

Drug Discovery on Alzheimer Disease Using Machine Learning

Parth Bagal

*Dept of CSE DS KIT's
College of Engineering
(Autonomous)*

Kolhapur, Maharashtra, India
parthbagal20@gmail.com

Shreyas Bagave

*Dept of CSE DS KIT's
College of Engineering
(Autonomous)*

Kolhapur, Maharashtra, India
shreyashbagave95@gmail.com

Prathamesh Nale

*Dept of CSE DS KIT's
College of Engineering
(Autonomous)*

Kolhapur, Maharashtra, India
prathameshnale7515@gmail.com

Prathamesh Garate

*Dept of CSE DS KIT's College of Engineering
(Autonomous)*

Kolhapur, Maharashtra, India prgarate15@gmail.com

Abstract—Alzheimer's drug discovery was so far constrained to focusing on well-studied therapeutic hypotheses. A more diverse, systems-integrated approach may uncover new therapeutic hypotheses based on the heterogeneity of AD processes and the requirement for understudied protein targets and biological mechanisms. Target enabling packages are the development of high quality experimental reagents and informatic outputs and would accelerate the rapid evaluation of new emerging systems-integrated targets in AD. Drug discovery has witnessed dramatic change recently. Many AI/ML technologies have been applied in the recent past. The enhanced nature of these models has led to a growing need for transparency and interpretability. In this paper, we make a review of the XAI approach, which provides a novel means to achieve this goal for a more interpretable understanding of machine learning predictions. To improve the classification performance of a problem by overcoming the limitations of traditional statistical models and conventional machine learning techniques in handling complex molecular datasets. To do this, a dataset of 7298 compounds extracted from the ChEMBL database along with molecular descriptors were used

Index Terms—Alzheimer's disease, drug development, drug targets

I. INTRODUCTION

Alzheimer's disease is one of the most sensitive medical challenges of our time: the devastating cognitive decline of these patients and the complete loss of their memories and identities. Increasingly, as populations age in every country and continent, the need for effective treatments intensifies, striding researchers through the complicated journey of drug discovery.

At the heart of Alzheimer's disease, the accumulation of beta-amyloid plaques and tau protein tangles in the brain leads to neuronal dysfunction and eventually cell death. Although we have gained much insight into the nature of this disease,

translating this knowledge into effective therapies poses a major challenge.

Drug discovery in Alzheimer's disease is multi-faceted, from target identification to the direct clinical evaluation of therapeutic candidates in patients. Early efforts typically focus on identifying the cellular and molecular mechanisms behind neurodegeneration and therefore identifying key molecular targets that are implicated in the disease pathology. With novel approaches such as genomics, proteomics, and advanced imaging techniques, this search unravels the complex molecular landscape of Alzheimer's and elucidates drug-targetable processes.

When potential targets are identified, the next step is the development and optimization of therapeutic compounds that aim to modulate these targets. This is a delicate dance in biology: on one hand, efficacy and safety are the main parameters; on the other hand, pharmacokinetic properties of the compounds need to be considered. Working in an intricate playground of biological systems in the brain, the researchers test drug candidates in cell and animal models to assess their efficacy and safety profile.

Indeed, despite the challenges of clinical trials, the field of Alzheimer's drug discovery has witnessed unprecedented collaboration and innovation. From traditional small molecules to biologics and gene therapies, drug discovery researchers have explored a range of promising therapeutic approaches in the fight against this neurodegenerative course.

This paper explores the complexities of drug discovery in Alzheimer's disease, discussing recent advances, existing challenges, and future directions in the field. We hope that by opening the black box of how things happen from bench to bedside, we can infuse hope into our readers and stimulate efforts to find effective therapies against this robust adversary.

II. RELATED WORK

Alzheimer’s disease (AD) is a neurological illness that worsens with time and is marked by functional impairment, memory loss, and cognitive decline. Adverse disease treatments that

work are still elusive despite substantial research efforts. The purpose of this review of the literature is to give a broad picture of the state of drug discovery for AD, emphasising important obstacles, viable strategies, and potential paths forward.I.

Aspect	Description
Challenges in AD Drug Discovery	
Slow and expensive process	Traditional drug discovery for AD is slow, expensive, and has high failure rates in clinical trials. .
Limited success in clinical trials	Clinical trials of novel drugs for AD have yielded limited success.
Data Quality and Availability	Limited access to high-quality data impedes ML-based drug discovery for AD. Efforts are needed to enhance data sharing and standardization.
Model Interpretability	Interpretability issues with ML models raise concerns about reliability and trustworthiness. Developing interpretable algorithms and explanation techniques is crucial
Biological Complexity	AD is multifactorial, involving complex interactions. Integrating diverse data sources and biological knowledge into ML models is crucial for accuracy.
The Promise of Machine Learning	
Drug target identification	ML can analyze vast datasets to identify potential drug targets and repurpose existing drugs for AD treatment.

Techniques used	Support vector machines, deep neural networks, and recurrent neural networks are employed for bioactivity prediction and target identification.
Analysis of biological mechanisms	ML can dissect complex biological mechanisms underlying AD, leading to more targeted therapies.
Key Applications	
Drug Repurposing	ML models identify existing drugs with potential to treat AD based on molecular properties and disease associations.
Target Identification	ML analyzes gene expression data and protein-protein interactions to identify novel drug targets in AD pathogenesis.
In-silico Drug Design	ML virtually screens vast libraries of compounds to predict those with potential therapeutic effects for AD
Benefits	
Streamlining drug discovery	ML reduces time and cost in the drug discovery process.
Exploration of biological pathways	ML facilitates exploration of intricate biological pathways in AD
Limitations and Future Directions	
Data quality dependence	ML models heavily depend on the quality and comprehensiveness of training data.

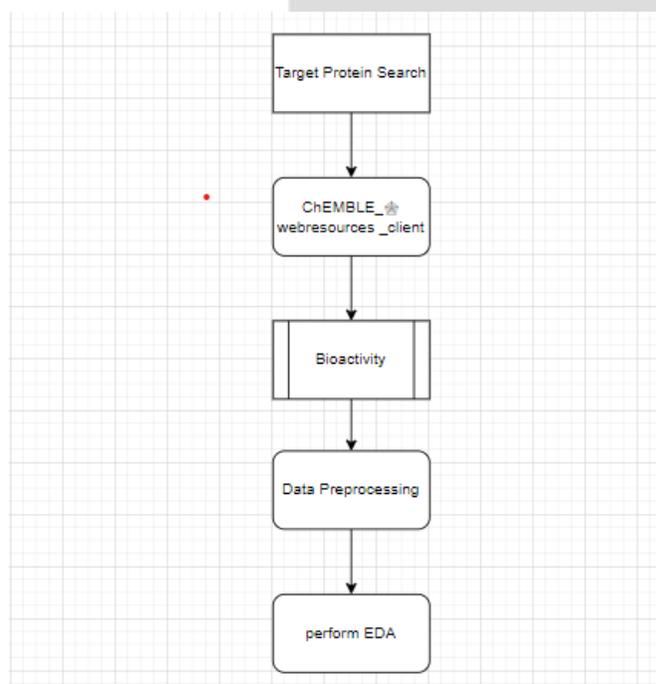
Need for validation	In-silico predictions require validation through further wet-lab experiments.
Exploration of biological pathways	ML facilitates exploration of intricate biological pathways in AD

TABLE I : Literature Review

III PROPOSED METHOD

Data collection

Alzheimer's disease (AD) is a multifaceted neurodegenerative condition marked by memory loss and a steady deterioration in cognitive function. Adverse disease treatments that work are still elusive despite substantial research efforts. This work offers a thorough method for gathering data for AD medication discovery using the ChEMBL library, a useful tool for obtaining bioactivity data and chemical details on compounds examined against pertinent AD-related protein targets. Researchers can find a plethora of information by using the ChEMBL library to find promising drug candidates, rank compounds for additional testing, and learn more about the molecular pathways behind AD pathogenesis.



In this work, we gathered information pertinent to AD medication discovery using the ChEMBL collection. First, we used terms like "Alzheimer" and "amyloid beta" to scan the ChEMBL database for protein targets linked to AD. After that, we gathered pertinent chemical data, such as compound structures and bioactivity data, from the compounds that had been evaluated against these targets. Prioritising compounds for additional research and identifying possible therapeutic

candidates were achieved through the processing and analysis of the gathered data.

A comprehensive dataset of chemicals evaluated against protein targets associated to AD was obtained from the ChEMBL library as a result of our data collection efforts. Promising candidates with favourable bioactivity profiles and structural features were identified by dataset analysis for future assessment in preclinical and clinical research. Furthermore, the molecular pathways underpinning AD pathogenesis were clarified by our investigation, which made it easier to identify new treatment targets.

A. EDA

PaDEL Descriptor: A software programme called PaDEL (Prediction of Activity Spectra for Substances) computes a variety of molecular descriptors for chemical substances.

Chemical compounds are quantitatively represented by molecular descriptors, which encompass a range of structural, physicochemical, and electrical properties. More than 1800 descriptors, such as topological, quantum chemical, constitutional, geometrical, and electrotopological descriptors, can be computed using the PaDEL Descriptor. These descriptors are employed in virtual screening, compound clustering, and quantitative structure-activity relationship (QSAR) modelling, among other tasks.

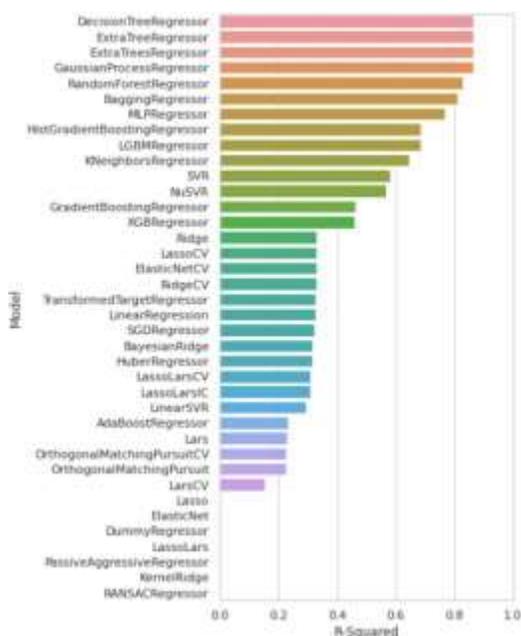
article graphicx

Fingerprints: Chemical compounds can be represented as binary or bit-string representations in molecular fingerprints, which compactly encode structural information. Usually, they are produced by encoding a molecule's existence or lack of specific chemical substructures or characteristics. In cheminformatics and drug development, fingerprints are frequently utilised for machine learning, grouping, virtual screening, and similarity finding.

Lipinski's Rule of Five (Ro5): The Lipinski's Rule of Five (Ro5) was put forth by Christopher A. Lipinski in 1997 and is an often followed guideline in medicinal chemistry and drug discovery. It aids in determining whether a chemical substance is drug-like based on its physicochemical characteristics. According to Lipinski's Rule of Five, Molecular Weight (MW): The chemical should have a molecular weight of no more than 500 daltons. Lipophilicity (LogP): Less than five should be the computed octanol-water partition coefficient (LogP). Hydrogen Bond Donors (HBD): There should be less than or equal to five hydrogen bond donor groups (such as the NH and OH groups). Hydrogen Bond Acceptors (HBA): Less than or equal to ten hydrogen bond acceptor groups, such as N and O atoms, should exist.

Calculate Lipinski descriptors: Do a Lipinski descriptor calculation. Pfizer scientist Christopher Lipinski developed a set of guidelines for determining whether a molecule is druglike or not. The pharmacokinetic profile, often referred to as the Absorption, Distribution, Metabolism, and Excretion (ADME), serves as the foundation for this druglikeness. Lip-

Model Selection



Model:- Random Forest Regression model

Suitable for regression problems, the Random Forest Regression model is a potent machine learning technique that predicts continuous outcomes. During training, a large number of decision trees are built using this ensemble learning technique, which generates the average forecast of each individual tree. A thorough explanation of the Random Forest Regression model may be found here:

Ensemble Learning: Random Forest Regression is a type of ensemble learning, which integrates several separate models to create a final model that is more reliable and powerful. Decision trees are the individual models in the Random Forest scenario.

Bootstrap Aggregating (Bagging): Random Forest trains each decision tree using a bootstrapped sample of the training data, a process known as bagging. The variety that is introduced among the trees by this method helps to minimize overfitting. Selecting Features at Random: Apart than utilizing bootstrapped

III. CONCLUSION

This conclusion highlights the key findings of our project are ChEMBL as a Valuable Resource: This research leveraged the ChEMBL database to efficiently extract raw data for target protein acetylcholinesterase. This highlights the importance of ChEMBL as a public resource for facilitating drug discovery projects. Lazypredict for Model Comparison: The utilization of the Lazypredict library for model comparison is a noteworthy aspect. It facilitated the evaluation of various models based on accuracy (identifying the best algorithm) and root mean squared error (RMSE) to assess prediction quality. This demonstrates a data-driven and rigorous approach to model selection. Random Forest Regression: Superior Performance: The research identified the Random Forest regression model as the optimal choice based on its impressive 86% accuracy and low RMSE. This emphasizes the model's effectiveness for predicting the target property (likely binding affinity) for acetylcholinesterase inhibitors. Beyond Accuracy: Descriptor and Fingerprint Generation: This project went beyond simply identifying the best model. The generation of descriptors and fingerprints using the chosen model opens doors for further virtual screening and lead optimization efforts.

This research project demonstrates the successful application of a multi-faceted in silico approach for identifying promising acetylcholinesterase inhibitors. By leveraging ChEMBL, PaDEL descriptors, Lipinski's rule, Lazypredict, and Random Forest regression, the project provides valuable insights for future drug discovery efforts targeting this crucial enzyme.

Prediction output

	molecule_name	pIC50
0	(CHEMBL1870889)	5.1726
1	(CHEMBL1809052)	5.1815
2	(CHEMBL1629979)	6.109
3	(CHEMBL1626128)	5.1726
4	(CHEMBL1618111)	5.1388

Fig: OUTPUT

REFERENCE

S

- [1] A. Gholami, "Alzheimer's disease: The role of proteins in formation, mechanisms, and new therapeutic approaches," *Neurosci. Lett.*, vol. 817, p. 137532, Nov. 2023, doi: 10.1016/j.neulet.2023.137532.
- [2] A. Gustavsson et al., "Global estimates on the number of persons across the Alzheimer's disease continuum," *Alzheimer's Dement.*, vol. 19, no.2, pp. 658–670, Feb. 2023, doi: 10.1002/alz.12694.
- [3] Bellenguez C, Kucukali F, Jansen IE, et al. New insights into the genetic etiology of Alzheimer's disease and related dementias. *Nat Genet.* 2022;54:412-436.
- [4] Wingo AP, Liu Y, Gerasimov ES, et al. Integrating human brain proteomes with genome-wide association data implicates new proteins in Alzheimer's disease pathogenesis. *Nat Genet.* 2021;53:143-146
- [5] T. R. Noviandy et al., "Ensemble Machine learning Approach for Quantitative Structure Activity Relationship Based Drug Discovery: A Review," *Infolitika J. Data Sci.*, vol. 1, no. 1, pp. 32–41, Sep. 2023, doi: 10.60084/ijds.v1i1.91.
- [6] Arshad Z, Smith J, Roberts M, et al. Open access could transform drug discovery: a case study of JQ1. *Expert Opin Drug Discov.* 2016;11:321-332.

