library of high-level mathematical functions that operate on these arrays and matrices.
- 3) **Pandas:** Pandas is a Python library used for working with data sets. It has functions for analyzing, cleaning, exploring, and manipulating data.
- 4) **Keras:** Keras is a high-level, deep learning API developed by Google for implementing neural networks. It is written in Python and is used to make the implementation of neural networks easy. It also supports multiple backend neural network computation.
- 5) **Tensorflow:** The TensorFlow platform helps you implement best practices for data automation, model tracking, performance monitoring, and model retraining.
- 6) Some libraries are imported and used for implementation of machine learning algorithms.

5.4 Non-Functional Requirements

5.4.1 Performance Requirement:

1. Interaction should be deprecated as fast as possible. The latency and delay in the display must be as less as possible.
2. The system should be designed to perform efficiently and effectively, with acceptable response times for all interactions. It should be able to handle large datasets and complex computations without sacrificing performance.

5.4.2 Security Requirements:

1. **Data security:** The tool must comply with data security and privacy regulations, including secure storage and transmission of patient data.
2. **Adverse event reporting:** The tool should have a mechanism for healthcare professionals to report adverse events related to drug-drug interactions, which could help improve the tool's accuracy and safety over time.
3. **Systematic review:** The tool must undergo systematic reviews and evaluations to ensure its safety and effectiveness, which could help identify and address potential safety issues.

5.4.3 Safety requirements:

1. **Real-time updates:** The tool must be regularly updated with the latest drug information, interactions, and adverse event data to ensure the accuracy of its recommendation.
2. **Integration with EHRs:** The tool must integrate with EHRs securely to avoid errors and to provide accurate patient-specific recommendations.

5.5 Other Requirements

- 1. Scalability:** It should be able to handle a growing number of interactions and users without experiencing significant degradation in performance.
- 2. Maintainability:** The system should be designed to be easily maintained, upgraded, and enhanced over time. It should be well-documented and designed with modularity and extensibility in mind to facilitate future development and improvements.
- 3. Usability:** The system that is built should have essential key requirements that are simple to use, can be learnt easily, and navigated easily. It must be designed or decorated with the needs and skill level of the user in mind and provide a user-friendly interface that enables efficient drug target interaction.
- 4. Reliability:** The system should be dependable and operate without errors or interruptions. It should be designed to prevent data loss, ensure data integrity, and maintain consistent performance under varying conditions.
- 5. Interoperability:** It will be operable on multiple OS's as researchers use Windows and Linux more. So, main focus will be there.
- 6. Performance:** Speed at which the interaction depiction displayed to the user is of utmost importance. Multiple users accessing the site at a time should not impact the performance.
- 7. Security:** Researchers performing critical research can have utmost privacy of their work and will not be made public until they owe to.
- 8. Reusability:** Certain parts of code as interaction can be made up for reuse. We also permit the system to be used up in De novo Drug Design applications if created in future to use up our methodology for Interaction prediction.
- 9. Application Compatibility:** We highly motivate ourselves to use Browser-stack or Cross Browser Testing tools to make our application compatible across different platforms.

CHAPTER 6

SYSTEM DESIGN

6.1 Introduction

Drug target interaction is the process by which a drug molecule interacts with a specific biological target, such as a protein or receptor, in order to elicit a therapeutic effect. Drug repurposing, also known as drug repositioning, is the process of identifying new therapeutic uses for existing drugs that were originally developed for a different purpose. Drug repurposing has become an attractive strategy in drug discovery, as it can offer several advantages over traditional drug development, including reduced development time and costs, as well as increased safety and efficacy compared to newly developed drugs.

In our project, we are going to use a deep learning approach, a type of machine learning that uses neural networks to analyse large-scale data sets. The following documentation will discuss design considerations that include existing designs, design details, and architecture style used and talk about reliability and performance of the system. The general high level design diagram, architecture style diagram, logical data flow diagram, proposed approach diagram, as well as a UI diagram will be discussed.

6.2 Current System

Predicting whether two drugs or drug and protein will interact with each other. We used Random forest algorithm for prediction as other algorithms were having less accuracy. It will consider two drugbank_id's as input and returns whether the interaction is there or not. If interaction found returns 1, else 0 based on drug features, along with the most feature importance.

6.2.1 Decision Tree

A supervised learning method used for classification and regression. A Decision Tree is a tree in which each internal node represents a feature test, each branch a possible test outcome, and

each terminal node a class label. The optimum splitting criteria for each node must be determined at each node, potentially making decision trees computationally expensive. Decision Trees are valuable for visualizing the decision-making process and identifying key factors influencing student performance. Pruning is a method of reducing the size of a decision tree. By restricting the size of the tree, it decreases the possibility of overfitting.

6.2.2 Random Forest

A random forest is a grouping of multiple individual decision trees that work together to create an ensemble model. Each decision tree in the random forest makes a class prediction, and the class with the most votes become the model's prediction.

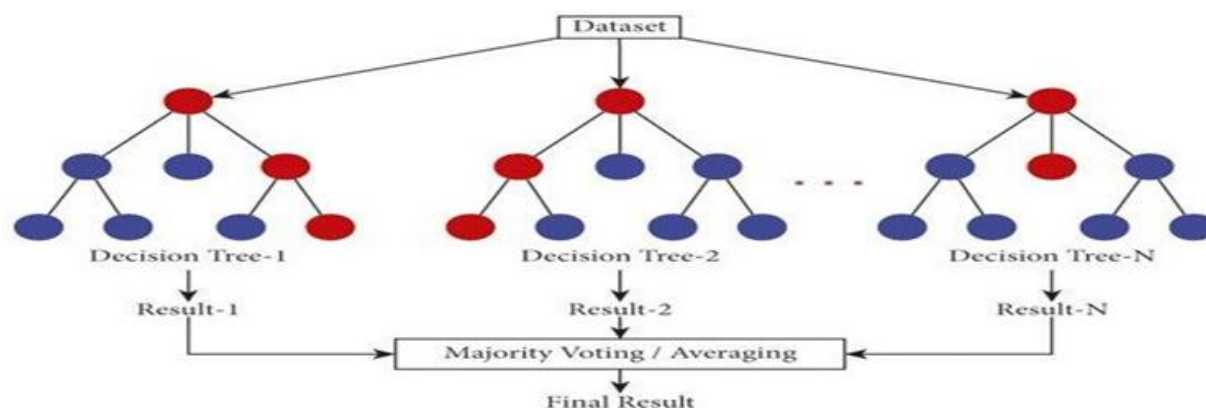


Fig 6.1: Graphical representation of Random forest

Random Forest introduces randomness both in the selection of features for each tree and the data points used to build them. By combining diverse trees, the algorithm provides robust and interpretable predictions, making it well-suited for visualizing the impact of various factors on student performance.

6.3 General Low level Design Diagram

The high level design aims to shed some light on the broad approach taken during the design of our system. Our design consists of a simple but efficient three module system. The basic components consists of Data preprocessing module, Descriptor (Characteristic property of the drug compounds) Derivation module and finally the Machine learning model.

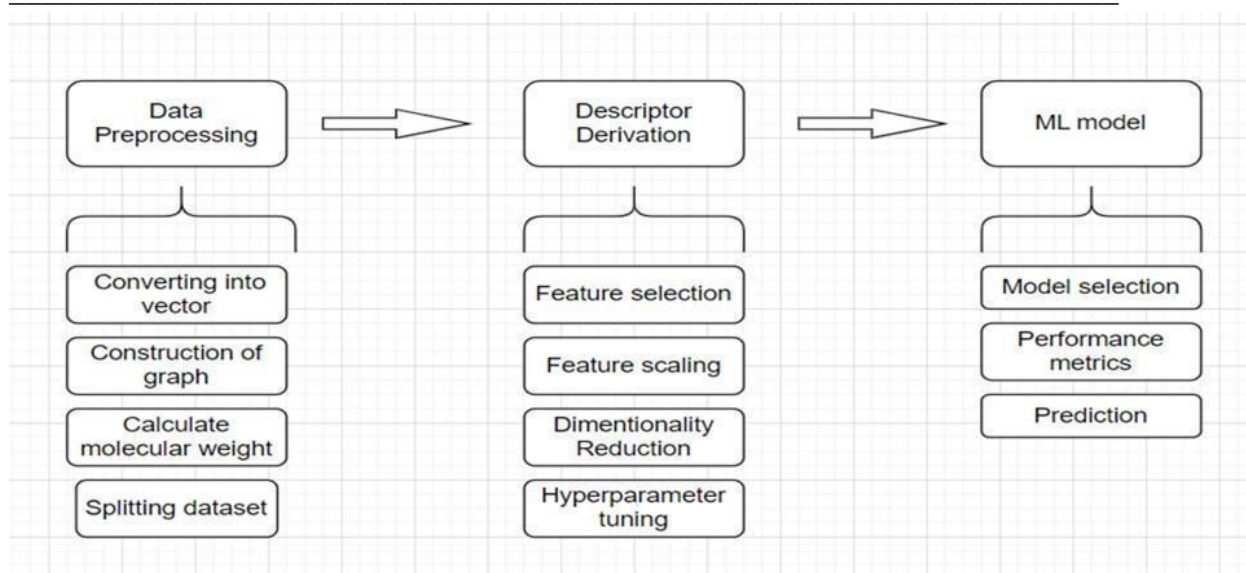


Fig 6.2: An overview of the low level design of the system.

6.3.1 Exploratory Data Analysis(Data Visualization)

It refers to the visual representation of data. It is an effective way of conveying the data when it is numerous. Various plots and graphs such as scatter plots, bar graphs, pie chart, feature importance charts, hyperplane animation, ROC curve and confusion matrix allow us to interpret the variability of the data and uncover patterns.

A Receiver Operating Characteristic (ROC) curve: is a graphical representation of a binary classification model's performance across different threshold settings. It plots the true positive rate against the false positive rate, providing insights into the trade-off between sensitivity and specificity. The area under the ROC curve (AUC) is a common metric to assess the overall performance of a classification model, with a higher AUC indicating better discrimination capability.

Confusion matrix: is a table used in machine learning and statistics to assess the performance of a classification model. It summarizes the results of classification by showing the counts of true positive, true negative, false positive, and false negative predictions.

Hyperplanes: are decision boundaries that help classify the data points. Data points falling on either side of the hyperplane can be attributed to different classes. Also, the dimension of the hyperplane depends upon the number of features.

6.3.2 Data Pre-Processing

Data pre-processing refers to various activities which prepare the raw dataset and make it suitable for training in a machine learning model. It is an essential step in machine learning. This step ensures that the input data for algorithm visualization is clean and ready for analysis.

6.3.3 Feature Engineering

Feature engineering includes feature selection, feature extraction and column transformations. Relevant features might include previous exam scores, study hours, participation in extracurricular activities, etc.

Feature selection for a ML model is imperative to the final prediction as in this step the highest contributing features are selected based on a metric known as feature importance. Features with low feature importance will inevitably be dropped.

The column transformations can include encoding categorical data into numerical data that the machine learning model will be able to interpret with techniques such as one-hot encoding, binary encoding, target encoding, frequency encoding, label-encoding etc.

6.4 ML Model Development

6.4.1 Train and Test Split

One way to evaluate the performance of machine learning algorithms is to split training and testing.

It can be applied to tasks involving classification, regression, and any supervised learning technique.

This step requires splitting the data set into 80% training subset and 20% testing subset. The training set is the initial subset used to fit the model. The model was not trained on the test subset. Instead, it is fed into an input array of datasets and its predictions are constructed and compared to the expected values. This contrast highlights the model's performance when given previously unseen data.

Train Dataset: To fit machine learning model.

Test Dataset: To evaluate the fit machine learning model.

The premise of the split is to evaluate the performance of machine learning model using unused data to train the model.

6.4.2 Evaluate Performance

ML model performance is evaluated using metrics such as accuracy, precision, recall, F1 score and area under ROC curve. Accuracy measures the overall correctness of the model's predictions. Precision assesses the accuracy of positive predictions, while recall evaluates the model's ability to capture all relevant instances. F1-score is the harmonic mean of precision and recall. The ROC curve provides a graphical representation of the trade-off between true positive rate and false positive rate, and the area under the curve (AUC) quantifies the model's ability to distinguish between classes. These metrics collectively offer a comprehensive evaluation of an ML model's performance in predicting student outcomes, providing insights into its strengths and areas for improvement.

6.5 ML GUI Deployment

Tkinter is a GUI application framework library in Python. It will use the ML models in backend to predict interaction and gives the most important feature which is responsible for drug interaction.

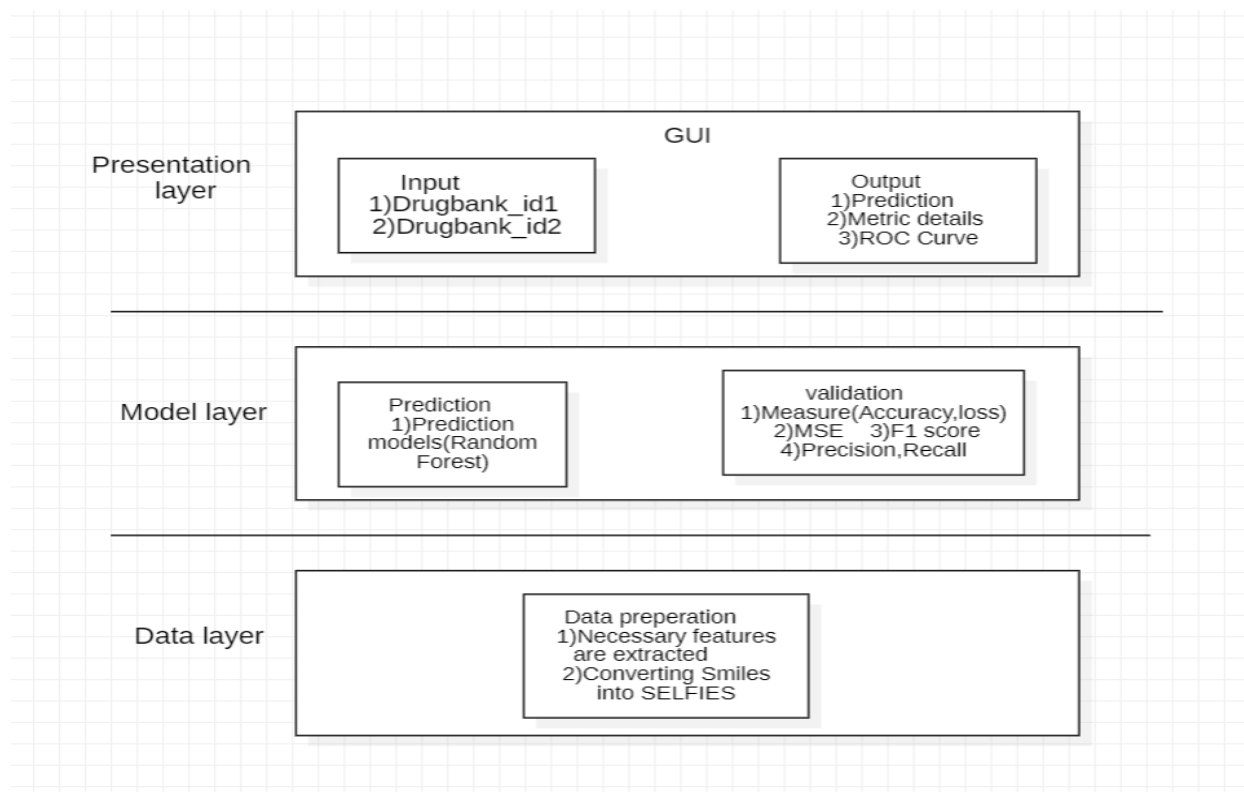


Fig 6.3: An overview of the design of the system.

6.6 Architecture Choice

Among all the available classification methods, random forests provide the highest accuracy. The random forest technique can also handle big data with numerous variables running into thousands. It can automatically balance data sets when a class is more infrequent than other classes in the data. Random forest is a commonly-used machine learning algorithm which combines the output of multiple decision trees to reach a single result. Its ease of use and flexibility have fueled its adoption, as it handles both classification and regression problems.

6.7 Constraints, Assumptions And Dependencies

Like any other machine learning method, the use of Random forest for drug target interaction has certain constraints, assumptions, and dependencies. Here are a few:

Constraints:

Data availability: The effectiveness of GCNs in drug target interaction prediction is dependent on the availability and quality of data on drug-target interactions, which can be limited in some cases.

Model complexity: GCNs can be computationally expensive, which can be a constraint when dealing with large datasets or limited computing resources.

Interpretability: GCNs are sometimes viewed as a "black box" method, making it difficult to understand how the model arrived at its predictions.

Assumptions:

- There should be some actual values in the feature variables of the dataset, which will give the classifier a better chance to predict accurate results, rather than provide an estimation.
- The predictions from each tree must have very low correlations.

Dependencies:

Data quality: The accuracy of Random forest predictions depends on the quality of the data used to train the model. Noisy or incomplete data can lead to inaccurate predictions.

Hyper-parameter tuning: The performance of Random forest can be sensitive to the choice of hyper-parameters, such as learning rate and number of layers. Appropriate hyper-parameter tuning is crucial to achieve optimal performance.

Model architecture: The effectiveness of Random forest for drug target interaction prediction is also dependent on the specific architecture chosen for the model, as different architectures may have varying levels of performance on different datasets.

CHAPTER 7

PROPOSED METHODOLOGY

We considered various algorithm to predict interaction between two drugs, Random forest and XgBoost gave high accuracy compared to other models.

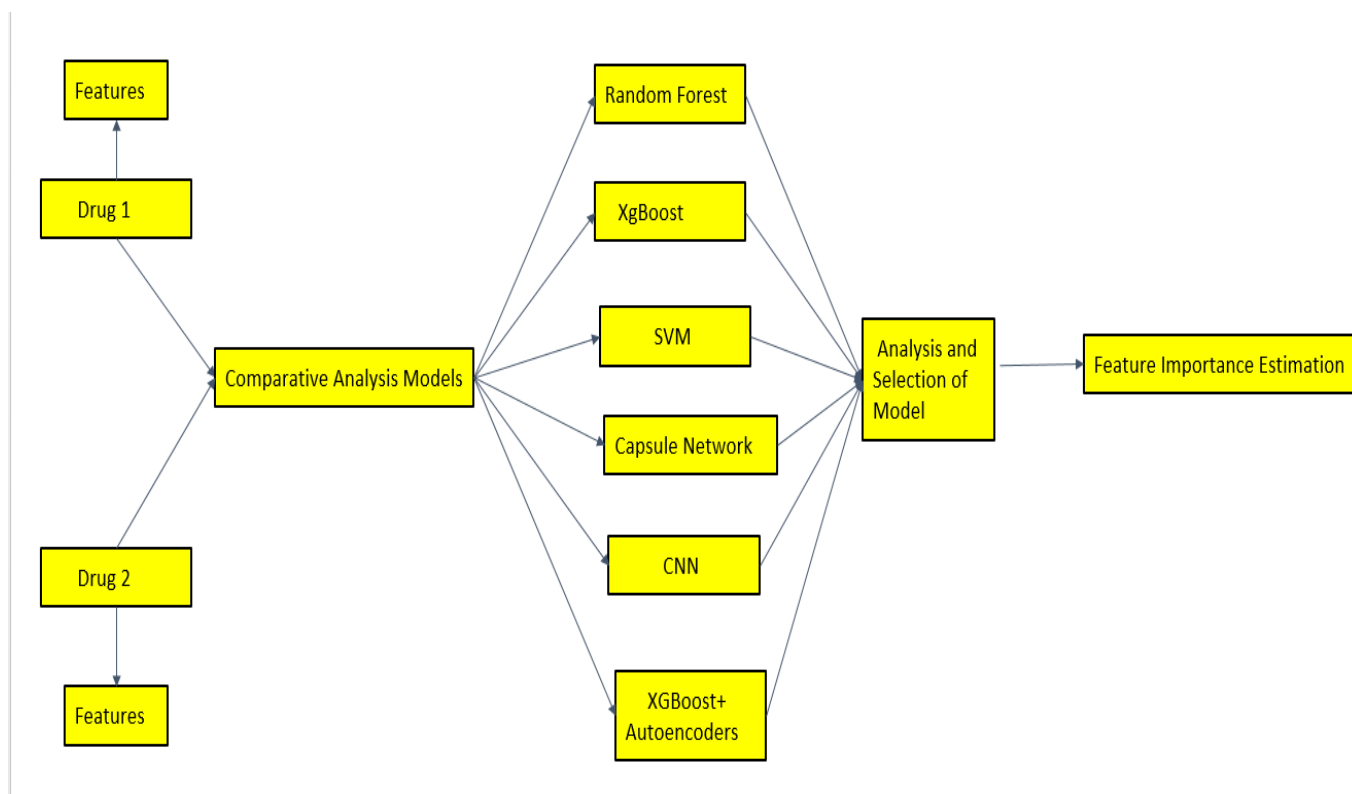


Fig 7.1: Figure represents an overview of the high level design of the system.

Finally we used Random forest algorithm to predict interaction between two drugs. Random forest uses group of decision trees to make decision. Random forest is a commonly-used machine learning algorithm which combines the output of multiple decision trees to reach a single result. Its ease of use and flexibility have fueled its adoption, as it handles both classification and regression problems.

Each drug is characterized by five features mainly dealing with solubility and permeability namely logP ChemAxon, logP ALOGPS, solubility, pKa acidic value and pKa basic value. We

take supervised learning approach to solve the problem. The data for source and target interaction was studied with features. Then the dataset was analyzed using multiple models.

A. Dataset Preprocessing

Individual preprocessing of the dataset is carried out.

Step 1: Removal of Unwanted Features:

We proceed with use of 'drugbankid' as our reference id for the identification of drug and 'cas' drug identifier column is unnecessary and hence is removed from Drug Features Dataset. 'SMILES' column present in Drug Names dataset provides no use to our estimation and is eliminated.

Step 2: Joining of datasets:

First part of this step includes inner join on two datasets i.e, Drug Features Dataset and Drug Names Dataset on the column 'drugbank id'. New dataset formed is named as the Drug Feat_Names dataset. In the next part we join the 'source' column from the Drug-Drug Interaction dataset with the 'drugbank id' column of Drug Feat Names dataset. Results yielding from this join are saved as Source Dataset. We rename columns from 'drugbank id', 'logALOGPS', 'logPChemAxom', 'solubility', 'LOGPS', 'pKa(strongest acidic)' and 'pKa(strongest basic)' to 'source drugbank id', 'source logP', 'source logP CA', 'source solubility', 'source acidic pKa' and 'source basic pKa'. Similarly we tend to join the 'target' column from the Drug-Drug Interaction dataset with the 'drugbank id' column of Drug Feat Names dataset. Results yielding from this join are saved as Target Dataset. We rename columns from 'drugbank id', 'logALOGPS', 'logPChemAxom', 'solubility', 'LOGPS', 'pKa(strongest acidic)' and 'pKa(strongest basic)' to 'target drugbank id', 'target logP', 'target logP CA', 'target solubility', 'target acidic pKa' and 'target basic pKa'. As a part of the last step we merge Source Dataset and Target Dataset based on Drug Drug Interaction dataset and save as Merged Dataset. This dataset comprises of columns 'source drugbank id', 'source logP', 'source logP CA', 'source solubility', 'source acidic pKa', 'source basic pKa', 'target drugbank id', 'target logP', 'target logP CA', 'target solubility', 'target acidic pKa' and 'target basic pKa'.

Step 3: Removal of Null values

Any row entries having NA values are removed in this step. After this step our dataset is ready to be utilized by the models.

B. Application of Models:

RandomForestClassifier: RandomForestClassification technique creates multiple decision tree estimators by the utilization of different dataset subsets. We split the dataset for training and testing by 80:20 ratio. And then train the model using 100 tree estimators to output the predictions.

XGBoost: XGBoost is the most powerful form of gradient boosting techniques. The model comprises decision tree estimators analyzed sequentially rather than parallel as in case of RandomForestClassification. Here again we make use of 100 decision tree estimators for classification and 80 data for training.

XGBoost + Autoencoders: This model takes the predictions from XGBoost and combines with the encoded features obtained from the autoencoders. A simple classifier layer is trained above the encoded features. Here again data is split in the ratio 80:20.

SVM: SVM draws a hyperplane separating classes. Here the hyperplane studies the drug features and solely separates into two parts i.e., the drug classes that interact and the ones which do not.

CNN: Input layer receives features and passes down to two layers of convolution and a flatten layer to extract essential features. Activation unit used in the Convolution layer is RELU. A dense layer added on top of it for classification purposes. Activation unit used in the dense layer is Sigmoid. Then we compile the model with binary cross entropy loss, Optimiser as Adam and activation unit as RELU and is made to run for 20 epochs to obtain desirable results.

Capsule Networks: Capsule Networks are one of the widely used Deep learning technique. This model comprises three layers. Input layer reads the features and is passed to the first layer of Capsule Network i.e. Primary Capsule layer which extracts essential features in the form of

vectors. Followed by that extracted features are sent to Digit Capsule Layer where Squash function is used to extract the spatial relationship amongst the vector data. The relationship learnt in this layer is utilized by Class Capsule layer to predict the results function $\text{squash}(\text{vectors})$ squared norm = $\text{sum}(\text{square}(\text{vectors}))$ scale = $\text{squared norm} / (1 + \text{squared norm}) / \text{sqrt}(\text{squared norm} + \text{epsilon})$ squashed vectors = scale * vectors.

C. Selection of model:

Metrics of performance are calculated on each of the models we discussed above for selection of good result yielding model. Metrics studied include Accuracy score, Precision Score, Recall Score, F1 score , Mean Squared Error (MSE) and ROC- AUC values.

D. Feature Importance Estimation:

We make use of the trained RandomForest model above and add upon an interface for predicting and testing the results. Tkinter library of python was used to build GUI. The GUI checks firstly whether an interaction exists between the drugs and if at all the interaction is found out then it plots the feature importance of the two interacting drugs. The GUI also outputs the highest feature importance.

CHAPTER 8

IMPLEMENTATION AND PSEUDOCODE

8.1 Random Forest

Random forest is a commonly-used machine learning algorithm which combines the output of multiple decision trees to reach a single result. Its ease of use and flexibility have fueled its adoption, as it handles both classification and regression problems.

```
from sklearn.ensemble import RandomForestClassifier
from sklearn.metrics import accuracy_score

# Train a Random Forest classifier
clf = RandomForestClassifier()
clf.fit(X_train, y_train)

# Predict on the test set
y_pred = clf.predict(X_test)

# Evaluate the model
accuracy = accuracy_score(y_test, y_pred)
print(f'Accuracy: {accuracy}')
```

Accuracy: 0.9536755110917791

```
from sklearn.metrics import accuracy_score, f1_score, precision_score, recall_score

# Predict on the test set
y_pred = clf.predict(X_test)

# Calculate F1 score, precision, and recall
f1 = f1_score(y_test, y_pred)
precision = precision_score(y_test, y_pred)
recall = recall_score(y_test, y_pred)

# Print the results
print(f'Accuracy: {accuracy}')
```

Fig 8.1 : Code for Implementation of Random forest

Uses in-built library i.e, RandomForestClassifier, then trains the model using training dataset. For prediction uses the test data. Calculates accuracy of the model and other metrics also.

8.2 SVM(Support Vector Machine)

A supervised learning algorithm that works by finding the optimal hyperplane that best separates data points of different classes in a high-dimensional space. SVM is particularly effective in scenarios with complex decision boundaries and high-dimensional feature spaces. The algorithm aims to maximize the margin, the distance between the hyperplane and the nearest data points of each class. This results in robust generalization to new, unseen data. SVM is versatile, accommodating linear and non-linear classification tasks through the use of different kernel functions.

```
from sklearn.svm import SVC

# Train an SVM model
svm = SVC()
svm.fit(X_train, y_train)

# Predict on the test set
y_pred_svm = svm.predict(X_test)

# Evaluate the model
accuracy_svm = accuracy_score(y_test, y_pred_svm)
print(f'SVM Accuracy: {accuracy_svm}')
```



```
#####
from sklearn.svm import SVC
from sklearn.metrics import accuracy_score

# Train an SVM model
svm = SVC()
svm.fit(X_train, y_train)

# Predict on the training and testing sets
y_pred_svm_train = svm.predict(X_train)
y_pred_svm_test = svm.predict(X_test)

# Calculate training and testing accuracy
accuracy_svm_train = accuracy_score(y_train, y_pred_svm_train)
accuracy_svm_test = accuracy_score(y_test, y_pred_svm_test)

print(f'SVM Training Accuracy: {accuracy_svm_train}')
```

8.2 Code for implementation of SVM

For implementing, we should use inbuilt library i.e, SVC. Then train model using training dataset. Predict output using testing data and then evaluate the model using different metrics.

8.3 XgBoost

XGBoost, which stands for Extreme Gradient Boosting, is a scalable, distributed gradient-boosted decision tree (GBDT) machine learning library. It provides parallel tree boosting and is the leading machine learning library for regression, classification, and ranking problems. Flexibility: XGBoost supports a variety of data types and objectives, including regression, classification, and ranking problems. Regularization: XGBoost incorporates regularization techniques to avoid overfitting and improve generalization performance.

```
from xgboost import XGBClassifier
from sklearn.metrics import accuracy_score, precision_score, recall_score, f1_score, roc_auc_score, roc_curve, auc
import matplotlib.pyplot as plt

# Train an XGBoost model
xgb = XGBClassifier()
xgb.fit(X_train, y_train)

# Predict on the test set
y_pred_xgb = xgb.predict(X_test)

# Accuracy
accuracy_xgb = accuracy_score(y_test, y_pred_xgb)
print(f'XGBoost Accuracy: {accuracy_xgb}')

# Precision
precision = precision_score(y_test, y_pred_xgb)
print(f'Precision: {precision}')

# Recall
recall = recall_score(y_test, y_pred_xgb)
print(f'Recall: {recall}')

# F1 Score
f1 = f1_score(y_test, y_pred_xgb)
print(f'F1 Score: {f1}')
```

[79]

```
... XGBoost Accuracy: 0.957155284906481
Precision: 0.9588208820882088
Recall: 0.9967251461988305
F1 Score: 0.9774056657873609
```

8.3 Code for implementation of XgBoost

For implementing , we should use in-built library i.e, XgBoostClassifier. Then train model using training dataset. Predict output using testing data and then evaluate the model using different metrics i.e, Accuracy, Precision, Recall, F1 Score.

8.4 XGBoost + AutoEncoders

An autoencoder could misclassify input errors that are different from those in the training set or changes in underlying relationships that a human would notice. Another drawback is you may eliminate the vital information in the input data. The XGBoost model outperforms the deep learning models on most datasets.

```
class GradientBoostingClassifier:
    def __init__(self, n_estimators=100, learning_rate=0.1):
        self.n_estimators = n_estimators
        self.learning_rate = learning_rate
        self.models = []

    def fit(self, X, y):
        y_pred = np.zeros(y.shape)
        for _ in range(self.n_estimators):
            residuals = y - y_pred
            tree = DecisionTreeRegressor(max_depth=3)
            tree.fit(X, residuals)
            update = self.learning_rate * tree.predict(X)
            y_pred += update
            self.models.append(tree)

    def predict(self, X):
        y_pred = np.zeros(X.shape[0])
        for tree in self.models:
            y_pred += self.learning_rate * tree.predict(X)
        return y_pred

# Train an XGBoost-like model
xgb = GradientBoostingClassifier(n_estimators=100, learning_rate=0.1)
xgb.fit(X_train, y_train)

# Predict on the test set
y_pred_xgb = xgb.predict(X_test)

# Combine predictions with encoded features
encoded_dim = 32 # Adjust the dimension based on your problem
encoder = tf.keras.Sequential([
    tf.keras.layers.Input(shape=(X_train.shape[1],)),
    tf.keras.layers.Dense(128, activation='relu'),
    tf.keras.layers.Dense(encoded_dim, activation='relu'),
])
decoder = tf.keras.Sequential([
    tf.keras.layers.Input(shape=(encoded_dim,)),
    tf.keras.layers.Dense(128, activation='relu'),
    tf.keras.layers.Dense(X_train.shape[1], activation='sigmoid'),
])
autoencoder = tf.keras.Sequential([encoder, decoder])
autoencoder.compile(optimizer='adam', loss='mean_squared_error')

# Fit the autoencoder on the training data
autoencoder.fit(X_train, X_train, epochs=20, batch_size=64, validation_split=0.1, verbose=2)

# Encode the features
encoded_features = encoder.predict(X_train)
encoded_features_test = encoder.predict(X_test)

# Train a simple classifier on top of the encoded features
classifier = tf.keras.Sequential([
    tf.keras.layers.Input(shape=(encoded_dim,)),
    tf.keras.layers.Dense(64, activation='relu'),
    tf.keras.layers.Dense(1, activation='sigmoid')
])
classifier.compile(optimizer='adam', loss='binary_crossentropy', metrics=['accuracy'])

# Fit the classifier on the encoded features
classifier.fit(encoded_features, y_train, epochs=20, batch_size=64, validation_split=0.1, verbose=2)

# Predict using the classifier
y_pred_classifier = classifier.predict(encoded_features_test)
y_pred_classifier = (y_pred_classifier > 0.6).astype(int)

# Calculate metrics
accuracy_combined = accuracy_score(y_test, y_pred_classifier)
precision = precision_score(y_test, y_pred_classifier)
```

8.4 Code for implementation of XgBoost + Autoencoder

For implementing, we should use in-built library i.e., GradientBoostingClassifier and used encoder and decoder. Then train model using training dataset. Predict output using testing data and then evaluate the model using different metrics i.e., Accuracy, Precision, Recall, F1 Score.

8.5 Capsule Network

In a capsule network, each capsule is made up of a group of neurons with each neuron's output representing a different property of the same feature. This provides the advantage of recognizing the whole entity by first recognizing its parts. The input to a capsule is the output (or features) from a CNN. Capsules leverage vectors for more detailed representation, rather than scalars.

```
import numpy as np
from sklearn.model_selection import train_test_split
from sklearn.preprocessing import StandardScaler
from tensorflow import keras
from tensorflow.keras.layers import Input, Conv1D, Flatten, Dense, Reshape, Lambda
from tensorflow.keras import backend as K

# Split your data into training and testing sets
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)

# Standardize your data
scaler = StandardScaler()
X_train = scaler.fit_transform(X_train)
X_test = scaler.transform(X_test)

# Define the Capsule Network architecture
input_dim = X_train.shape[1]

# Primary Capsules
x = Input(shape=(input_dim,))
x_resaped = Reshape(target_shape=(input_dim, 1))(x)
primary_capsules = Conv1D(filters=32, kernel_size=3)(x_resaped)

# Digit Capsules
def squash(vectors, axis=-2):
    s_squared_norm = K.sum(K.square(vectors), axis, keepdims=True)
    scale = s_squared_norm / (1 + s_squared_norm) / K.sqrt(s_squared_norm + K.epsilon())
    return scale * vectors

digit_capsules = Lambda(squash)(primary_capsules)

# Class Capsules
num_classes = 1 # Binary classification
class_capsules = Conv1D(filters=num_classes, kernel_size=input_dim - 2)(digit_capsules)
class_capsules_resaped = Reshape(target_shape=(num_classes,))(class_capsules)

# Define the Capsule Network model
capsule_model = keras.models.Model(x, class_capsules_resaped)

# Compile the model
capsule_model.compile(optimizer='adam', loss='binary_crossentropy', metrics=['accuracy'])

# Train the Capsule Network
capsule_model.fit(X_train, y_train, epochs=10, batch_size=64)

# Evaluate the model on the test data
accuracy = capsule_model.evaluate(X_test, y_test)
print(f'Accuracy: {accuracy[1]}')
```

8.5 Code for implementation of Capsule Network

The above figure represents the implementation of capsule network in python. It uses CNN and squash definition to scale the vectors. It compiles the model and evaluate the model. Calculates the model accuracy.

8.6 Convolutional Neural Network

CNN classifier for image classification is a CNN-based model specifically designed to classify images into different predefined classes. It learns to extract relevant features from input images and map them to the corresponding classes, enabling accurate image classification.

```
# Split your data into training and testing sets
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)

# Standardize your data
scaler = StandardScaler()
X_train = scaler.fit_transform(X_train)
X_test = scaler.transform(X_test)

# Define the Convolutional Neural Network architecture
input_dim = X_train.shape[1]

model = keras.Sequential()
model.add(Input(shape=(input_dim, 1)))
model.add(Conv1D(filters=32, kernel_size=3, activation='relu'))
model.add(Flatten())
model.add(Dense(64, activation='relu'))
model.add(Dense(1, activation='sigmoid')) # Output layer for binary classification

# Compile the model
model.compile(optimizer='adam', loss='binary_crossentropy', metrics=['accuracy'])

# Define a callback for ROC AUC calculation
class ROCCallback(keras.callbacks.Callback):
    def on_train_begin(self, logs={}):
        self.aucs = []

    def on_epoch_end(self, epoch, logs={}):
        if self.validation_data is None:
            return

        X_val, y_val = self.validation_data[0], self.validation_data[1]
        y_pred = self.model.predict(X_val)
        auc = roc_auc_score(y_val, y_pred)
        self.aucs.append(auc)
        print(f'Epoch {epoch + 1} - AUC: {auc:.4f}')

roc_callback = ROCCallback()

# Train the CNN model
history = model.fit(X_train, y_train, epochs=10, batch_size=64, validation_data=(X_test, y_test), callbacks=[roc_callback])

# Evaluate the model on the test data
y_pred = model.predict(X_test)
y_pred_binary = (y_pred > 0.5).astype(int)

accuracy = accuracy_score(y_test, y_pred_binary)
precision = precision_score(y_test, y_pred_binary)
recall = recall_score(y_test, y_pred_binary)
f1 = f1_score(y_test, y_pred_binary)

print(f'Testing Accuracy: {accuracy:.4f}')
print(f'Precision: {precision:.4f}')
print(f'Recall: {recall:.4f}')
print(f'F1 Score: {f1:.4f}')

Epoch 1/10
8/288 [=====] - 2s 4ms/step - loss: 0.2545 - accuracy: 0.9265 - val_loss: 0.2295 - val_accuracy: 0.9302
Epoch 2/10
```

8.5 Code for implementation of Convolution Neural Network

The above figure represents the python code for implementation of CNN. It uses convolutional layer and pooling layer, fully connected layer, dense layer. It uses adam as optimizer, RELU as activation layer.

CHAPTER 9

RESULTS AND DISCUSSION

9.1 Random forest

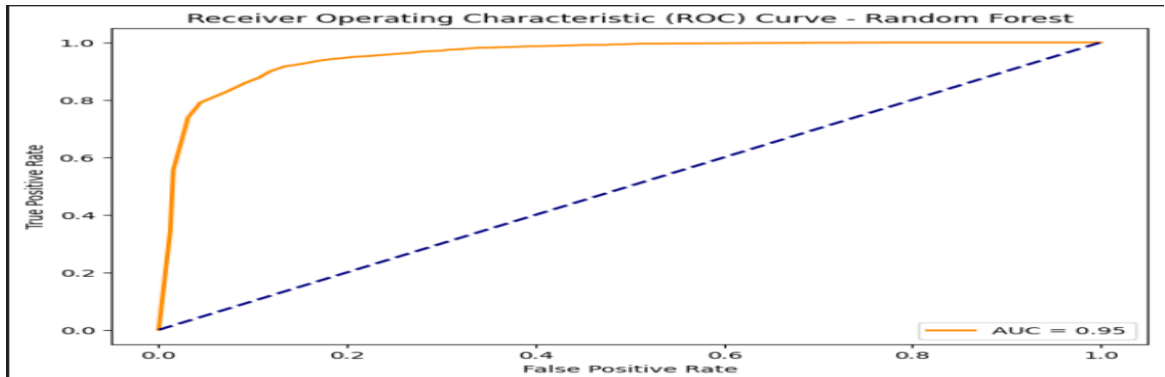


Fig. 9.1: ROC Curve for Random Forest

AUC of 0.95 indicates a very good model performance. i.e; it has high ability to distinguish between positive and negative instances.

9.2 SVM

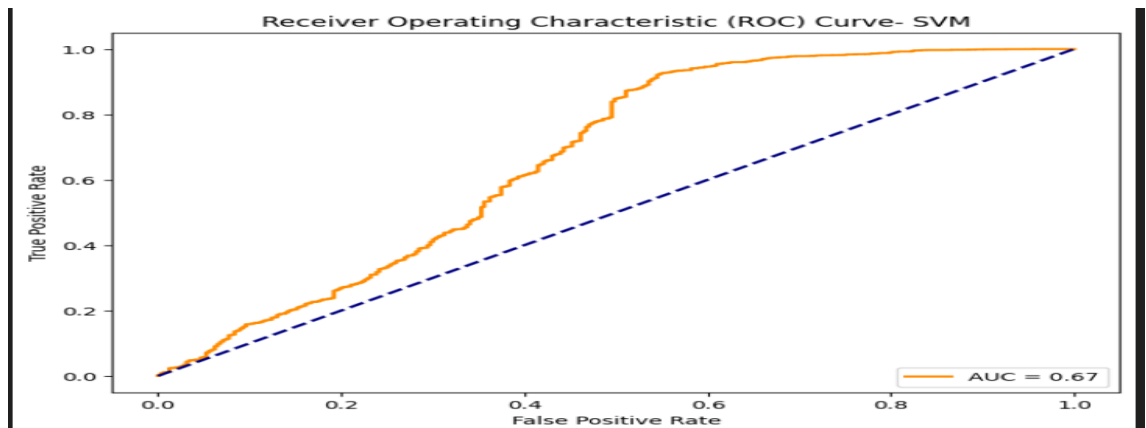


Fig. 9.2: ROC Curve for SVM

AUC of 0.67 indicates moderate performance. i.e, there is a good balance between true positive rate and false positive rate.

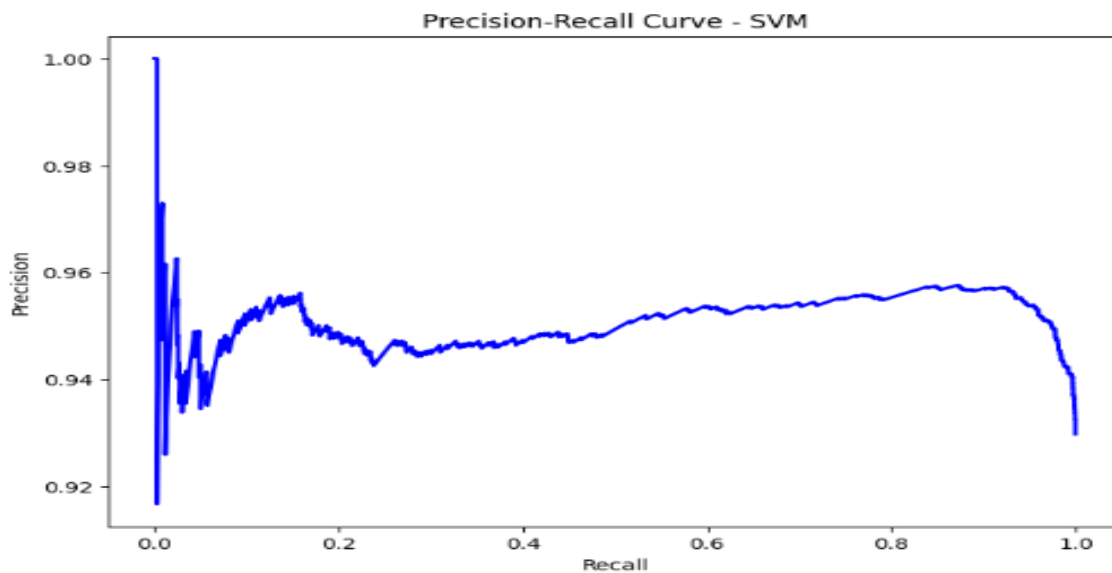


Fig. 9.3: Precision-Recall Curve for SVM

With a precision and recall of 93.3% and 99.9% respectively indicates that the instances predicted by the model as positive are indeed positives.

9.3 XGBoost

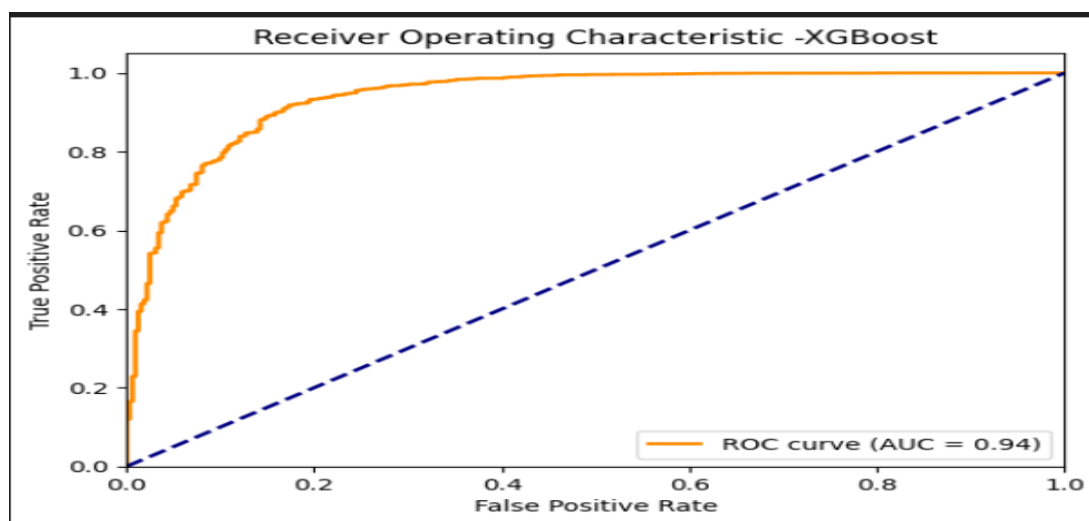


Fig. 9.4: ROC Curve for XGBoost

With an AUC of 0.94 this model is also a very good model like Random Forest Classifier. i.e, it has high ability to distinguish between positive and negative instances.

9.4 XGBoost + Autoencoders

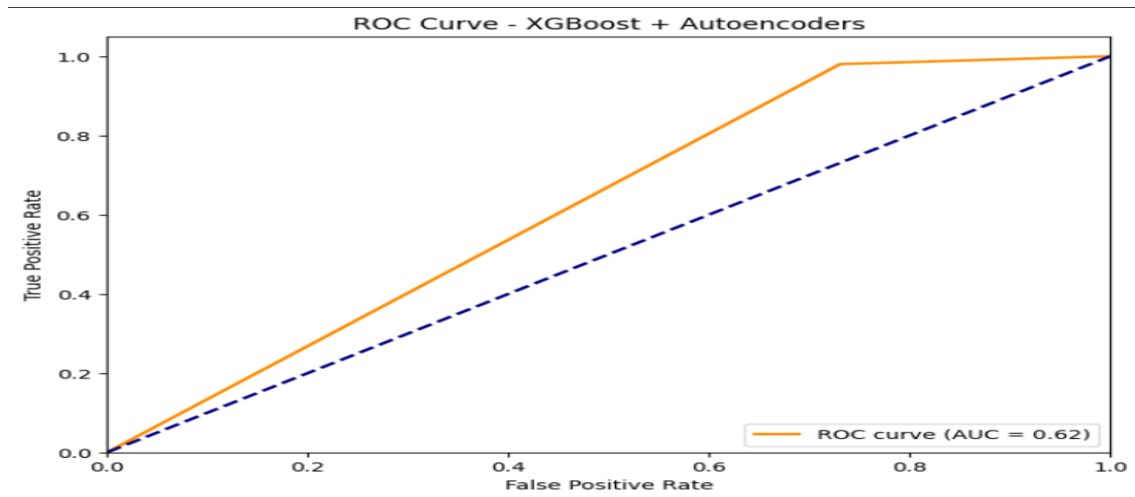


Fig. 9.5: ROC Curve for XGBoost + Autoencoders

AUC of 0.62 indicates lower performance compared to all the other models.

9.5 Capsule Networks

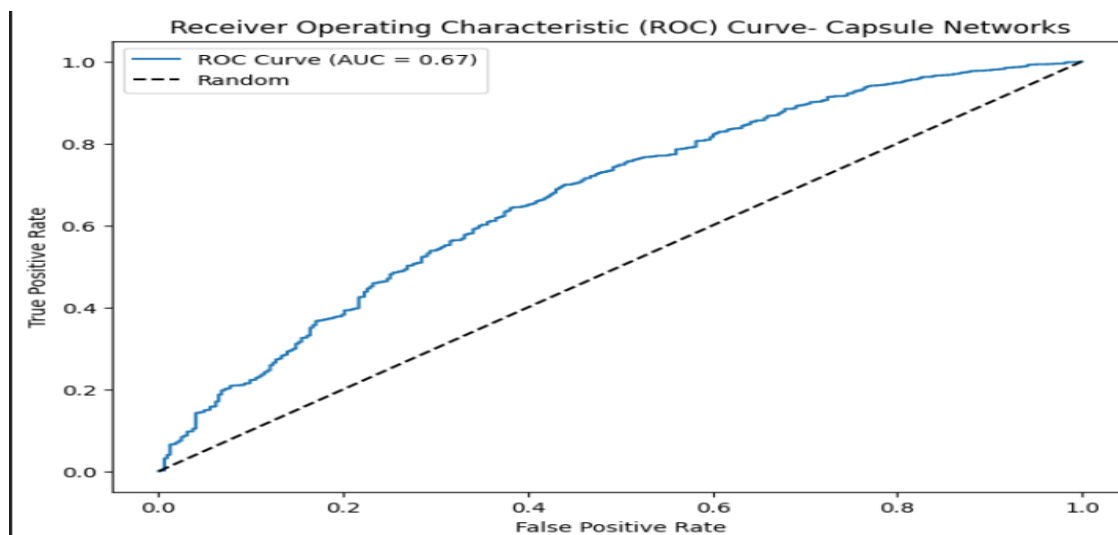


Fig. 9.6: ROC Curve for Capsule Networks

AUC of 0.67 is similar to SVM suggesting comparable performance.

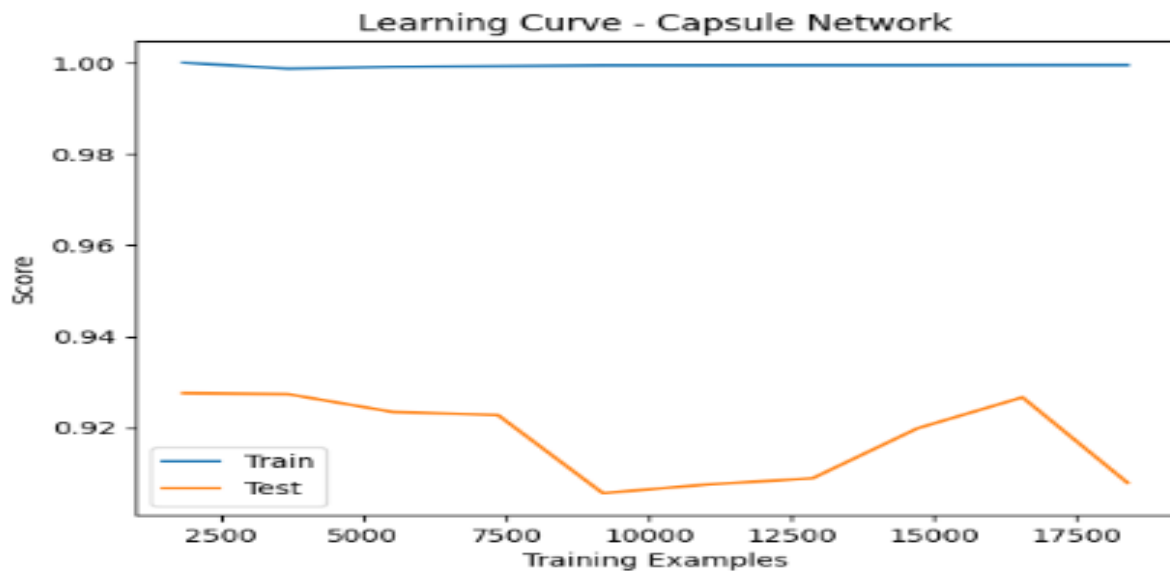


Fig. 9.7: Learning Curve for Capsule Networks

The results of the learning curve for capsule network in Fig 9.7 shows that the model is a strong performer and could be useful in applications where capturing positive instances is crucial.

9.6 CNN

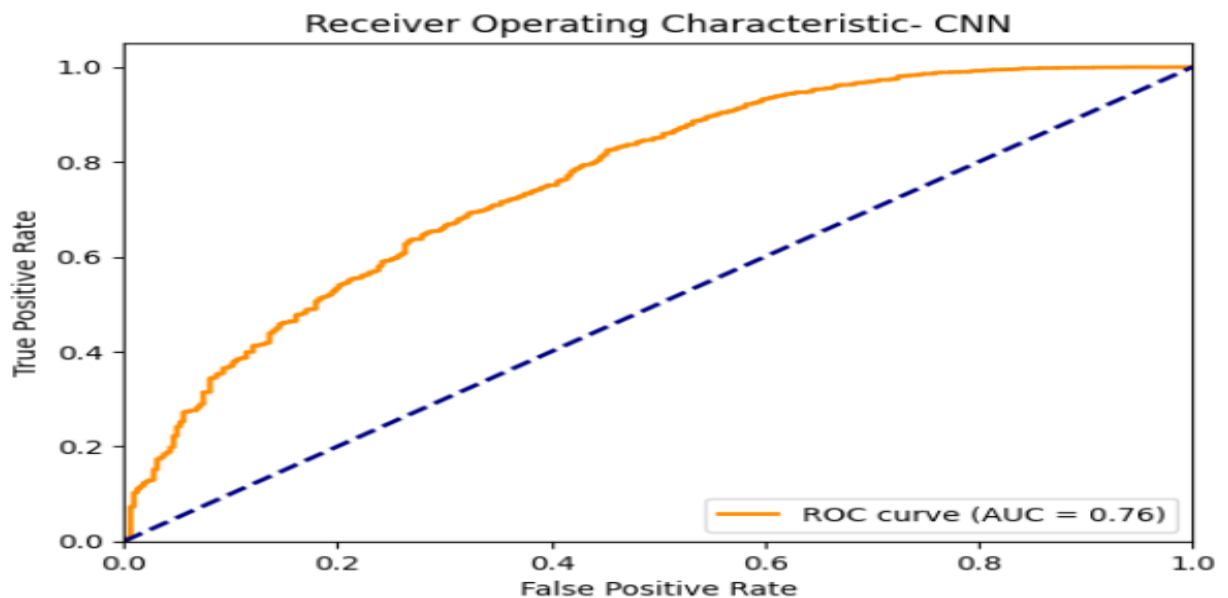


Fig. 9.8: ROC Curve for CNN

AUC of 0.76 suggests that the model performed reasonably well i.e; there is a good balance between true positive rate and false positive rate.

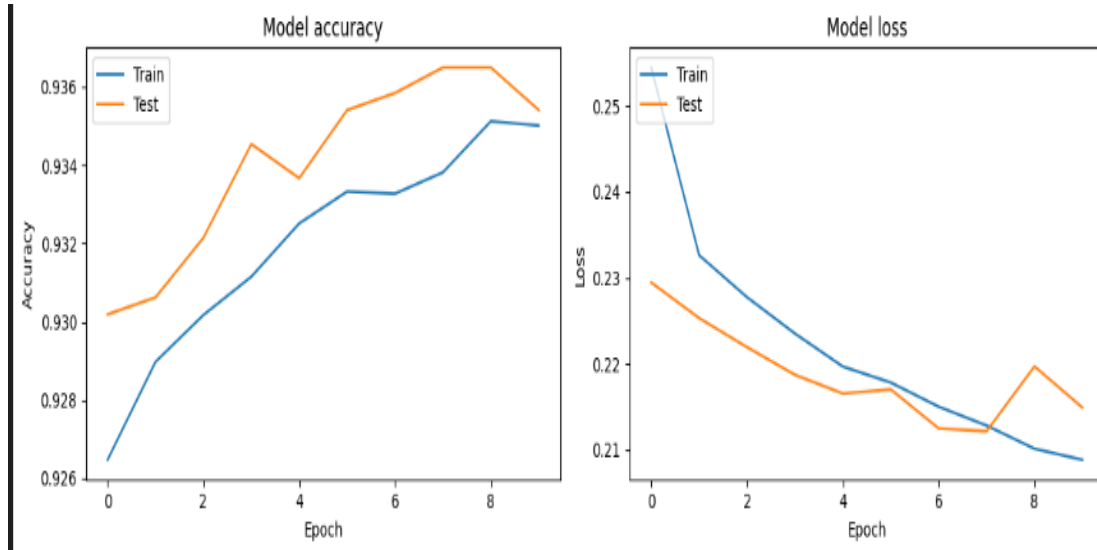


Fig. 9.9: Model accuracy and loss Curve for CNN

With an accuracy of 93.61% the correctly predicts the class labels for a substantial portion of the dataset.

9.7 Table for Comparison

Algorithm	Accuracy	Precision	Recall	F1 Score	ROC-AUC
Random Forest	95.60	95.63	99.83	0.97	0.96
<u>XGBoost</u>	95.71	95.88	99.67	0.97	0.94
SVM	93.32	93.31	99.97	0.96	0.67
Capsule Networks	92.90	92.97	99.92	0.96	0.67
CNN	93.61	94.03	99.44	0.96	0.77
<u>XGBoost+Autoencoder</u>	93.30	94.42	98.61	0.96	0.61

Fig. 9.10: Metrics analyzed for different algorithm

Now, let us make conclusions on the results that we have obtained.

Accuracy which provides understanding on overall correct prediction, is fundamental for any comparison. Here again we see that XGBoost and Random Forest Classifier performed slightly better than other models with accuracy values 95.7 and 95.4 respectively. However the models Capsule

Networks, CNN, SVM and AutoEncoders+XGBoost were good enough proving their values 93.9, 93.6, 93.3 and 95.7 respectively. Performance of all models were nearly equivalent.

Accuracy alone cannot justify the comparison. Hence other comparison metrics like Precision, F1 score and Recall also contribute to effective estimation.

Now let us analyze the **Precision score** among models, which measures the accuracy of the positive predictions made by the model. It is the ratio of true positive predictions to the total number of positive predictions made by the model. Capsule network predicted accurately with the score of 99.8 followed by XGboost and Random Forest Classifier with scores of 95.8 and 95.5 respectively. XGboost+Autoencoders, CNN and SVM also predicted well with precision values 94.7, 93.8 and 93.3 respectively.

Recall also goes hand in hand with precision metrics which measure out of relevant features how many did we get right? Here we see that SVM, Random Forest Classifier, Capsule Networks and CNN played its game well with scores of 99.9, 99.8, 99.8 and 99.7 respectively. XGBoost and XGBoost +Autoencoders were in no lag with values 95.8 and 94.7 respectively.

It also becomes essential to calculate **F1 score** when we have both Recall and Precision score. The results here also shows that performance of models were nearly equivalent. XGBoost, Random Forest Classifier, CNN, Capsule Network and SVM had their scores 97.7, 97.6, 96.7, 96.3 and 96.3 respectively. XGBoost + Autoencoders had a score of 95.8.

Receiver Operating Characteristic Curve (ROC) which provides insights on true positive rate and false positive rates was plotted for all models and was found that Random Forest Classifier and XGBoost performed better with AUC 0.95 and 0.94 respectively. The ROC for models CNN, SVM, Capsule Network and AutoEncoders + XGboost performed above satisfactory with AUC values 0.76, 0.67, 0.67 and 0.62 respectively.

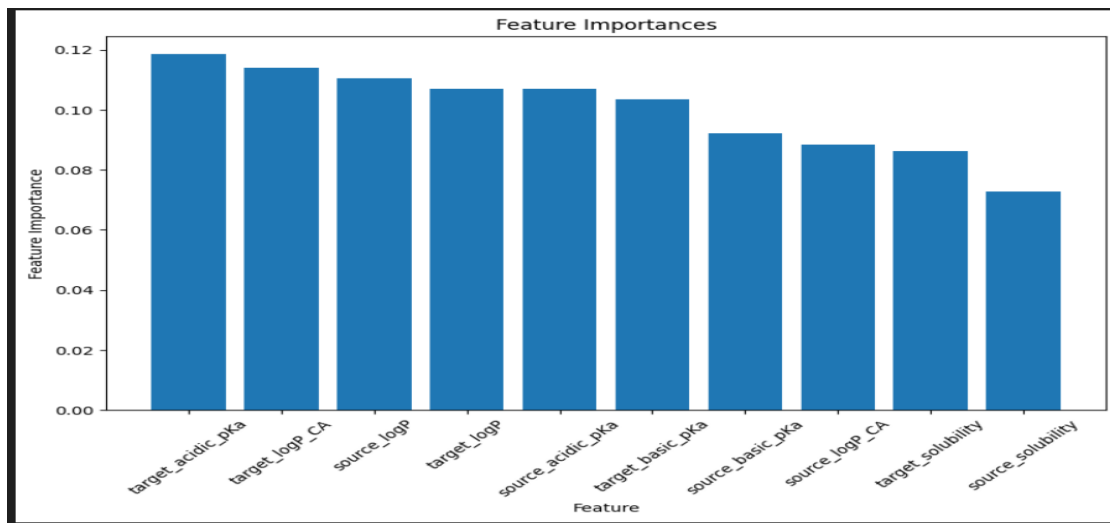


Fig. 9.11: importance of features

At last to predict the importance of features, we used Random Forest Classification as the model proved better prediction and the calculated result is displayed above.

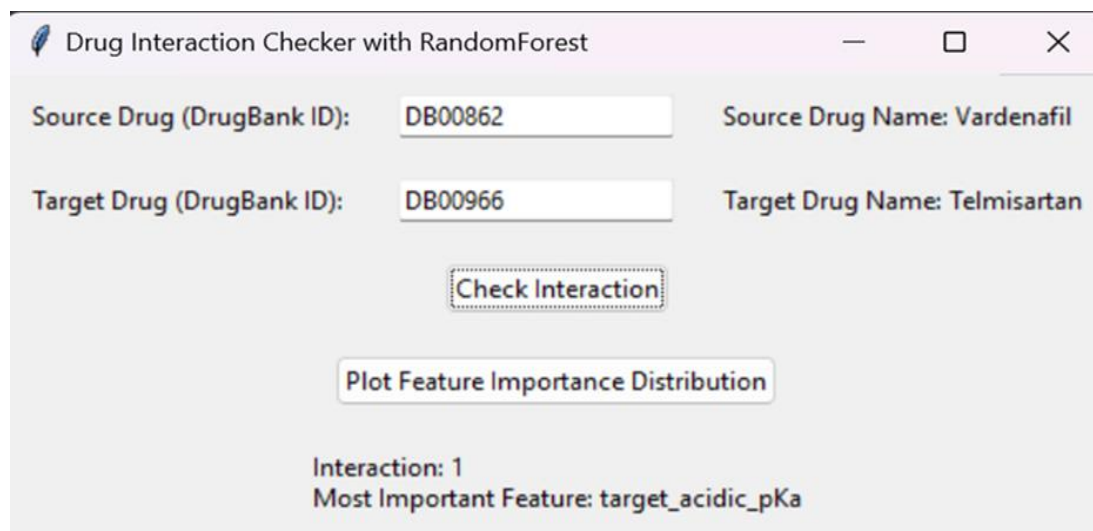


Fig 9.12: GUI of our system when interaction is present.

The above figure represents it will take input as two drugbank_id's of drugs. When you will click on check interaction, it will print 1 as interaction is present between them. And returns most important feature which causes interaction between those two drugs.

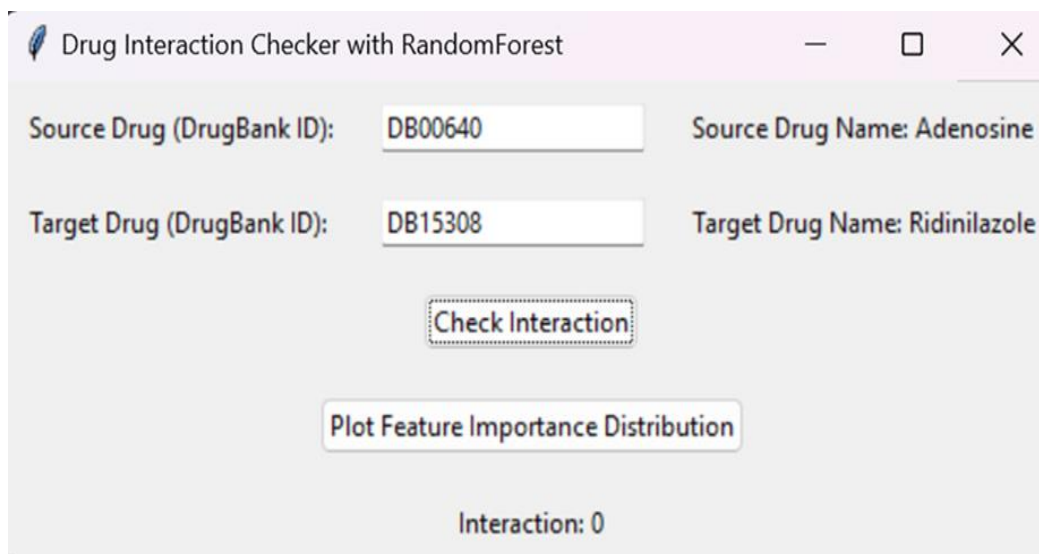


Fig 9.13: GUI of our system when there is no interaction.

The above figure represents it will take input as two drugbank_id's of drugs. When you will click on check interaction, it will print 0 as interaction is not present between them. And it does not returns most important feature as there is no interaction between them.

REFERENCES

- [1] Xiaoqiang Xu; Ping Xuan; Tiangang Zhang; Bingxu Chen; Nan Sheng, "Inferring Drug-Target Interactions based on Random Walk and Convolutional Neural Network" in IEEE/ACM Transactions on Computational Biology and Bioinformatics, Volume: 19, Issue: 4, 01 July-Aug. 2022.
- [2] Shaghayegh Sadeghi , Alioune Ngom, "DDI PRED:Graph Convolutional Network based Drug-Drug Interactions Prediction using Drug Chemical Structure Embedding" in IEEE, 2022.
- [3] Nilay Fatma , Yıldız Alper Özcan, "Graph Convolutional Autoencoder and Generative Adversarial Network-Based Method for Predicting Drug-Target Interactions " published in 2022.
- [4] Xu Gong, Maotao Liu, Haichao Sun, Min Li, Qun Liu, "HS-DTI: Drug-Target Interaction Prediction based on Hierarchical Networks and Multi-Order Sequence Effect" published in 2022.
- [5] Tianjung , Wang , Xin , Liu, "A Graph Convolution-Transformer Neural Network for Drug-Target Interaction Prediction", ICBBT 2022, May 27–29, 2022.
- [6] Chang Sun, Ping Xuan , Tiangang Zhang and Yilin Ye, "Using Supervised Machine Learning Algorithms for Drug-Target Interaction Prediction" in IEEE/ACM Transactions On Computational Biology And Bioinformatics, Vol. 19, No. 1, January/February 2022.
- [7] Mohammad A. Rezaei , Yanjun Li , Dapeng Wu, Xiaolin Li and Chenglong Li, "Deep Learning in Drug Design : Protein-Ligand Binding Affinity Prediction" in IEEE/ACM Transactions on Computational Biology and Bioinformatics, Volume: 19, Issue: 1, 01 Jan.-Feb. 2022.

- [8] Nelson R. C. Monteiro , Bernardete Ribeiro , and Joel P. Arrais, "Drug Target Interaction Prediction :end-to-end Deep Learning Approach" in IEEE/ACM Transactions on Computational Biology and Bioinformatics, Volume: 18, Issue: 6 , 01 Nov.-Dec. 2021.
- [9] Liu S., An J., Zhao J., Zhao S., Lv, H. & Wang S, "Drug-Target Interaction Prediction based on Multisource Information Weighted Fusion" from Drug-Target Interaction Prediction Based on Multisource Information Weighted Fusion. Contrast media & molecular imaging, 6044256, 2021.
- [10] Chengkun Wu , Wei Wang , Xi Yang, Canqun Yang, "Drug-Drug Interaction extraction from Biomedical Texts based on Multi-Attention Mechanism", ICBRA, 2021.
- [11] Mongia A, Majumdar A, "Drug-Target Interaction Prediction using Multi Graph Regularized Nuclear Norm Minimization", from PLoS ONE 15(1): e0226484, 2020.
- [12] Xinyu Hao , Jiaying You , PingZhao Hu, "Predicting Drug-Drug Interactions using Deep Neural Networks", ICMLC '19, February 22–24, 2019.

APPENDIX A

ACRONYMS AND ABBREVIATIONS

1. **SMILES** - Simplified Molecular Input Line Entry System
2. **SELFIES** - Self-Referencing Embedded Strings
3. **DDI's** - Drug Drug Interactions
4. **HS-DTI** - Heterogeneous drug target interaction
5. **GCN** - Graph Convolution Network
6. **AUROC** - Area Under the Receiver Operating Characteristics
7. **AUPRC** - Area under Precision-Recall curve
8. **SVM** - Support vector machine
9. **CNN** - Convolutional Neural Network
10. **ML/DL** - Machine learning / Deep learning
11. **MGRNNM** - Multi-Graph Regularized Nuclear Norm Minimization
12. **WNN-GIP** - Weighted Nearest Neighbours – Gaussian interaction profile