# Identification of Drug Compound Bio-Activities Through Artificial Intelligence

Rohit Rastogi, ABES Engineering College, India*

https://orcid.org/0000-0002-6402-7638

Yash Rastogi, ABES Engineering College, India

Saurav Kumar Rathaur, ABES Engineering College, India

Vaibhav Srivastava, ABES Engineering College, India

## ABSTRACT

In the fields of drug discovery and development, machine learning techniques have been used for the development of novel drug candidates. The methods for designing drug targets and novel drug discovery now routinely combine machine learning algorithms such as regression and classification models to enhance the efficiency, efficacy, and quality of developed outputs. Applying machine learning model for drug discovery on different diseases that exists already, the author team fetched the datasets from the ChEMBL database that contain the bio-activity data, after preprocessing the data according to the bioactivity threshold in order to obtain a curated bio-activity data. Therefore, structural analogs of the drugs that bind to the target are selected as drug candidates. However, even though compounds are not structural analogs, they may achieve the desired response. A new drug discovery method based on drug response, which can complement the structure-based methods, is needed. Present manuscript is an effort for same.

## KEYWORDS:

Bio-Activity, canonical_smiles, ChemBL, Chemical Compounds, Drugs, Machine Building, Protein, QSAR, standard_values, Web Service

## MOTIVATION

For every pharmaceutical company, creating a medicine and determining its bioactivity are more difficult tasks. To create an efficient drug, selling it, and keeping it efficient and successful need significant time, effort, and financial investment on the part of these corporations. Therefore, the research team came up with the concept to use our machine learning model to simplify this difficult task while saving time and resources.

## SCOPE OF THE STUDY

ML algorithm requires the large dataset to train the model for which different diversified sources can provide the data from which required data are extracted and used in training the model. The Author and his team are fetching the relevant dataset based on the target protein (acetylcholinesterase) and some required features through the cHembl database.

## TOPIC ORGANIZATIONS

The motive of this paper is to take the attention of the reader to how the infrastructure of pharmaceutical companies works on new technologies. How the modern world tackled the failure of the formation of vaccines in the phase of covid–19 viruses. For the underwriting of the review, the author team did a literary survey and explored ten research papers concerning points. This literary survey provides deep knowledge about drug discovery using various methods, such as machine learning and Bio-Activity detection.

The author and his team looked at how the manufacturing process for generating medications and preventative chemicals takes longer for each ailment. The researchers went on to explain that in order to successfully complete the medication development phase, target identification must come first. In addition to concentrating on some target variables like proteins, possible drugs (Bio-Activity), and DNA mutations, this phase must concentrate on the key chemicals (chemical compounds). For both machines and people, creating pharmaceuticals is not an easy task. creating a chemical compound combination and creating a potent and effective medicine.

The author team has described the procedure in which they have addressed the techniques utilized for the review. This study utilized the ML and AI-based basic examination of the data collected from the CHEMBL Database. The author and his team take the relevant data which is useful for their analysis and create a dataset and use it in their model.

### Ethical Committee and Funding

The experiments don't include any human-related experiments and so no ethical constraints have been violated. Though the subjects performing the study were humans and air quality directly affects them but the study doesn't violate any health-related measures. The Project is not funded by any agency.

## ROLE OF AUTHORS

Dr. Rohit Rastogi acted as team leader and coordinated among all co-authors. He prepared the topic introduction and background study; also contributed to experiments. He also prepared the structure of the manuscript and ensured the quality of the content along with all co-authors. Mr. Yash Rastogi did the data analysis. Mr. Vaibhav Srivastav did the experimental Analysis along with concluding remarks. Mr. Saurav Kumar Rathaur contributed to the results and discussions along with concluding remarks.

## 1. INTRODUCTION

In the current scenario, India's health department faces a double challenge. The first one is that they have to wrestle with diseases such as diarrhea, and lower respiratory infections(Asthma, Pneumonia, etc.) and another one is they have to tackle non-communicable conditions like cancer, and diabetes. There is no permanent medicine and treatment for these diseases available. For any disease, medicine takes a long time to manufacture and takes enough resources. Using a machine learning Author team builds a model which helps to manufacture or test the Bioactivity of drugs easily.

## 1.1 Indian and Global Health Scenario

Health is one of the most indispensable and significant areas in any sort of economy. There are many investigations that were directed at the various parts of health benefits and related regions. The world is confronting a 'triple weight of sickness' comprising of the fragmented arrangement of adaptable diseases, as of late emerging and returning contaminations as well as the remarkable rising of noncommunicable continuous sicknesses. The factors which assist progress and improvement nowadays with preferring globalization of trade, urbanization, effortlessness of overall travel, popular developments, etc, go probably as a two-sided bargain as they lead to positive wellbeing results on one hand and augmentation the shortcoming to persistent slightness of course as these add to fixed lifestyles and unfortunate dietary examples.

There is a high ordinariness of tobacco use close by extension in undesirable dietary practices and decreasing in real work adding to an expansion in regular bet factors which in this manner prompts an expansion in noncommunicable diseases (NCD). Below Figure1 delineates how life-related issues are increasing in NCDs (Kumar, S. et al., 2014) (pl. Refer Figure 1).

(Source: https://www.ijcm.org.in/article.asp?issn=0970-0218;year=2012;volume=37;issue=1; spage=5;epage=12;aulast=Kumar)

URL: https://www.ijcm.org.in/article.asp?issn=0970-0218;year=2012;volume=37;issue=1;sp age=5;epage=12;aulast=Kumar

## 1.2 Need of Drug Discovery and Healthcare Analysis

Scientists experimented on the drug detection and founded out different approaches to develop the drug in the health and clinic field and tested on humans and it might have to go through many clinical trials, in which much time are waste and also the resulted drug cannot give precise output. The Need of automated model required that applied the machine learning model that predict the effective output in very short time for diseases. In the process of manual drug discovery is very tedious task to work on and hence there should a drug detecting model should be present (Patel, L., et al., 2020) (pl. Refer Figure 2).

(Source:-https://www.ncbi.nlm.nih.gov/pmc/articles/PMC5725284/figure/fig5/)

## 1.3 Chemical Compounds and Variation in Drugs through Quantity

The process of drug formation begins with the development of a medical particle that has shown helpful value to fight, control, check or fix illnesses. The union and portrayal of such atoms which are likewise called active pharmaceutical ingredients (APIs) and their examination to make starter security and medicinal ampleness data are necessary to identify prescription opportunities for extra low down assessments.(Rastogi, R., et. al., 2021).

**Figure 1. How Lifestyle-Related Issues Contribute to an Increase in Noncommunicable Diseases**
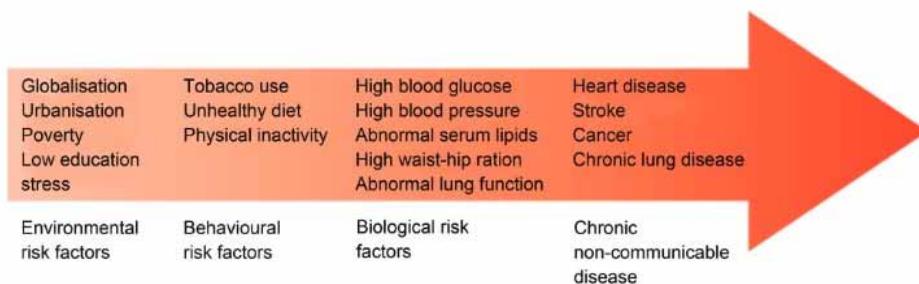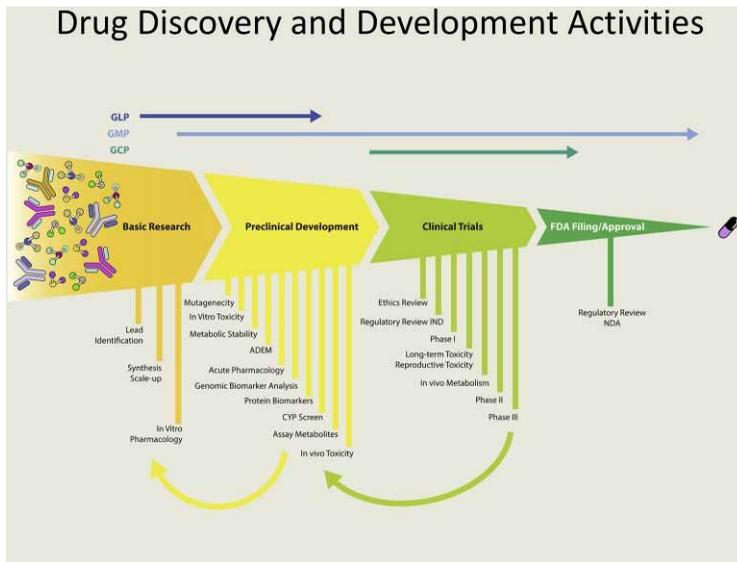
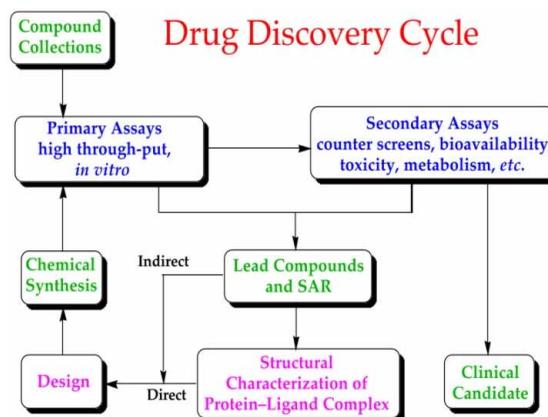**Figure 2. Drug Discovery And Development Activities**



Proteins are significant in any creature and are liable for basic usefulness in any being. Protein capacity is reliant on 3-Dimensional construction. By changing this construction, the usefulness of protein can be changed, and this is a significant component in drug disclosure. A large portion of the medications are intended to tie to the particular protein. So this is the significant element in drug conveyance to decide whether the protein can be tied to the medication or not. This is called the drug target cooperation forecast (Manne, R., et. al., 2021) (pl. Refer Figure 3).

(Source:-https://commons.wikimedia.org/wiki/File:Drug_discovery_cycle.svg)

## 1.4 Detection of Chemical Compound and Investigation in Pharmaceutical Labs

The "compound" which is set to turn into the medical molecule goes through security tests and a progression of trials to demonstrate that it is invested in the circulation system, appropriated to the legitimate site of activity in the body, used adequately and exhibits its non-harmfulness consequently,

**Figure 3. Drug Discovery Cycle**

can be viewed as protected and fruitful. When the compound is settled, the preclinical examination, for example in vitro investigations followed by the creature testing to check energy, poisonousness, and cancer-causing nature tests are performed. In the wake of finishing the preclinical assessments, the administrative specialist's award consent for the clinical preliminaries. The clinical preliminaries check regardless of whether the medication is working in the proposed component, its ideal portion, and timetable while the last two phases create genuinely significant information about adequacy, security and generally speaking the advantage risk relationship of the medication. In this stage, the likely connection of the medication with different not entirely settled and screened medication's drawn-out adequacy. After an effective finishing of the clinical preliminaries, the medications are sent off on the lookout for patients (Siddiqui, M., et. al., 2017) (pl. Refer Figure 4).

(Source:-https://www.sciencedirect.com/science/article/pii/S1878535213001056)

## 1.5 Usage of AI and ML in Drug Detection

The present technology is totally based on data science technique that played important part in reducing human work overload and improving the quality of life. The Author Team used integration of AI in ML to augment drug identification and development to make them accurate and efficient. (Rastogi, R., et. al. 2021).

AI and ML have a significant role to train the dataset on the basis of previous learning of data. The trained model are utilized in discovery, detection of drugs that filters out the target and studies

Figure 4. Summary of Phase-Wise Clinical Trial

| Phase of clinical trail | Number and type of subjects | Investigation |
|---|---|---|
| Phase 1 | 50–200 healthy subjects (usually) or patients who are not expected to benefit from the IMP | Is the IMP safe in humans? |
| | | What does the body do to the IMP? (pharmacokinetics) |
| | | What does the IMP do to the body? (pharmacodynamics) |
| | | Will the IMP work in patients? |
| Phase 2 | 100–400 patients with the target disease | Is the IMP* safe in patients? |
| | | Does the IMP seem to work in patients? |
| Phase 3 | 1000–5000 patients with the target disease | Is the IMP really safe in patients? |
| | | Does the IMP really work in patients |
| Phase 4 | Many thousands or millions of patients with the target disease | Just how safe is the new medicine? (pharmacovigilance) |
| | | How does the new medicine compare with similar medicines? |

its chemical structure and chemical combination on the basis of previous learning, and processes it through the ML and AI model that tests the drug with the target protein and evaluates the consequences in an approximately precise manner(Patel, V. et al., 2021) (pl. Refer Figure 5).

(Source:- https://www.sciencedirect.com/science/article/pii/S2667102621001066#bib0012)

## 1.6 Current Progress Through the advancement of Technology

The last ten years have been very progressive regarding technological advancements and new innovations and transformations. We can notice unprecedented evaluation in all the aspects of science and technology because of a solid outlay in research and development, education, entrepreneurship, and innovation all throughout the world. All the new technologies are making things easier and more systematic and even now we are presumably to see troublesome transformations considered as before. (Rastogi, R., et. al. 2021).

The view of the technological confluence is the framework of the various new products, procedures, operations, and services and softens the existence of old-fashioned technologies due to innovations, which is observed as the innovative demolition of today's era. The current modern era of the industrial revolution with rising appearance and qualified technologies and system arrangements, for example, we can say 5G, Artificial Intelligence, Machine Learning, VR (Virtual Reality) and Augmented Reality, IoT, BigData, cloud computing, and the blockchain and also cybersecurity is leading the way extreme positive results on upgrading the standard (quality) of life and experiences (Khan, M.K. et al., 2020) (pl. Refer Figure 6).

**(Source:-** https://www.go-rbcs.com/articles/the-increasing-pace-of-technology-advancement**)**

## 2. LITERATURE REVIEW

To make it easier for the reader to navigate this report, this chapter gives a brief overview of the content of the individual chapters. The literature review contains the aim and background for the research, describing a flowchart of the process, the conclusion summarizes the general findings of the literature review, and the conclusions of the three research papers are presented in a narrative form. This chapter concludes with some general perspectives on considerations and needs for further research based on this report.

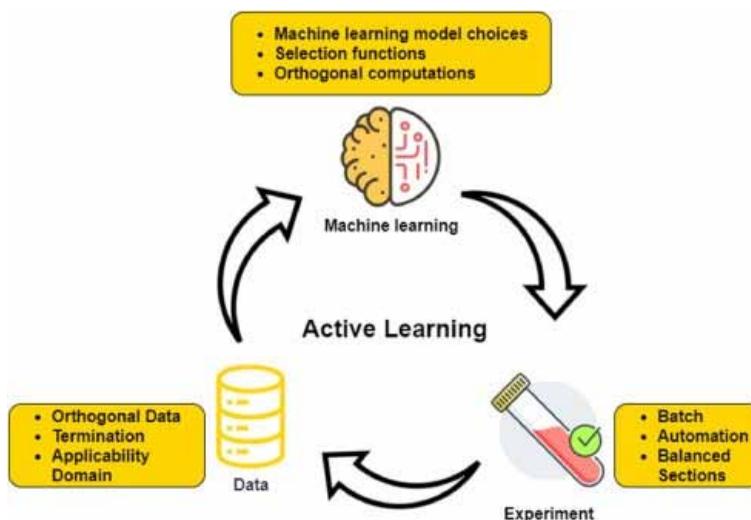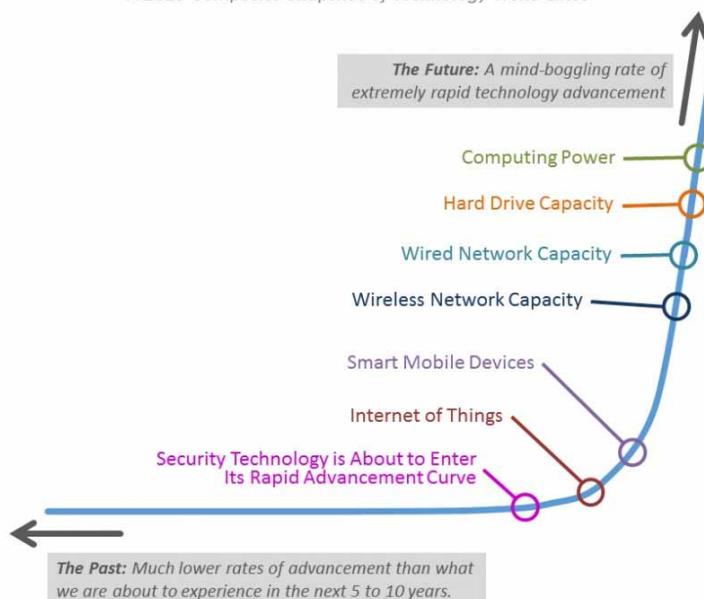**Figure 5. ML and AI in Drug Discovery**

**Figure 6. Security Technology is About to Enter a Period of Very Rapid Advancement**



Reboredo P. et al., (2021) and the team worked in drug detecting technique by applying different ML approaches and to find particular chemical component with varying chemical properties to cure diseases. The previous experimental data can be the basis for designing the model and feature extraction would help in developing new drug compound.

The machine learning models are trained on the basis of the structure of the molecules by descriptors(chemical characteristics of a molecule in numerical form,) which are able to capture the properties and characteristics. In the generating of the new drug may take up to 12 years, a fact that the cost exponentially grows up to billions euros in the launching of the drug into the market (Rastogi, R., et. al. 2021).

In the completion of task, Machine Learning and ML algorithm that are designed from raw and unprocessed data are used .Different machine learning model such as regression model, classification model, clustering model are applied for large dataset that covers the different variety of data and hence covers large aspects of pharmaceutical industry for prediction of new chemical bioactivities, and interaction between the target and diseases. The method like Naive Bayes and Support Vector Machine used to classify the drugs . The ML methodology should be reusable in all field where drug related task is performed. It should be the general model that can be used and applied in drug discovery. The Cheminformatics data repositories provides the data related with the drug, which can be tested and implemented on the model to get the result in precise manner. Different public Repositories like Drug Bank, PubChem, CheMBL are present (Reboredo, P.C., et al., 2021).

Park, J. et al., (2022) and the team elaborated on utilization of machine learning in detecting the bioactivity of compounds, Bioactive compounds are those chemical compounds that are used to cure different diseases. It is used as a curing agent after predicting using ML approaches. This research introduces how ML methodology are used in detecting and evaluating of novel bioactive

compounds. The identification of new bioactive compounds for a given target protein requires not only chemical information, but also detailed molecular information about the target protein, such as amino acid sequence, domain, and three-dimensional structural information. In designing the model, it can accurately classify the sample model into classes by applying regression to predict the value for a dataset containing labeled samples having many features could be categorize into three sub-datasets, called training, validation, and test datasets. The best models are selected through extensive training steps using training and validation datasets for which performance is evaluated using a test dataset.

There are some limitations even in applying ML-based approaches that are to be applied in NP-related research (Rastogi, R., et. al. 2021).

Experimentally verified bioactivity of natural products is insufficient to build an ML model. To process the lead optimization task more efficiently, many machine learning approaches have been studied recently: (i) atom modification reinforcement learning models that add or delete atoms or bonds, (ii) generative reinforcement learning which generates similar but modified structures, (iii) generative machine learning with controlled chemical properties that also generates similar modified structure with preserved predictive properties, and (iv) a 3D structure-based ligand design model that uses a 3D crystal structure of protein and ligand to generate novel molecules (Park, J., et al., 2022).

Gupta, R.et al.,(2021) and the team stated that AI and deep learning algorithms have been applied in an alternate piece of medication revelation processes like peptide, synthesis, structure-based virtual screening, ligand-based virtual screening, harmfulness forecast, drug checking, and quantitative design action relationship. Scientists had been challenges facing tough challenges in developing efficient and minimum risks drugs for two decades. The cost of discovering new drugs as therapeutic agents were exorbitant. The gene articulation method is broadly utilized in separating genes to comprehend illness mechanisms. Microarray and RNA-seq technologies creates a lot of gene articulation information for different problems. NCBI, TCGA, and Molecular Diversity are the huge storehouses that contain the component of gene articulation information, to figure out the objective qualities, liable for various disorders. In the location of large information, the enormous synthetic data set is expected from which finding ideal medications for a particular objective can be accomplished efficiently. Likewise, PubChem, ChemBL is uninhibitedly and an open available separately substance data set that contains information on different substance structure, including their organic, physical, compound, harmful properties, and a data set containing information on various bioactive compounds. The ChEMBL data set likewise contains data on retention, dispersion, digestion, and discharge (ADME), poisonousness properties of these mixtures, and, surprisingly, their target interactions. (Rastogi, R., et. al. 2021).

The course of drug screening incorporates the characterization and arranging of cells by picture examination through AI technology. ML models calculations are utilized in perceiving pictures with extraordinary exactness yet are not helpful while dissecting huge information. To order the objective cell, the ML model requires preparing so it can distinguish the phone and its elements, which is essentially finished by differentiating the picture of the designated cells, what isolates it from the foundation. Pictures with shifting finished highlights like wavelet-based surface elements and Tamura surface highlights are removed, and Molecular Diversity (2021) 25:1315-1360 1325 1 3 is additionally decreased in aspects through head part examination (PCA). A study proposes that most least-square SVM (LS-SVM) showed the most elevated grouping exactness of 95.34%. Regarding cell sorting, the machine should be quick to isolate out the designated cell type from the assigned example. Proof proposes that picture initiated cell arranging (IACS) is the most advance gadget that could measure the optical, electrical, and mechanical properties of the cell. (Gupta, R. et al., 2021).

Dara, S. et al.,(2021) and the team discussed Drug Detection using Machine Learning Tools and approaches particularly imposed in each phase of drug development to revive the exploration movement and close the chance of cost in remote (clinical) trials. There are many AI methods that work on the examination of medications information across the different order like Quantitative Structure-Activity Relationship (QSAR) investigation, de novo drug, hit revelations (a hit compound is

a particle that shows the ideal sort of movement in a screening measure) to recuperate exact outcome. (Rastogi, R., et. al. 2021).

The movement of protein structure is treated as a capability in drug plan. As it's undeniably true that there are numerous pollutants that have showed up in the human body as protein dysfunctions. There are Structural Drug Design outlines that are utilized to separate little particles in the protein targets. Protein organization in the 3D configuration required more cash and time for projecting the 3D construction and after this likewise it dealt with numerous issues, for example, making more exact again anticipating in the 3D design. Through profound learning and component extraction apparatuses, it is necessary to gauge the lower-level (here lower-level means secondary) structure and involve the protein's contacts. It acquires statistical data points on the connection between the construction and request from the element extractions.

The crucial steps of drug improvement in covering contain numerous biological sources for anticipating drug-protein interactivities. The hardships should be visible in the colossal forecasts which rely upon the various secret interactivities. Therefore, semi-administered preparing methods should be utilized to address these natural (unlabeled) and named date intricacies. For the most part, just marked information will make prevalent results. Also, the semi-managed innovation blends synthetic design, drug-protein intuitiveness network information, and hereditary material arranged information (Dara, S., et al., 2021).

Zhang, Y. et al.,(2021) and the team elaborated that Deep Learning notably accelerates the medication revelation process and adds to worldwide endeavors to sever the spread of infectious sicknesses. Aside from improving the effectiveness of screening of pharmacological mixtures against an expansive scope of irresistible infections, profound learning has likewise the possibility to dependably and productively perceive the medication competitors instead of Severe Acute Respiratory Syndrome Coronavirus 2 (SARS-CoV-2). For the recognizable proof of various expected drugs against SARS-CoV-2. This additionally incorporates Kaletra, Atazanavir, Remdesivir, Veneto Clax and some more. Preparing of Deep learning model to redress atoms dynamic contrary to anti-infection safe microscopic organisms lead to the disclosure of 2-amino-5-((5-nitro-2-yl)Thio)-1,3,4-thiadiazol (Helicine) with eight extra potential anti-infection agents from the ZINC data set. Strangely, these mixtures recognized by profound gaining are in a general sense divergent from the ordinary and standard antibiotics.

In sponger (parasite) research, deep learning has been placed in to foresee new mepacrine(antimalarial) drug applicants. Neves et al., (2020) in deep learning out how to acquire double, continuous Quantitative Structure-Activity Relationship (QSAR) models using datasets pulled from the ChEMBL data set. Various Antiviral mixtures are related to the assistance of Deep Learning. Two researchers Timmons and Hewage fostered a clever strategy called ENNAVIA, which connects profound learning and chemoinformatic, to perceive polypeptides with short, lethal, and magnificent organic action.(Rastogi, R., et. al. 2021).

The author also demonstrated the way that Deep Learning can bring down the time and cost of the medication discovery process, particularly in its high level stages. Subsequently profound learning-based approaches have been effectively used to perceive story (novel) antibacterial mixtures rather than a wide variety of infectious microorganisms, including organisms (bacteria), ameba, whip, and viruses. Distinguishing particles dynamic against anti-microbial safe microscopic organisms prompts the expectation of Helicin and an additional eight potential anti-microbials from the ZINC dataset. (Zhang,Y. et al.,2021).

Brogi, S. et al.,(2020) and the team stated that numerous philosophies are playing a steadily developing job in drug identification that is decrying in the conservative acknowledgment of promising medication candidates. These sorts of computational strategies are material in limiting the utilization of creature models in pharmaceutics research, helping the sensible plan of novel and safe medication candidates moving promoted medications, and supporting clinical scientific experts and pharmacologists all through the medication revelation course. As it is realized that there are numerous

pertinent quantities of papers that are centered around the different ligand-and construction put together methodologies or with respect to an association whereof to perceive promising particles for a given objective. Correspondingly Velázquez-Libera et al., expounds on a joined design and ligand-based way to deal with investigate the primary basics overseeing compassion of a progression of particles for the human Sigma1 receptor (S1R). This kind of receptor goes about as a significant medication focus for serving neuro-mental issues. (Rastogi, R., et. al. 2021).

The authors come across an effective S1R ally named RC-33 as a promising neuroprotective specialist. Likewise depicted a computational way for picking lead compounds from the colossal datasets of the relative multitude of synthetic outfits, got by HTS (high-throughput screening). Heavenly body Plots as broad strategies for seeing different and fitting sub-atomic portrayals to build the data contained in a noticeable portrayal and review of synthetic space are likewise presented by the creators. This kind of technique consolidates sub-structure-based portrayal and classification of the particles with a "customary '' coordinate-based portrayal of synthetic space. A commonplace result of the referenced procedures is that orchestrating the particles in a direct or resemble series prompts the development of classifications of the compound, which is otherwise called the "Constellations", in chemical space.

Team also developed and systematically approved an in silico method valuable for the hit to lead streamlining. In particular, from the micromolar HIV integrase requirements, the creator expounds computational usefulness in view of a silicon structure-based combinative library planning technique. The given philosophy is extremely valuable for joining the plan of the combinatorial library and side-chain holding with Quantum Polarized Ligand Docking (QLPD) and (MD) reproductions. The creator likewise expressed that in silico way to deal with investigate drug blend collaboration by using the tremendous accessible dataset enumerating synergism of anticancer medications (NCI-ALMANAC, with more than 290,000 synergy determinations) and also two AI systems, irregular backwoods and the outrageous angle supporting were used on the selected dataset (Brogi,S., et al.,2020).

Siddiqui, M.R. (2017) et. al. and the team stated their work in A comparison of machine learning techniques for detection of drug target articles. Author and team analyzed that Lately, significant advancement in treating infections, for example, malignant growth, AIDS, or Parkinson's illness, among numerous others, has been conceivable on account of the ID of medication targets connected to these illnesses. The momentum drug revelation process is essentially centered around the pursuit and approval of medication applicants that follow up on a specific helpful objective. Initially, the course of a specific infection is contemplated and its physiologic not entirely set in stone to recognize the medication targets connected with this illness. Then, at that point, new medications are intended to follow up on these objectives. Because of the significant expense furthermore, the significant time-frame expected by the medication advancement process, drug industry requirements to work on the systems for focusing on targets and medication applicants in the medication disclosure process. (Rastogi, R., et. al. 2021).

Here the author and his team works on classification machine learning based models detecting the elements that have drug target compounds and particles. The task they are facing is biomedical knowledge and resources. Author used the PubMed service to access the dataset which is article based.They created the corpus of negative and positive drug targets from Drug Bank and PubMed. They work with the NLP-like algorithm; they extract the feature in the given articles and find the drug target through this. Here the main task is feature extraction which is based on NLP tokenization and regular expression based and also works with UMLS MMTx. MMTx investigates the message grammatically to part it into parts of various syntactic levels: sentences, phrases, lexical components and tokens. DrugNer broadens the data given by MMTx, by the utilization of the classification rules suggested by the WHO International Nonproprietary Names (INNs) Program10 to distinguish and group drug substances.

A few investigations were completed to approve the proposed classifier for drug-target articles. Since the noticed proportion between the quantity of positive and negative models is profoundly uneven, we have concentrated on the impact of utilizing unique extents in sure and negative models in the preparation set. Consequently, we have considered 4 preparation datasets containing 5% (genuine

appropriation), 10%, 20% and half of positive models separately, in which the different preparation sets share as numerous models as could be expected. This arrangement decreases the chance of trivial outcomes because of contrasts in preparing information (Siddiqui, M. R., et al., 2017).

Manne, R., (2021) elaborated in his work in Machine Learning Techniques in Drug discovery and development. The author represented a model that every real-world problem in this world is solved using technology and machine learning techniques so they work in machine learning techniques and try to solve problems in pharmaceutical companies(all stages of drug discovery).

The drugs which were endorsed in the year 2005 to 2006 took a normal clinical advancement season of six and half years, also, from 2008 to 2012 took a normal season of 9.1 years. In the later phase of clinical testing, the pace of drug failure has increased. The failure of the drug and the time-consuming cycle, as it takes truly significant stretches alongside enormous costs can be baffling, particularly when trials were insufficient fruitful. The author describes all the machine learning algorithms but they focus on the Naive Bayes algorithm. According to him - Naive Bayesian classifiers are utilized in cheminformatics to foresee organic properties instead of physicochemical properties. This is applied in foreseeing the harmfulness of the compound, protein target, and bioactivity grouping for drug-like particles.

Naive B For the informational collections that are recovered Credulous base work on the precision. In the order apparatuses for biomedical information, NB calculations have shown extraordinary guarantee despite the fact that information is filled with undesirable information called clamor. NB method additionally demonstrated to play a significant part in ligand-target collaboration expectations, which is additionally an extraordinary venture towards lead revelation.

The nearby design of the information K-NN calculation is touchy, thus it is ideal to compute properties with solid regions, as is the situation with protein work expectation. The K-NN approach doesn't have limits. Random Forest has properties that work on the forecast of QSAR information. These properties are underlying descriptor determination, high exactness of expectation. Creator has distributed a technique which was applied to mining estrogen receptors from a dataset of 57000 particles and this strategy utilizes an alternate arrangement of descriptors to fabricate choice tree model which is precise.

Author has proposed a system which depends on profound learning for drug target corporation forecasts. These profound learning systems which are utilized for DTI forecasts take the two mixtures and protein data as info (Manne, R., 2021).

Rollinger, J. M. et. al., (2008) and the team demonstrated their work in Virtual screening for the discovery of bioactive natural products. Main concept of the team is to find the chemical (drug) component in the natural products. They give some basic ideas, limitations and requirements for machine learning strategies and support their thoughts in Natural Product research with already performed studies.

The normal thought of all computational methodologies inside the early drug discovery process is to mine pretty much enormous compound data sets in silico and to choose a predetermined number of competitors proposed to have the ideal natural movement. For this interaction the term 'data mining' was instituted in 1996, which was compactly characterized by Gasteiger and co-creators: 'to extract information from a large arrangement of data(training) to make forecasts of new events. Inside the lead revelation process, virtual screening advances have to a great extent improved the effect of computational science and these days chemo informatics assumes a dominating part in beginning stage drug research. The vital objective of the utilization of such strategies is to diminish the general cost related to the revelation and advancement of another medication, by recognizing the most encouraging contender to concentrate the trial endeavors on. As of late, distributed books and surveys on the effect of computational science for lead structure assurance feature these works.

Assuming the 3D design of the organic objective is known, high throughput docking ended up being a significant construction based virtual screening strategy to be utilized. Inside this unique situation, the scoring of hits recovered still remains an inquiry that is frequently talked about. At present, truth be told, the significant shortcoming of docking programs lies not in the docking

algorithms themselves but rather still in the mistake of the capacities that are utilized to assess the liking among ligands and focus on the alleged scoring capacities.

Relating to the medication revelation from nature we are confronting two realities:

(i)measurements show that the bunch of primarily assorted regular mixtures are the most preferred wellspring of new medications for clinical use;

    (ii)  the medication revelation process has moved towards more level headed ideas in view of the expanding comprehension of the atomic standards of protein-ligand collaborations. Restricted endeavors applying creative techniques in silico apparatuses in NP research are sought after up to this point, on the grounds that the quest for bioactive mixtures is a complex and multidisciplinary challenge (Rollinger, J. M., et al., 2008). Pl. refer Table 1 for literature review summary.

## 3. METHODOLOGY AND SETUP OF EXPERIMENT

Author and his team used the ChEMBL database which contains multiple row and column values (the actual shape of the database is 4695 rows, 882 columns) but it has taken the relevant attributes values like pointing out ligands and protein. Accordingly, the authors first cleaned the data, found the missing values and maintained these missing values and performed some EDA (Exploratory Data Analysis) on it, analyzed the dataset, understood the dataset, found some correlation between the data attributes and then created the model. In model creation, regression algorithms were used like random forest, decision tree, QSAR, Y-Square and many more machine learning Algorithms. While comparing these algorithms and checking the accuracy of the model, the QSAR model provided a good accuracy which is good for the proposed prediction model (as per Figure 7).

## 4. SETUP

### 4.1 Name of Algorithms Used

Authors used Random Forest, Decision Tree, QSAR, Y-Square Regression Machine Learning Model and Checked the Model accuracy.

### 4.2 Types of Databases

The one chemical database that contains the information of chemical compounds used in Drug Detection is ChEMBL. ChEMBL is a dataset of bioactive drug-like little particles, that comprises 2-D designs, determined properties (for example logP, Molecular Weight, Lipinski Parameters, and so on) and preoccupied bioactivities (for example restricting constants, pharmacology and ADMET data).It is a physically organized data set of bioactive particles with drug-like properties. It unites substance, bioactivity and genomic information to help the interpretation of genomic data into viable new drugs. The ChEMBL Database is a dataset that contains organized bioactivity information of 2 million mixtures. It is gathered from in excess of 76,000 archives, 1.2 million measures and the information spans 13,000 targets and 1,800 cells and 33,000 signs (Data as of March 25, 2020; ChEMBL variant 26).
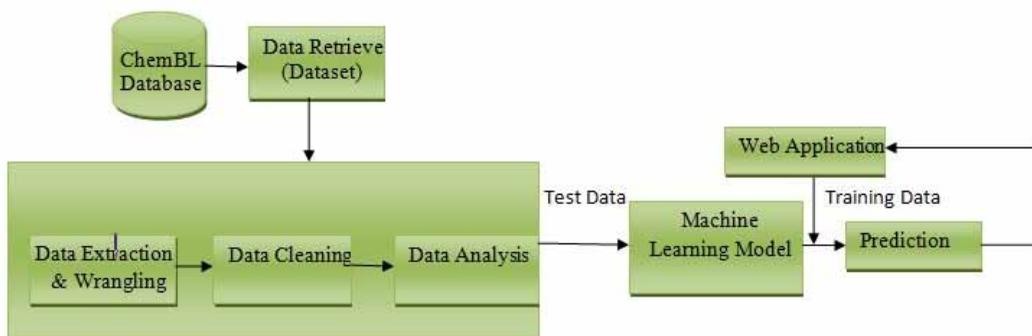
### 4.3 SCOPE AND ACCESS:

CheMBL database stores the different bioactivity related data of different chemical compounds for drug object. Data can be fetched from here and apply machine learning model, the data can be trained according to ML model training and data are analyzed in order to develop compounds for screening and identification during drug discovery. Clients can look for intensifies utilizing a keyword search

**Table 1. Tabular Summary for Literature Review Based Papers**

| S.No. | Paper Name | Summary | Methodology, dataset, Algo | Concluding Remarks |
|---|---|---|---|---|
| 1 | A review on machine learning approaches and trends in drug discovery (Reboredo, P. C. et al., 2021) | This review will focus in chiefly on the techniques used to demonstrate the atomic information, as well as the organic issues tended to and the Machine Learning calculations utilized for drug revelation as of late | CADD,chemspider, PubChem ChemBL Binding DB ZINC(Chemical Compound Database).Compound Selection, Experimental Validation. | 1. Binding Energy Calculation 2.Active sites Identification 3.Identify Disease-Drug Interaction. |
| 2 | A Brief Review of Machine Learning-Based Bioactive Compound Research(Park, J., et al., 2022) | This review presents how AI approaches can be utilized for the identification and assessment of bioactive mixtures. The identification of new bioactive compounds for a given target protein requires not only chemical information, but also detailed molecular information about the target protein, such as amino acid sequence, domain, and three-dimensional structural information. | Machine Learns from Data and Creates a Model for a Task Using Machine Learning Algorithm.To build a model that can accurately classify samples into classes (classification) or predict values for samples (regression), a dataset, which contains a large number of correctly labeled-samples with many features, should be divided into three sub-datasets, called training, validation, and test datasets. | Machine Learning Application of NP(Natural Product) or NP-Like Chemical Compounds Discovery for Cardiovascular and Metabolic Diseases. |
| 3 | Artificial intelligence to deep learning: machine intelligence approach for drug discovery(Gupta, R., et al., 2021) | Machine learning and deep learning algorithms have been applied in different parts of drug discovery processes like peptide, synthesis, structure-based virtual screening, ligand-based virtual screening, toxicity prediction, drug monitoring and quantitative structure -activity relationship. | Machine alongside algorithms like Naïve Bayes, decision tree (DT), hidden Markov models (HMM) .Some Machine learning models like Support Vector Machine, Neural Network, Clustering. Data set location from chemspider, PubChem ChemBL Binding DB ZINC(Chemical Compound Database). | Artificial intelligence consciousness in essential and auxiliary medication screening. Peptide blend and little molecule design. Recognizing Disease tweaking Target Protein(Using Bioinformatics Analysis). |
| 4 | A review on Machine Learning in Drug Discovery.,Artificial Intelligence Review(Dara,S., et al., 2021) | Using AI techniques authors want to solve the problem of drug development using QSAR analysis and the protein composition in 3D Format | Quantitative Structure-Activity Relationship (QSAR), Protein composition in 3D format, semi-supervised training, the dataset from ChEMBL. | Machine learning techniques are best to solve all the problems in the discovery of drugs.It can better resolve all the molecular activities by forecasting the data from unsupervised learning. |
| 5 | Deep Learning-Driven Drug Discovery, Tackling Severe Acute Respiratory Syndrome Coronavirus 2 (Zhang,Y. et al.,2021) | Giving some basic ideas of the DL used in drug discovery. Deep learning is very accurate to tackle all the severe acute respiratory syndrome COVID-19. Authors solved the hidden problem with the help of DL. | DL along with algorithms like CNN, and Naïve Bayes. Utilized Deep learning out how to acquire double, continuous Quantitative Structure-Activity Relationship (QSAR) models using datasets pulled from the ChEMBL information base. | Deep learning can lower the time and cost of the drug detection process, especially in its advanced phases. Hence deep learning-based approaches have been successfully utilized to recognize story (novel) antibacterial compounds opposed to a wide variation of contagious microorganisms, including microbes (bacteria), ameba, flagellate, and viruses |
| 6 | Editorial:In *silico* Methods for Drug Design and Discovery(Brogi,S., et al.,2020) | Performed an experiment proving the economical recognition of promising drug applicants. These types of computational methods are applicable in restricting the use of animal models in pharmaceutics research, helping the logical design of novel and safe drug applicants shifting marketed drugs, and supporting medical chemists and pharmacologists throughout the drug discovery course. | CYP-mediated reaction rules and the site of metabolism (SoM) prediction. Two machine learning (ML) procedures, Random Forest (RF), and Extreme Gradient Boosting (XGBoost) were employed on the selected dataset. dataset reporting synergism of anticancer drugs (NCI-ALMANAC, with over 290,000 synergy determinations) | 1. reproving in the economical recognition of promising drug applicants. 2. restricting the use of animal models in pharmaceutics research, for helping the logical design of novel and safe drug applicants shifting marketed drugs, supporting medical chemists and pharmacologists throughout the drug discovery course. |
| 7 | A comparison of machine learning techniques for detection of drug target articles(Siddiqui, M. R., et al., 2017) | Performed an extensive experimental analysis using a combination of techniques for feature selection and the most important machine learning algorithms for text classification. | MMTx, UMLs, Drugner, Classification(SVM, C4.5, Bayesian Statistics) Machine learning algorithms, MedLine, PubMed, DrugBank(Drug Target Article) datasets. | A few investigations were completed to approve the proposed classifier for drug-targeted articles.Using NLP and MMTx to get the accuracy. |
| 8 | Machine Learning Techniques in Drug discovery and development (Manne, R., 2021) | Using Machine Learning and Deep learning, authors try to solve pharma based problems. | Chemical Compound related dataset from Kaggle and UCI ML. K-NN, Naive Bayes, Random Forest. | Author has proposed a system which depends on profound learning for drug target corporation forecasts. |
| 9 | Virtual screening for the discovery of bioactive natural products(Rollinger, J. M., et al., 2008). | They give some basic ideas, limitations, and requirements of machine learning strategies and support their thoughts in Natural Product research with already performed studies. | 3d database of natural compounds, Silico, Vitro, LC-MS, LC-NMR, GC-MS, etc. | Only restricted endeavors applying creative in silico apparatuses in NP research are sought after up to this point, on the grounds that the quest for bioactive mixtures is a complex and multidisciplinary challenge. |

**Figure 7. Drug Detection Process and Model**



with names/equivalent words or ChEMBL identifiers. However, a more compelling procedure will frequently be to look through by compound design.

### 4.3.1 Accessing Resources

In addition to the dataset, the ChEMBL bunch has created devices and assets for information mining. These incorporate Kinase SARfari, a coordinated chemogenomic workbench focused on kinases. The framework integrates and connects succession, construction, mixtures and screening information.

## 4.4 Dataset

The researcher team searched in the web and found different sets of datasets some of them are mentioned below:

1. ChEMBL Database
2. UCI ML Drug Review Dataset
3. QM9 drug discovery dataset (QM9 provides quantum chemical properties).

But the research team used the ChEMBL Database of Drug Compound information as it is free to download. Author filtered the database and fetched their relevant information through this. So they created a total of 2818 rows and 4 columns of data Set, which the author called BioActivity Dataset.

### 4.4.1 BioActivity DataSet Attributes

molecule_chembl_id:- For every molecular compound a unique ChEMBL id is used.

canonical_smiles:- this attribute denotes the complex and critical structure of compound into a linear text format.

bioactivity_class:- This attributes inherited from the standard value of compound here mainly 3 values are stored if the standard value of compound was less than 1000 author denoted it as active state, if standard value of compound was greater than 10000 author denoted it as inactive state and if standard value between range of 1000 to 10000 then author denoted it as intermediate state.

Standard_value:- Bioactive peptides (BP) are regular substances formed by amino acids joined by covalent bonds known as amide or peptide bonds.

### 4.4.2 ChEMBL Source and URL

! pip install chembl_webresource_client

Using Above line of code, the Author installed the ChEML web service package.

from chembl_webresource_client.new_client import new_client

Importing the above library, the Author used the ChEMBL database and filtered the dataset according to their model.

## 4.5 METADATA

Protein's (standard_value):The Information of the target Protein is necessary as it hits the molecules,and alternates the chemical structure, so in order to discover a drug target protein identification is necessary. Target expectation with AI calculations can assist with speeding up this inquiry, restricting the quantity of required tests. In any case, Drug-Target Interactions data sets utilized for preparing present a high measurable predisposition, introduced to countless bogus up-sides, hence expanding the time and cost of exploratory approval promotions.

QSAR Dataset: Quantitative design action relationship (QSAR) is a computational or numerical displaying technique to uncover connections between natural exercises and the underlying properties of synthetic mixtures.

The dataset in the wake of applying numerical activity and Computational demonstrating gives a QSAR dataset (Pl. refer Figure 8, Figure 9, Figure 10).

## 4.6 ChEMBL Database and Data Set One Sample Image

**Figure 8. ChEMBL Database Main Class Data**

| | cross_references | organism | pref_name | score | species_group_flag | target_chembl_id |
|---|---|---|---|---|---|---|
| 0 | [{'xref_id': 'P11511', 'xref_name': None, 'xre... | Homo sapiens | Cytochrome P450 19A1 | 20.0 | False | CHEMBL1978 |
| 1 | [{'xref_id': 'P22443', 'xref_name': None, 'xre... | Rattus norvegicus | Cytochrome P450 19A1 | 20.0 | False | CHEMBL3859 |

**Figure 9. ChEMBL Database**

| | B | C | D | E | F | G | H | I | J | K | L | M | N | O | P | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | 82585 [] | | CHEMBL6667! | Inhibition of Cy | B | | | BAO_000C | BAO_000C | single prot | CC12CCC(O)CC1=CCC1C2CCC2( | CHEMBL1: | J. Med. Ch |
| | 94540 [] | | CHEMBL6667! | Inhibition of Cy | B | | | BAO_000C | BAO_000C | single prot | C[C@]12CC[C@H]3[C@@H](CC | CHEMBL1: | J. Med. Ch |
| | 112960 [] | | CHEMBL6617C | In vitro inhibiti | B | | | BAO_000C | BAO_000C | single prot | CCn1c(C(c2ccc(F)cc2)n2ccnc2)c | CHEMBL1: | Bioorg. Me |
| | 116766 [] | | CHEMBL6617C | In vitro inhibiti | B | | | BAO_000C | BAO_000C | single prot | CCn1cc(C(c2ccc(F)cc2)n2ccnc2) | CHEMBL1: | Bioorg. Me |
| | 118017 [] | | CHEMBL6617C | In vitro inhibiti | B | | | BAO_000C | BAO_000C | single prot | Clc1ccccc1Cn1cc(Cn2ccnc2)c2c | CHEMBL1: | Bioorg. Me |
| | 118020 [] | | CHEMBL6617C | In vitro inhibiti | B | | | BAO_000C | BAO_000C | single prot | Cc1ccc(S(=O)(=O)n2cc(C(c3cccc | CHEMBL1: | Bioorg. Me |
| | 126976 [] | | CHEMBL6617C | In vitro inhibiti | B | | | BAO_000C | BAO_000C | single prot | CCn1ccc2cc(C(c3ccc(F)cc3)n3cc | CHEMBL1: | Bioorg. Me |
| | 135646 [] | | CHEMBL6617C | In vitro inhibiti | B | | | BAO_000C | BAO_000C | single prot | Cn1cc(C(c2ccc(F)cc2)n2ccnc2)c | CHEMBL1: | Bioorg. Me |
| | 138013 [] | | CHEMBL6617C | In vitro inhibiti | B | | | BAO_000C | BAO_000C | single prot | CCn1cc(C(c2ccc(F)cc2)n2ccnc2) | CHEMBL1: | Bioorg. Me |
| | 138016 [] | | CHEMBL6617C | In vitro inhibiti | B | | | BAO_000C | BAO_000C | single prot | CCn1ccc2cc(C(c3ccccc3)n3ccnc | CHEMBL1: | Bioorg. Me |
| | 139306 [] | | CHEMBL6617C | In vitro inhibiti | B | | | BAO_000C | BAO_000C | single prot | N#Cc1ccc(Cn2cc(Cn3ccnc3)c3cc | CHEMBL1: | Bioorg. Me |
| | 184046 [] | | CHEMBL6667! | Inhibitory activ | B | | | BAO_000C | BAO_000C | assay form | CCCCCCN1C(=O)CCC(CC)(c2ccn | CHEMBL1: | J. Med. Ch |
| | 184487 [] | | CHEMBL6667! | In vitro inhibito | B | | | BAO_000C | BAO_000C | assay form | c1ccc2c(c1)CCC1(c3cc[nH]n3)( | CHEMBL1: | J. Med. Ch |
| | 185411 [] | | CHEMBL6667! | Inhibitory activ | B | | | BAO_000C | BAO_000C | assay form | CCCCCCCC1(c2ccncc2)CCC(=O)l | CHEMBL1: | J. Med. Ch |
| | 185984 [] | | CHEMBL6616! | In vitro inhibito | B | | | BAO_000C | BAO_000C | assay form | O=C1/C(=c Outside ty Values for | CHEMBL1: | J. Med. Ch |
| | 185988 [] | | CHEMBL6616! | In vitro inhibito | B | | | BAO_000C | BAO_000C | assay form | O=C1/C(=c Outside ty Values for | CHEMBL1: | J. Med. Ch |
| | 186032 [] | | CHEMBL6667! | Inhibition of ar | B | | | BAO_000C | BAO_000C | microsom | C[C@]12CC[C@H]3[C@@H](CC | CHEMBL1: | J. Med. Ch |
| | 186034 [] | | CHEMBL6667! | Inhibition of ar | B | | | BAO_000C | BAO_000C | microsom | C[C@]12CC[C@H]3[C@@H](C[ | CHEMBL1: | J. Med. Ch |
| | 187134 [] | | CHEMBL6616! | In vitro inhibito | B | | | BAO_000C | BAO_000C | assay form | O=C1/C(=C\c2c[nH]cn2)CCc2cc | CHEMBL1: | J. Med. Ch |
| | 188443 [] | | CHEMBL6616! | In vitro inhibito | B | | | BAO_000C | BAO_000C | assay form | O=C1/C(=C/c2ccnc2)CCc2ccccc | CHEMBL1: | J. Med. Ch |
| | 188448 [] | | CHEMBL6616! | In vitro inhibito | B | | | BAO_000C | BAO_000C | assay form | O=C1/C(=C/c2ccnc2)CCc2ccccc | CHEMBL1: | J. Med. Ch |
| | 188467 [] | | CHEMBL6667! | Inhibition of ar | B | | | BAO_000C | BAO_000C | microsom | C[C@]12CC[C@H]3[C@@H](CC | CHEMBL1: | J. Med. Ch |
| | 189643 [] | | CHEMBL6616! | In vitro inhibito | B | | | BAO_000C | BAO_000C | assay form | c1ccc2c(c Outside ty Values for | CHEMBL1: | J. Med. Ch |

**Figure 10. BioActivity Dataset Fetched from ChEMBL Database**

| A | B | C |
|---|---|---|
| molecule_chembl_id | canonical_smiles | standard_value |
| CHEMBL341591 | CC12CCC(O)CC1=CCC1C2CCC2(C)C(CC3 | 7100 |
| CHEMBL2111947 | C[C@]12CC[C@H]3[C@@H](CC=C4C[C | 50000 |
| CHEMBL431859 | CCn1c(C(c2ccc(F)cc2)n2ccnc2)c(C)c2cc( | 238 |
| CHEMBL113637 | CCn1cc(C(c2ccc(F)cc2)n2ccnc2)c2ccccc2 | 57 |
| CHEMBL112021 | Clc1ccccc1Cn1cc(Cn2ccnc2)c2ccccc21 | 54 |
| CHEMBL324070 | Cc1ccc(S(=O)(=O)n2cc(C(c3ccccc3)n3cc | 5400 |
| CHEMBL41761 | CCn1ccc2cc(C(c3ccc(F)cc3)n3ccnc3)ccc2 | 41 |
| CHEMBL111868 | Cn1cc(C(c2ccc(F)cc2)n2ccnc2)c2cc(Br)c | 78.5 |
| CHEMBL111888 | CCn1cc(C(c2ccc(F)cc2)n2ccnc2)c2cc(Br) | 51.8 |
| CHEMBL112074 | CCn1ccc2cc(C(c3ccccc3)n3ccnc3)ccc21 | 205 |
| CHEMBL324326 | N#Cc1ccc(Cn2cc(Cn3ccnc3)c3ccccc32)cc | 50 |
| CHEMBL37321 | CCCCCCN1C(=O)CCC(CC)(c2ccncc2)C1= | 6600 |
| CHEMBL353068 | c1ccc2c(c1)CCC1C(c3cc[nH]n3)C21 | 51000 |
| CHEMBL41066 | CCCCCCCC1(c2ccncc2)CCC(=O)NC1=O | 3200 |
| CHEMBL166709 | O=C1/C(=C/c2cccnn2)CCc2ccccc21 | 250000 |
| CHEMBL424556 | O=C1/C(=C/c2ccnnc2)CCc2ccccc21 | 103000 |
| CHEMBL1630273 | C[C@]12CC[C@H]3[C@@H](CC=C4[C@ | 6800 |
| CHEMBL1630261 | C[C@]12CC[C@H]3[C@@H](C[C@@H] | 50 |
| CHEMBL169251 | O=C1/C(=C\c2c[nH]cn2)CCc2ccccc21 | 170 |
| CHEMBL168636 | O=C1/C(=C/c2cccnc2)CCc2ccccc21 | 9200 |
| CHEMBL90585 | O=C1/C(=C/c2ccncc2)CCc2ccccc21 | 4600 |
| CHEMBL1629805 | C[C@]12CC[C@H]3[C@@H](CCC4=CCC | 49 |
| CHEMBL433728 | c1ccc2c(c1)CCC1C(c3ccsc3)C21 | 250000 |
| CHEMBL38877 | CCCN1C(=O)CCC(CC)(c2ccncc2)C1=O | 31000 |

## 4.7 SIZE (STORAGE SPACE)

Here, the author and his team worked with the ChEMBL Database author and had 43 columns in the database they Select and recover bioactivity information for the Inhibition of Cytochrome P450. Here, we will retrieve just bioactivity information for Covid 3C-like proteinase (CHEMBL3927) that are reported as IC50 values in nanomolar units. The dataset contains a total of 2818 rows and 4 cols, which is a total of 2MB.

## 4.8 HARDWARE REQUIREMENTS

The specifications of the Hardware which are required for the smooth working of this Bio-Activity Detection system are:
RAM (Memory): 2 GB and above.
RAM Frequency: 800MHz in older systems and, DDR2 modules up to 5200MHz in DDR5. Current generation DDR4 modules usually need 3200MHz.
Storage Capacity:

a) For Application: 100 MB and above.
b) For Database: 800 MB and above.

Processor:

a)  Type: Dual Core and above versions of the Intel Core processors.
b)  Speed: 1.7 GHz Upto Max Turbo Frequency at 4.1 GHz.

## 4.9 SOFTWARE REQUIREMENTS

Along with the hardware requirements there is some software that is responsible to take over the hardware used. Here in this Drug detection system, the software required is listed below:

Front-end: Flask - (a python framework for designing UI for the whole system)

Back-end: Python 3.0 and all above versions, ChEMBL Database.

Libraries Used (additional): NumPy, pandas, matplotlib, matplotlib inline, sklearn, pillow, pickle.

### 4.9.1 OS REQUIREMENTS

Windows 10: for the backend work windows 10 is required. The versions are lower than this such that Windows 7, 8, or 11 which is the latest version can also be used.

Windows is used to provide the space or background to run all the software of Python required for the Bio-Activity Detection System and also helps in downloading the datasets.

Here, other operating systems such as Linux, and macOS can also be used.

### 4.9.2 DATABASE REQUIREMENTS

Database is required so that we can easily store the information regarding bioactivity of molecular compounds at a particular instant of time and can access it later without any problem. The research team uses SMILES, this SMILES stands for Simplified Molecular Input Line Entry System using SQL which is a programming language that develops databases.

### 4.9.3 STORAGE REQUIREMENTS

System required a total of approx. 10-50 MB storage to store the dataset and the author require a minimum 4 GB RAM system to build their model.

## 4.10 FRONT-END

The Researcher team is looking to design an attractive and user-friendly application for this Bio-Activity Drug Detection Model. To fulfill this criterion, the research team is using Flask (a python framework) for designing an application.

Flask framework is used for designing and developing all web applications by using Python as the backend, which is fully implemented on Werkzeug (Debugger) and Jinja2 (a fast, expressive, extensible templating engine).

There are advantages of using Flask as a Front-end, as Flask provides a built-in development server with a fast debugger. Flask is lightweight in comparison to other frameworks such as Django, and Kivy. Here in Flask, secure cookies are also supported along with the template using Jinja2. The main advantage is that support for the unit testing is built-in.

## 4.11 BACK-END

For Back-end, Python 3.0 and above all versions are used. Python is a broadly useful interpreted, object-situated, intelligent, and significant level programming language. As it is interpreted as a language, it tends to be very much utilized in AI.

Along with the interpreted it has additional many libraries such as NumPy which is used to work with arrays and multidimensional arrays, pandas for the use in data science or data analysis and AI tasks, Matplot - it is cross-platform used for data visualization and graphical plotting, matplot inline is used to set backend of matplot to the 'inline backend', sklearn provides the selection of the systematic tools for the ML and statistical modeling including regressions, classifications, etc .. Pillow is an imaging library in python used for opening, manipulating and also saving different image file formats.

## 4.12 STEPS OF EXECUTION

1.First Author installed the ChEMBL web service library to load the database in the local machine.

2.  Filter the database on the basis of targeted_chembl_id and selected_type.
3.  Save the dataset.
4.  Clean the dataset and Analyze the dataset.
5.  Feature Selection Machine building.
6.  Applying Machine learning Algorithms.
7.  Check the accuracy of the model.

## 4.13 DIAGRAMS

### 4.13.1 Flow Chart

It depicts the flowchart for BioActivity drug detection of compounds. In the above figure the author described each and every way they use the database and fetched the relevant dataset through this. The author described every step of their project mechanism in the above flow chart(As per Figure 11).

### 4.13.2 ER Diagram

It depicts how the model had been trained using a dataset and applying ML algorithm, that contains how using ChembL Database the data is abstracted and bioactivity is measured by preparing the dataset how it affects the target protein and on the basis of QSAR model, the bioactivity of the drug is detected(As Per Figure 12)

Figure 11. Flowchart for BioActivity Detection of Drug Compound

**Figure 12. Entity Relationship DataModel of BioActivity Detection of Drug Compounds**



### 4.13.3 Block Diagram

The diagram shows the pipeline of the Bio-Activity Detection of Drug compounds from the starting phase to the end phase. First of all, all the target structures are collected to determine the output of the product, and then according to the character of that dataset and features, small packets are to be bound towards the development phase. A study of the QSAR leads to optimization of that output and then the best drug is selected(As per Figure 13)

### 4.13.4 Object Oriented Class Diagram

It shows the method of objects by which the sentiments are integrated.It describes the whole process from scratch to end(As per Figure 14).

**Figure 13. Block Diagram OF Bio-Activity Detection of Drug Compounds**

Figure 14. A Sentiment Analysis Method of Objects by Integrating Sentiments from Tweets



## 4.13.5 Data Flow Diagram (level 0, level 1, level2)

It depicted working flow and steps of bioactivity detection of drugs in the Model, that initially required the input dataset abstracted from the CHembL Database and trained the dataset by applying the Machine learning model on the dataset, a model is prepared which detects bioactivity of drugs(As per Figure 15).

It shows data flow diagram of the prediction of Bioactivity of Drug level. This shows the deeper process of the system, in which the first step is a compound selection that is relevant to the dataset and model the data through a regression model and feature selection model that can efficiently detect and predict the bioactivity(As per Figure 16).

Figure 15. Data Flow Diagram(Level 0)

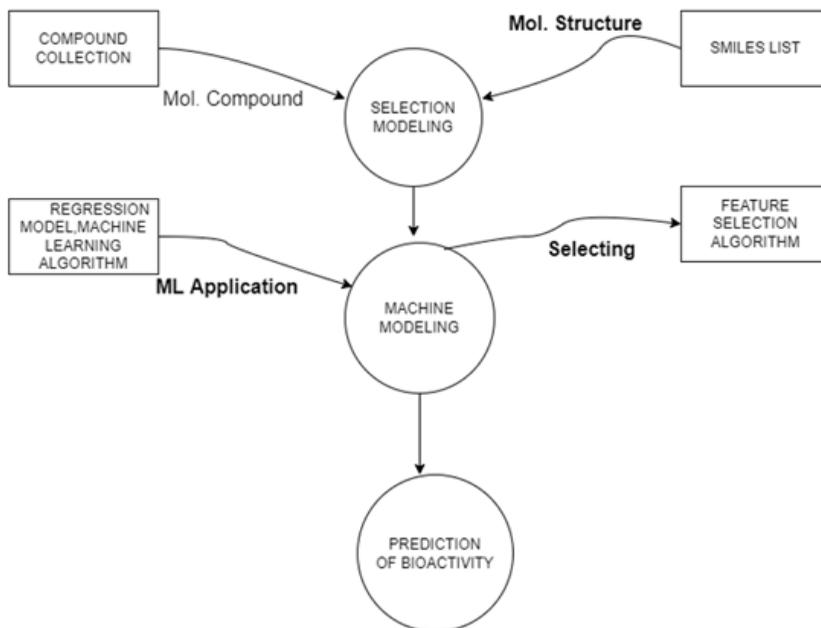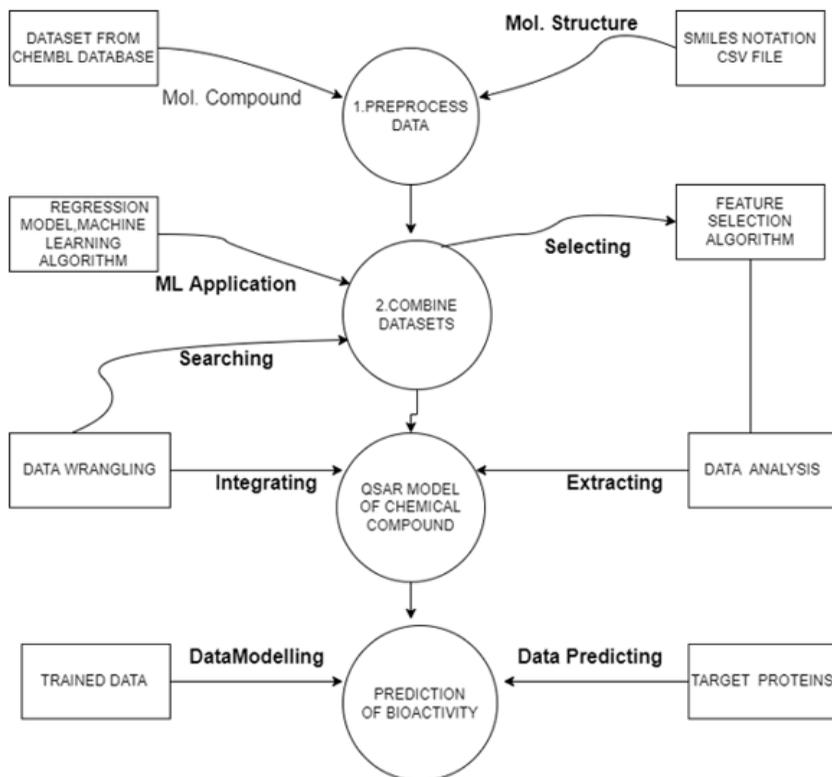**Figure 16. Data Flow Diagram(Level 1)**



**Figure 17. Data Flow Diagram(Level 2)**

This above figure shows the data flow diagram of Drug Detection level 2. This shows the deeper process of the system. In this figure, the textual content which is given by the user is analyzed in processes. First, this data is preprocessed, then processed data is further performed in the feature selection process. After that, the classifier classifies the Drug Data by applying Feature Selection and returns the result (as per Figure 17).

### 4.13.6 Use case diagram

The above figure represents the use case of the Drug Detection System . Author had plotted all the uses of the system according to the user or we can say actors . Actors here are actually the user who is responsible for his respective tasks throughout the process . In this system there are 5 actors who are responsible for the whole process . Each one of the actors has their own task to perform (as per Figure 18).
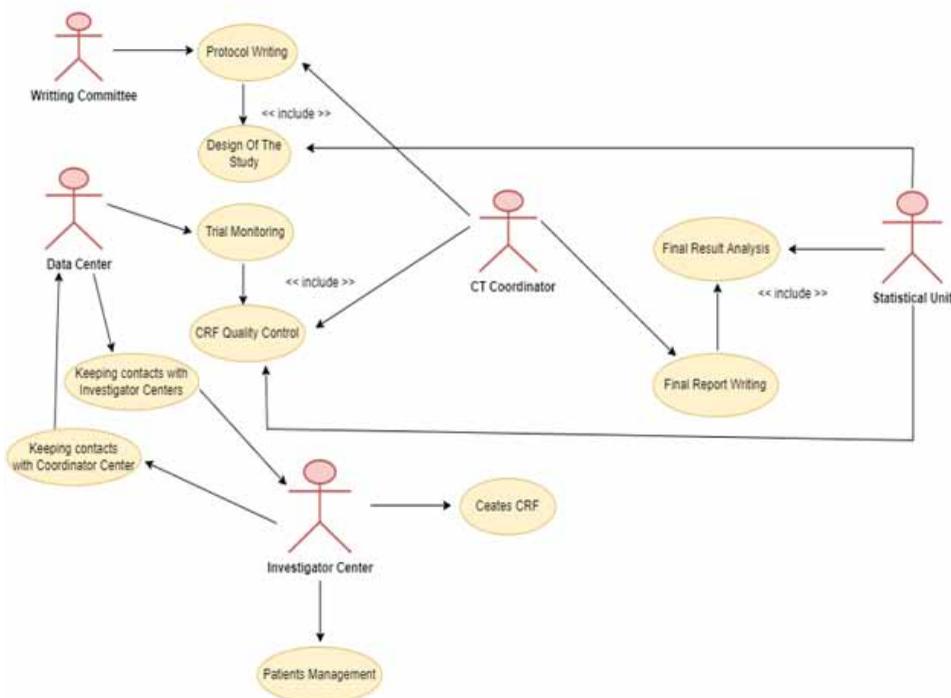
## 5. RESULTS AND DISCUSSIONS

As per the analysis done, following comparisons we get from that, all the images of the experimental values and the results are given below

Comparison between bioactivity and Molecular weight of Chemical Compounds

In the above give, graph author analyzed that, the bio-activity of a chemical compound to its molecular weight. The author fetched the molecular weight of AChE from the rdkit library. The blue and orange colors respectively showed active and Inactive IC50 values. It tends to be seen that the 2 bioactivity classes are traversing comparable synthetic spaces as apparent by the scatter plot of MW vs LogP (bioactive IC50 value) (Pl. refer Figure 19).

**Figure 18. Use Case Diagram**

The above graph shows the comparison between the experimental IC50 values and experimental IC50 values these predictions were done by modelling the Random forest regression and decision tree classifier machine learning algorithm. The above model is consist of 80% of dataset for training data and 20% of dataset for test data which are used respectively for experimental and prediction. Through the above graph, Author analyze and came to end with the result that as experimented value is high the prediction becomes much more accurate (Pl. refer Figure 21, 22).

**Figure 19. Comparison between bioactivity and Molecular Weight of the Chemical Compound**
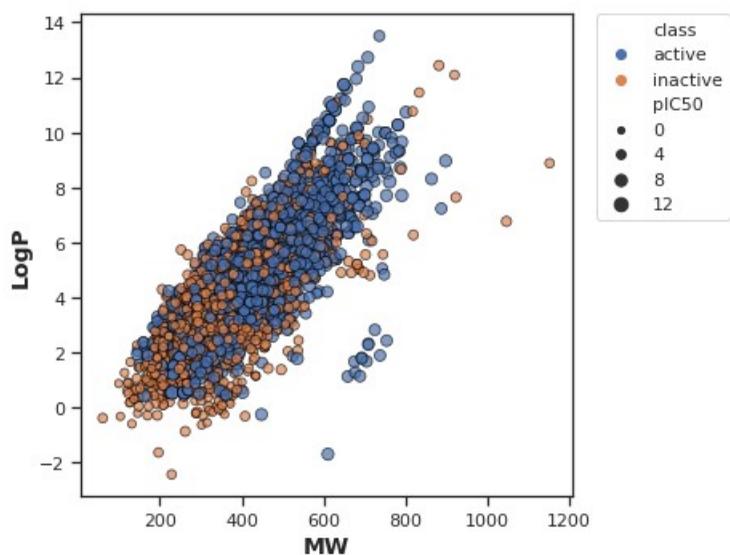


**Figure 20. Comparison between Experimental Values or Predicted Values (Decision Tree)**
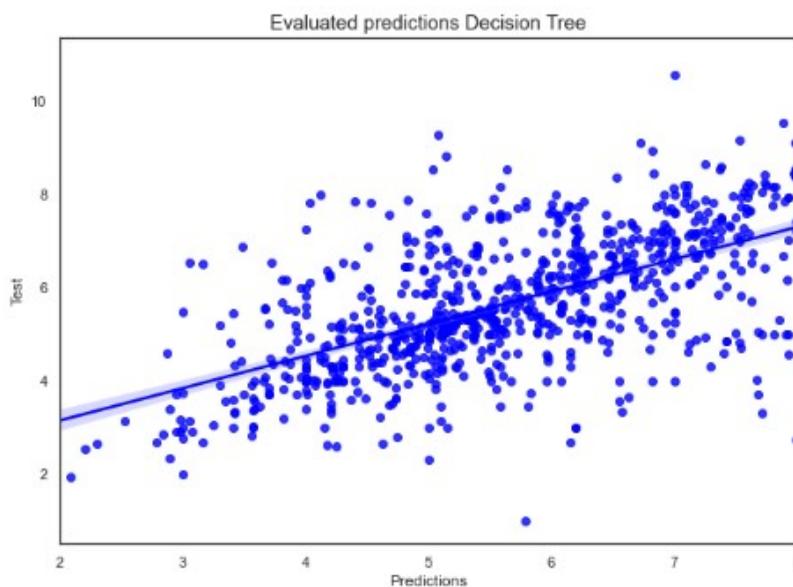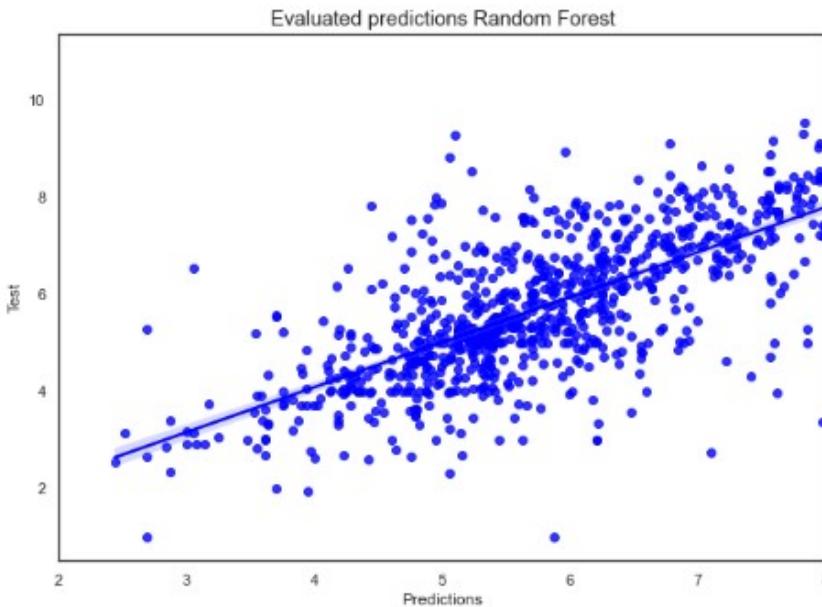
**Figure 21. Comparison between Experimental Values or Predicted Values (Random Forest)**



## 6. ANALYSIS OF IMPLEMENTED ML MODELS

## 6.1 The model used and experimental result for Train data (80% dataset)

Table 2. The model used and experimental result for Train data (80% dataset)

| Model | R-Squared | RMSE | Time Taken |
|---|---|---|---|
| Decision Tree Regressor | 0.86 | 0.57 | 0.16 |
| Random forest Regressor | 0.83 | 0.64 | 4.67 |
| K Neighbors | 0.65 | 0.92 | 2.98 |
| Linear Regression | 0.33 | 1.27 | 0.09 |

## 6.2 The model Used and Experimental Results for Test data (20% dataset)

Table 3. The model used and experimental result for Test data (20% dataset)

| Model | R-Squared | RMSE | Time Taken |
|---|---|---|---|
| Decision Tree Regressor | 0.28 | 1.31 | 0.16 |
| Random forest Regressor | 0.52 | 1.08 | 4.56 |
| K Neighbors | 0.46 | 1.14 | 0.97 |
| Linear Regression | 0.31 | 1.29 | 0.08 |

## 7. RESULTS

The above-given graph shows the implementation and accuracy of their model implementation of Decision Tree, Linear Regression, and Random Forest (Pl. refer Figure 22, 23, and Figure 24) (Pl. refer Table 2 and Table 3).

## 8. ALGORITHM COMPARISON

Random forest (RF) is a combinational classifier that is made out of a few decision trees. Momentarily, the principal thought behind RF is that as opposed to building a profound decision tree with a consistently developing number of hubs, which might be in danger of overfitting and overtraining the information, rather numerous trees are created to limit the change as opposed to boosting the precision. In that capacity, the outcomes will be noisier when contrasted with a thoroughly prepared decision tree, yet these outcomes are typically solid and robust, which is a quick execution of the RF calculation that was utilized for building the models.(Pl. refer Figure 25)

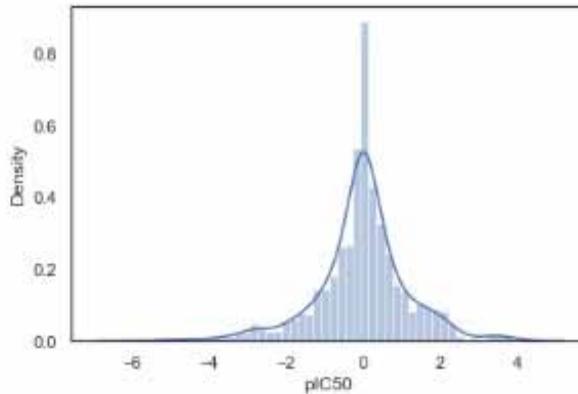**Figure 22. Implemented Decision Tree Model-1**
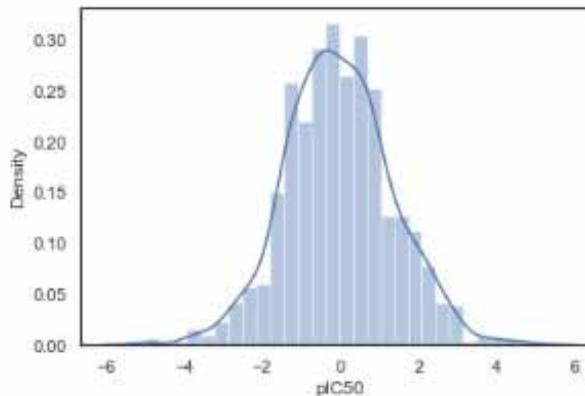


**Figure 23. Implemented Decision Tree Model-2**

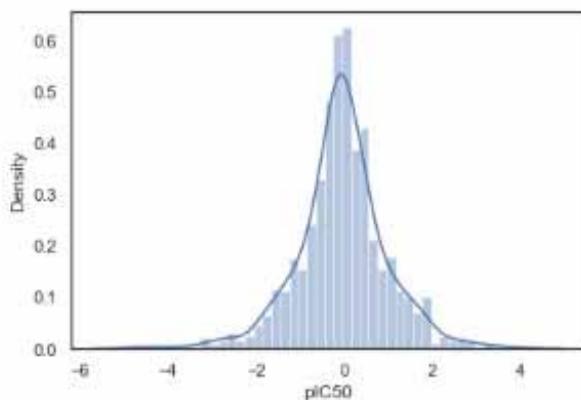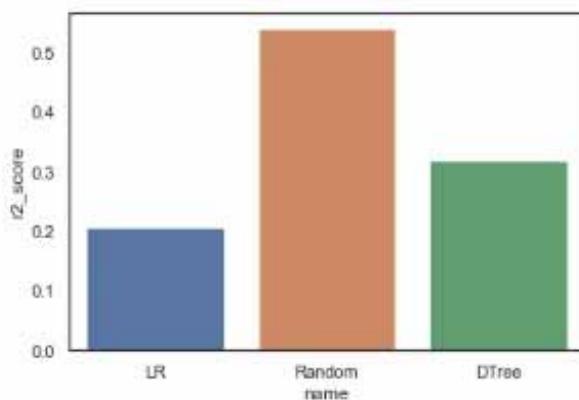**Figure 24. Implemented Linear Regression**



**Figure 25. Comparison between implemented machine learning model**



## 9. NOVELTIES

The essential point of drug detection is to observe novel atoms that are dynamic against an objective of restorative importance and that are not covered by any current licenses(patent). Because of the rising expense of innovative work in the later phases of drug discovery, and the expansion in drug competitors falling flat at these stages, there is a craving to choose the most assorted set of dynamic particles at the earliest phase of drug revelation, to expand the possibility observing an atom that can be enhanced into an effective drugs. Computational strategies that are both precise and proficient are one way to deal with this issue and can expand and analyze approaches in choosing which particles to take forward. There are many techniques of drug detection but yet now there has not been a finest and perfect method to detect the drug molecules within a less amount of time and having less cost overall and thus this Bio-Activity Detection of Drugs Compound system will help the pharmaceutical company to discover drugs very fast in easy and finest way to get accurate output without investing more money as compare to previous techniques.

## 10. RECOMMENDATIONS

In order to predict the bioactivity of drugs efficiently the cheminformatics organization should provide a large dataset of every year and also they should have mutual data provider so that more analysis can be done on the dataset so that chembl database can have more drug related data than can be used in predicting drug activity, drug analysis on different topics and can be develop similar and more usable model that can directly help to pharmaceutical industry and as a result a establishing of better health system in world for human development . At last, and fully intent on featuring this point, a joint exertion should be made in the search for and utilization of normalized systems. This point is significant for the fast change of scholarly outcomes to the business. Without a prosperity standard set in the cycles and procedures utilized, the outcomes got can't be stretched out to genuinely clinical errands. In this manner, the use of AI strategies should involve a powerful plan of the trials for their replicability by various scientists.

## 11. FUTURE RESEARCH DIRECTIONS AND LIMITATIONS

### 11.1 Limitations

Here, the author and their team work on the only ChEMBL id for CHEMBL1978 which represents the data for Inhibition of Cytochrome P450 19A1. We have limited chemical compounds we have to make predictions and do analysis only on them. The Author used the predefined dataset which we fetched from the ChEMBL database not using their own dataset. Numerous boundaries are changed during the preparation time of brain networks however a few hypothetical and useful systems are far off to upgrade these models.

### 11.2 Future Directions

- The presentation of deep learning strategies can straightforwardly impact the advancement of information mining in light of the fact that various deep neural networks are really prepared on an enormous volume of information. The principle point is to handle the exchange learning program issue.
- Web advancement was incorporated with clinical science to work on prescient power in independent direction and deep learning calculations about biomarkers, incidental effects in treatments, and restorative advantages.
- In clinical preliminaries, achievement is accomplished through the use of specific applications.

## 12. CONCLUSION

The advancement in the process of discovery of new drugs and prediction of bioactivity of drugs, the latest technology applied in this field is Machine Learning. Machine learning provides the feature of selection of data and extraction of data that might be fully relevant in order to extract the bioactivity properties of drugs. In the field of the pharmaceutical industry, there are different approaches to solving the problem of detecting the bioactivity of drugs. Previously, the main apparatus was the utilization of descriptors produced from the construction of little particles or peptides, that can take many trials and tremendous expenses at many levels yet after the bioactivity recognition utilizing ML would diminish the expense and time and furthermore decrease the stages and gives compelling outcome that will radically help in further developing the medical services framework.

A dataset collection involving 2,570 mixtures was utilized for the development of QSAR models. Especially, twelve arrangements of unique finger impression descriptors were benchmarked to find the best performing set. Before displaying, highlight determination was applied to eliminate collinear descriptors. Every one of the twelve models was then constructed utilizing an information split

proportion of 80/20 in which 80% of the informational index was utilized as the inner set and 20% as the outside set. This technique was iteratively acted in which every one of the 100 autonomous information parts were utilized for model development and the presentation results given in are the mean and standard deviation values got from these runs.

## ADDITIONAL READINGS

## KEY TERMS AND DEFINITIONS

**BioActivity-**Bioactivity is defined as the property of materials to develop a direct, adherent, and strong bond with bone tissue.

**Protein -** Protein is an important part of a healthy diet. Proteins are made up of chemical 'building blocks' called amino acids. Your body uses amino acids to build and repair muscles and bones and to make hormones and enzymes..

**Standard Value-** Bioactive peptides (BP) are organic substances formed by amino acids joined by covalent bonds known as amide or peptide bonds. Although some BP exists free in its natural source, the vast majority of known BP are encrypted in the structure of the parent proteins and are released mainly by enzymatic processes.

**Drug Discovery:-**Drug discovery is the process through which potential new medicines are identified. It involves a wide range of scientific disciplines, including biology, chemistry and pharmacology.

**Chemical Compound:-**A chemical compound is a chemical substance composed of many identical molecules (or molecular entities) composed of atoms from more than one element held together by chemical bonds.

**QSAR:-**Quantitative structure-activity relationship (QSAR) is a computational modeling method for revealing relationships between structural properties of chemical compounds and biological activities.

**rdkit:-** RDKit is a collection of cheminformatics and machine-learning software written in C++ and Python.

**AChE:-** Acetylcholinesterase inhibition

## ACKNOWLEDGEMENTS

# REFERENCES

ChEMBL. (2022). ChEMBL Data web services. ChEMBL. https://chembl.gitbook.io/chembl-interface-documentation/web-services/chembl-data-web-services

Chen, M., Yang, F., Kang, J., Gan, H., Lai, X., & Gao, Y. (2018). Discovery of molecular mechanism of a clinical herbal formula upregulating serum HDL-c levels in treatment of metabolic sy*ndrome by in vivo and computational stu*dies. Bioorganic & Medicinal Chemistry Letters, 28(2). https://www.sciencedirect.com/science/article/abs/pii/S0960894X17311320?via%3Dihub

Fang, J., Li, Y., Liu, R., Pang, X., Li, C., Yang, R., He, Y., Lian, W., Liu, A., & Du, G. (2015). Discovery of Multitarget-Directed Ligands against Alzheimer's Disease through Systematic Prediction of Chemical–Protein Interactions. American Chemical Society. https://pubs.acs.org/doi/10.1021/ci500574n

Goldblum, A., Yoshimoto, M., *& Hansch, C. (1981). Quantitative structure–activity rela*tionship of phenyl N-methylcarbamate inhibition of acetylcholinesterase. Journal of Agriculture and Food Chemicals. https://pubs.acs.org/doi/abs/10.1021/jf00104a017

Mohs, R. & Greig, N. (20*17). Drug discovery and development: Role of basic biological research. Alzheimer*s Dement, 3(4). https://www.ncbi.nlm.nih.gov/pmc/articles/PMC5725284/

Prado-Prado, F. J., Escobar, M., & Garcia-Mera, X. (2013). Review of Bioinformatics and Theoretical Studies of Acetylcholinesterase Inhibitors. *Current Biometrics, 8(4). http://w*ww.eurekaselect.com/article/54951

Racchi, M., Mazzuccehlli, M., Porrello, E., Lanni, C., & Govoni, S. (2004). Acetylcholinesterase inhibitors: novel activities of old molecules. Pharmacological Research, 50(4). https://www.sciencedirect.com/science/article/abs/pii/S1043661804000908?via%3Dihub

Steinbeck, C. *(2006). Recent Developments of The Chemistry Development Kit (CDK). Current Pharmaceuti*cal Design, 12(17). https://www.researchgate.net/publication/6987061_Recent_Developments_of_The_Chemistry_Development_Kit_CDK_-_An_Open-Source_Java_Library_for_Chemo-_and_Bioinformatics

Brogi, S., Ramalho, T. C., Kuca, K., Franco, J. M., & Marian, V. (Au*g., 2020). Editoria*l:*In s*ilico Methods for Drug Design and Discovery. doi:10.3389/fchem.2020.00612

Choudhary, A., Prakash, A., Nand, P., & Jain, V. (2021).Chapter 8. Examining the Effect of Ashes of Vedic Homa and Its Scientific Impacts on AQI with Social Perspective*s: An ML and CPS Based Experimental Study* for Delhi-NCR Zone amidst Pandemic Threats. In Intelligent Information Retrieval for Healthcare Systems. NOVA. . https://novapublishers.com/shop/intelligent-infor*mation-retrieval-for-*healthcare-systems/10.52305/ULIQ1795

Dara, S., Dhamercherla, S., Jadav, S.S., Babu, CH.M., Ahsan, M.J., (Aug., 2021). A review on Machine Learning in Drug *Discovery. Artificial Intelligence Review*., 55, 1947–1999 (2022). .10.1007/s10462-021-10058-4

Dara, S., Dhamercherla, S., Jaday, S., Babu, C.H., & Ashan, M. (2021). Machine Learning drug discovery: A review. Artificial Intelligence Review. https://l*ink.springer.com/article/10.1007/s10462-021-1*0058-4#Abs1

Eriksson, L. & Johansson, E. (1996). Multivariate design and modeling in QSAR. Chemometrics and Intelligent Laboratory Systems, 34(1), 1-19. https://www.sciencedirect.com/science/article/abs/pii/0169743996000238?via%3Dihu*b*

*Gupta, R., Srivastava, D., Sahu*, *M.*, Tiwari, S., Ambasta, R. K., & Kumar, P. (2021). Artificial intelligence to deep learning: Machine intelligence approach for drug discovery. Molecular Di*versity, 2021(25), 1315–1360. doi:10.1007/s11030-*021-10217-3 PMID:33844136

Khan, M. K. (2020). Technological advancements and 2020. Telecommun Syst., 73, 1–2, https://link.springer.com/article/10.1007/s11235-019-00647-8

Kumar, S. Gs., P., (*Jan. 2014). Health Promoti*on: An Effective Tool for **Global Health., Indian Journal of Community Medicine., 37(1):5-12., https://**www.ijcm.org.in/article.asp?issn=0970-0218;year=2012;volume=37;issue=1;spage=5;epage=12;aulast=Kumar

Manne.R., (April, 2021). Machine Learning Techniques in Drug Discovery and De*velopment, International Journal of Applied Research,, 7(4), 21-28. https://www.allresearchjournal.co*m/archives/?year=2021&vol=7&issue=4&part=A&ArticleId=8455

Moumtzoglou, A. (2022).. . International Journal of Reliable and Quality E-Healthcare, 10(1). Advance online publication. doi:10.4018/IJRQEH

Park, J., Beck, B. R., Kim, H., Lee, S., & Kang, K. (2022). A Brief Review of Machine Learning-Based Bioactive Compound Research. Applied Sciences (B*asel, Switzerland), 12(6), 2906. doi:10.3390/app12062906*

*Patel, L., Shukla*, T., Huang, X., Ussery, D., & Wang, S. (Nov., 2020), Machine Learning Methods in Drug Discovery. Multidisciplinary Digital Publishing Institute, 25(22). https://www.mdpi.com/1420-3049/25/22/5277

Patel, V., & Shah, M. (Nov., 2021). A comprehensive study on artificial intelligence and machine learning *in drug discovery and drug development. Intelligent Medicine., 10, (S2667*-1026(21)00106-6). https://www. sciencedirect.com/science/article/pii/S2667102621001066#bib0012

Rastogi, R. & Garg, P. (2021). Investigation of Air Quality Prediction and Analysis Amidst Pandemic challenges with Indian Science of Agnihotra: M*L Based Study for Healthcare 4.0. Interdisciplinary Environmental Review, 21(1), p*p. 66-85., https://www.inderscience.com/info/inarticle.php?artid=11528510.1504/IER.2021.115285

Rastogi, R. (2020) Yajna and Mantra Science Bringing Health and Comfort to Indo-Asian Public: A Healthcare 4.0 Approach and Computational Study. In: Jain V., Chatterjee J. (eds) Machine Learning with Health Care Perspective. Learning and Analytics in Intelligent Syst*ems, (vol* 13, pp. 357-390). Springer. https://link.springer. com/chapter/10.1007%2F978-3-030-40850-3_15

Rastogi, R., & Saxena, M., Chaturvedi*, D. K., Gupta, M., Rastogi, A. R., Rastogi, M., Sharma, A*., & Sagar, S. (2020). Kirlian Experimental Analysis & IoT: Part I. In Volume 10, Issue 1 in 2021 first Quarter, International Journal of Reliable and Quality E-Healthcare (IJRQEH), doi:, https://www.igi-global.com/journal/international-journal-reliable-quality-healthcare/4466010.*4018/IJRQEH*

*Rastogi, R., & Saxena, M. Chaturvedi, D.* K., Gupta, M., Rastogi, A. R., Rastogi, M., Sharma, A., Sagar, S. (2020). Kirlian Experimental Analysis & IoT: Part 2. International Journal of Reliable and Quality E-Healthcare (IJRQEH), 10(1). doi:, https://www.igi-global.com/journal/international-jour*nal-reliable-quality-healthcare/4466010.4018/IJRQEH*

*Rastogi, R.,* Chaturvedi, D. K., Saxena, M., Gupta, M., Guta, N., Yadav, S., Rustagi, D., & Sharma, P. (2021). Application of Kirlian Captures and Statistical Analysis of Human Bioelectricity and Energy of Different Organs: Observations and Graphical Notations. International Journal of He*alth Systems and Translational Medicine (IJHSTM), 1(1) doi: https://www*.igi-global.com/journal/international-journal-health-systems-translational/24119910.4018/IJHSTM

Rastogi, R., Sagar, S., Tandon, N., Garg, P., & Rastogi, M. (2021). Treatment Case Studies and Emissions Analysis of Wood in Yagya: Integrat*ing Spirituality and Healthca*re With Science. International Journal of Biomedical and Clinical Engineering (IJBCE), 10(2), 29-43. .10.4018/IJBCE.2021070103

Rastogi, R., Sagar, S., Tandon, N., Rajeshwari, T. (2022). Examining the Effect of Ashes of Vedic Homa a*nd its Scientific Impacts on AQI w*ith Social Perspectives: An ML and CPS based Experimental Study for Delhi-NCR Zone amidst Pandemic Threats. In V. Jain (ed.) Trends and Technologies in Ontology Based Information Retrieval for Health Care System. NOVA Science Publishers. *https:/easychair.org/con*ferences/?conf=ttobir2021.

Rastogi, R., Sagar, S., Tandon, N., Rajeshwari, T., & Singh, B. (2022). Computational Statistics on Stress Patients with happiness and Radiation Indices. Vedi*c Homa Therapy. https://scholar.google.com/*citations?user=-luGOX4AAAAJ&hl=en

Rastogi, R., Saxena, M., Chaturvedi, D.K., Gupta, M., Jain, P., Jain, R., Jain, M., Sharma, V., Sangam, U., Singhal, P. & Garg, P. *(2022). Indian Science of Yaj*na and Mantra to Cure Different Diseases: An Analysis Amidst Pandemic with a Simulated Approach. In A. Suresh, S. Vimal, Y.H. Robinson, D.K. Ramaswami and R. Udendhran (eds.) Bioinformatics and Medical Applications https://doi.org/10.1002/9781119792673.ch12

Rastogi, R., Saxena, M., Chaturvedi, *D. K., Gupta, M., Rastogi, M., Rustagi, D., Gaur,* V., Kohli, V., Srivastava, P., Jain, M., & Kumar, P. (2020). AI-Based Analysis for Novel Covid-19 and Its Treatment Through Yajna and Mantra Science. IJRQEH, 9(4), 75-98. https://www.igi-global.c*om/journal/international-journal-reliable-quality-healthcare/44660*

*Rastogi, R., Saxena, M., Sagar, S., Rajeshwari, T., Tandon, N*., Sharma, M., & Rastogi, M. (2022). Parameter Analysis of Electronic Gadgets during Homa: Smart City Healthcare Perspective with CPS and IoT. International Journal of Ap*plied Research on Pu*blic Health Management (IJARPHM), 7(1), 16. . https://www.igi-global.com/ journal/international-journal-applied-research-public/21490210.4018/IJARPHM.2022010101

Reboredo, P. C., Blanco, J. L., Fernandez, *N.R., Cedron, F., No*voa, F.J., Carballal, A., Maojo, V., Pajos, A., & Lojano, C.F., (Aug., 2021). A review on machine learning approaches and trends in drug discovery. Research Network of Computational and Structur*al Biotechnology, 19, (4538*-4558), .10.1016/j.csbj.2021.08.011

Rollinger, J. M., Stuppner, H., & Langer, T. (2008). Virtual screening for the discovery of bioactive natural products, Progress in Dr*ug Research, 65(211), 213-49. . h*ttps://link.springer.com/content/pdf/10.1007/978-3-7643-8117-2_6.pdf10.1007/978-3-7643-8117-2_6

Siddiqui, M., AlOthmank, Z., & Rahman, N. (Feb., 2017). Analytical Techniques in Pharmaceutical Analysis: A review., Arabian Journal of Chemistry., 10, (S1409-S1421)., https://www.sciencedirect.com/scien*ce/article/pii/S1878535213001056*

*Yan, A.* & Wang, K. (2012). Quantitative structure and bioactivity relationship study on human acetylcholinesterase inhibitors. Bioorganic & Medicinal Chemistry Letters, 22(9). https://www.sciencedirect.com/science/article/abs/pii/S0960894X12003332?via%3Dihub

Zhang, Y., Ye, T., Xi, H., & Juhas, M. (Oct. 2021). Deep Learning-Dr*iven Drug Discovery, Tackling Severe Acute R*espiratory Syndrome Coronavirus 2. Frontiers in Microbiology. 10.3389/fmicb.2021.739684

## APPENDIX

## Datasets

**Figure 26. Sample Dataset**

| molecule_chembl_id | canonical_smiles | standard_value |
|---|---|---|
| CHEMBL341591 | CC12CCC(O)CC1=CCC1C2CCC2(C)C(CC3CN3)C( | 7100 |
| CHEMBL2111947 | C[C@]12CC[C@H]3[C@@H](CC=C4C[C@@H]( | 50000 |
| CHEMBL431859 | CCn1c(C(c2ccc(F)cc2)n2ccnc2)c(C)c2cc(Br)ccc2 | 238 |
| CHEMBL113637 | CCn1cc(C(c2ccc(F)cc2)n2ccnc2)c2ccccc21 | 57 |
| CHEMBL112021 | Clc1ccccc1Cn1cc(Cn2ccnc2)c2ccccc21 | 54 |
| CHEMBL324070 | Cc1ccc(S(=O)(=O)n2cc(C(c3ccccc3)n3ccnc3)c3c | 5400 |
| CHEMBL41761 | CCn1ccc2cc(C(c3ccc(F)cc3)n3ccnc3)ccc21 | 41 |
| CHEMBL111868 | Cn1cc(C(c2ccc(F)cc2)n2ccnc2)c2cc(Br)ccc21 | 78.5 |
| CHEMBL111888 | CCn1cc(C(c2ccc(F)cc2)n2ccnc2)c2cc(Br)ccc21 | 51.8 |
| CHEMBL112074 | CCn1ccc2cc(C(c3ccccc3)n3ccnc3)ccc21 | 205 |
| CHEMBL324326 | N#Cc1ccc(Cn2cc(Cn3ccnc3)c3ccccc32)cc1 | 50 |
| CHEMBL37321 | CCCCCCN1C(=O)CCC(CC)(c2ccncc2)C1=O | 6600 |
| CHEMBL353068 | c1ccc2c(c1)CCC1C(c3cc[nH]n3)C21 | 51000 |
| CHEMBL41066 | CCCCCCCC1(c2ccncc2)CCC(=O)NC1=O | 3200 |
| CHEMBL166709 | O=C1/C(=C/c2cccnn2)CCc2ccccc21 | 250000 |
| CHEMBL424556 | O=C1/C(=C/c2ccnnc2)CCc2ccccc21 | 103000 |
| CHEMBL1630273 | C[C@]12CC[C@H]3[C@@H](CC=C4C[C@H](O)( | 6800 |
| CHEMBL1630261 | C[C@]12CC[C@H]3[C@@H](C[C@@H](O)C4= | 50 |
| CHEMBL169251 | O=C1/C(=C\c2c[nH]cn2)CCc2ccccc21 | 170 |
| CHEMBL168636 | O=C1/C(=C/c2cccnc2)CCc2ccccc21 | 9200 |
| CHEMBL90585 | O=C1/C(=C/c2ccncc2)CCc2ccccc21 | 4600 |
| CHEMBL1629805 | C[C@]12CC[C@H]3[C@@H](CCC4=CCCC[C@( | 49 |
| CHEMBL433728 | c1ccc2c(c1)CCC1C(c3ccsc3)C21 | 250000 |
| CHEMBL38877 | CCCN1C(=O)CCC(CC)(c2ccncc2)C1=O | 31000 |

**Figure 27. Sample Dataset**

| | | |
|---|---|---|
| CHEMBL38877 | CCCN1C(=O)CCC(CC)(c2ccncc2)C1=O | 31000 |
| CHEMBL169449 | O=C1c2ccccc2CCC1C(O)c1cnccn1 | 250000 |
| CHEMBL39275 | CCCCN1C(=O)CCC(CC)(c2ccncc2)C1=O | 40000 |
| CHEMBL39513 | CCCCCN1C(=O)CCC(CC)(c2ccncc2)C1=O | 21000 |
| CHEMBL2112738 | CC(=O)O[C@H]1CCC[C@@]2(CO)C1=CC[C@H | 11000 |
| CHEMBL289116 | CCC1(c2ccncc2)CCC(=O)NC1=O | 10000 |
| CHEMBL289116 | CCC1(c2ccncc2)CCC(=O)NC1=O | 45000 |
| CHEMBL39782 | CCC1(c2ccncc2)CCC(=O)N(C)C1=O | 30000 |
| CHEMBL304903 | c1ccc2c(c1)CCC1C(c3ccncc3)C21 | 370 |
| CHEMBL488 | CCC1(c2ccc(N)cc2)CCC(=O)NC1=O | 8000 |
| CHEMBL488 | CCC1(c2ccc(N)cc2)CCC(=O)NC1=O | 14000 |
| CHEMBL168434 | O=C1/C(=C/c2c[nH]cn2)CCc2ccccc21 | 260 |
| CHEMBL352645 | c1ccc2c(c1)CCC1C(c3c[nH]cn3)C21 | 600 |
| CHEMBL2112739 | CC(=O)O[C@H]1CCC[C@@]2(C)C1=CC[C@H]1 | 1000 |
| CHEMBL440930 | c1cncc(C2C3CCc4ccccc4C32)c1 | 13000 |
| CHEMBL1630274 | C[C@]12CC[C@H]3[C@@H](CC=C4[C@@H](C | 31 |
| CHEMBL1629804 | C[C@]12CC[C@H]3[C@@H](CCC4=CCCC[C@@ | 60 |
| CHEMBL39661 | CCCCCCCCN1C(=O)CCC(CC)(c2ccncc2)C1=O | 800 |
| CHEMBL39661 | CCCCCCCCN1C(=O)CCC(CC)(c2ccncc2)C1=O | 5000 |
| CHEMBL38550 | CCCCCCCCCN1C(=O)CCC(CC)(c2ccncc2)C1=O | 3000 |
| CHEMBL39152 | CCCCCCC1(c2ccncc2)CCC(=O)NC1=O | 3600 |
| CHEMBL166789 | c1ccc(C2C3CCc4ccccc4C32)nc1 | 250000 |
| CHEMBL3349856 | C[C@]12CC[C@H]3[C@@H](CC=C4[C@H](O)( | 860 |
| CHEMBL1630275 | CO[C@H]1CCC[C@@]2(CO)C1=CC[C@H]1[C@ | 150000 |
| CHEMBL1630267 | C[C@]12CC[C@H]3[C@@H](C[C@@H](O)C4= | 940 |

*Dr. Rohit Rastogi received his B.E. C. S. S. Univ. Meerut, 2003. Master's degree in CS of NITTTR-Chandigarh from Punjab University. He received his doctoral degree from the Dayal Bagh Educational Institute in Agra, India. He is serving as Associate Professor in the CSE department of ABES Engineering College, Ghaziabad, India. He has won awards in several areas, including improved education, significant contributions, human value promotion, and long-term service. He keeps himself engaged in various competition events, activities, webinars, seminars, workshops, projects and various other educational learning forums. He has guided around 40 B. Tech. students' projects and 5 M. Tech. Thesis. He is editor and reviewer member of several international Journals and conferences. He has 100+ publications in journals and conferences of International repute. He strongly believes that Transformation starts within self.*

*Yash Rastogi is engineering student in AKTU Univ. Presently he is a B.Tech. Final Year student of CSE in ABESEC, Ghaziabad. India. He is presently working on Yagya and Mantra therapy and its analysis by Machine Learning. He has keen interest in Google surfing. His hobbies is playing badminton and reading books. He is young, talented and dynamic. He is placed in a good IT company and strong interest in Data sciences. He is versatile and smart personality and wish to serve country through IT sector. He has developed some good analysis for different data science projects.*

*Saurav Kr. Rathauris an engineering student in Dr, APJ Abdul Kalam Technical University. Presently he is a Final Year student of Computer Science and Engineering in ABES Engineering College Ghaziabad, India. He has also completed his Diploma in Information Technology from Jawahar Lal Nehru Polytechnic, Mahmudabad-Sitapur, India. He has a keen interest in Web Development and Python Programming. He has also developed some good projects in the Web Development domain. He is placed in one of the top IT company in India. His hobbies are playing cricket, badminton and traveling. He is Ambivert, i.e., he has the qualities of both extrovert and introvert. He enjoys spending time with his family and being alone in Nature and searching for the betterment of himself.*

*Vaibhav Srivastavais a student of B.Tech. (CSE) in ABES Engineering College, Ghaziabad, which is affiliated to AKTU. He is currently in Final Year of the CSE branch of engineering .He is a dedicated and hardworking person having interest in learning new technology such as machine learning and also developed a good project, one is Clinic System app using NetBeans and Xampp server. He is placed in one of the top IT MNC companies in India .Beside this, he has a talent of writing poems and creative thinking of life that motivates everyone . He is looking for the best career ahead.*