

Improving Drug Discovery by Utilizing Machine Learning Algorithms and Biological

Activity Data of Target Proteins

Mathematics & Computer Science

Lagnajeet Panigrahi

Shrewsbury High School

Instructor: Catherine Phillips

75 Cypress Ave, Shrewsbury, MA 01545

Table of Contents

| | |
|-----------------------|----|
| Abstract | 3 |
| Introduction/Research | 4 |
| Experimental Design | 9 |
| Data Analysis/Results | 13 |
| Conclusion | 24 |
| Appendix | 25 |
| References | 30 |

Acknowledgement of help Received

This project was an independent project conducted at home and school. No outside assistance from professors or mentors was used during the development of this project. The only guidance received was from Shrewsbury High School Chemistry Teacher, Catherine Phillips during the creation of the research paper and abstract about structure and what information needed to be included.

Abstract

Humans have been battling the age-old battle against illness for centuries. Only recently have humans made significant progress in averting illness through the creation of drugs and the field of drug discovery. To create a drug and deem it usable requires a rigorous process that examines one attribute in particular: safety. Unfortunately, ninety percent of drug development fails despite the many successful strategies used. Current chemical methodologies rely on a hit-and-miss approach where large amounts of drugs are analyzed for their properties by hand. This process is expensive, time-consuming, and often inaccurate. This project aims to address this through the implementation of various machine-learning algorithms that can predict the IC₅₀ and pIC₅₀ values of possible drug candidates. These values are essential in determining the quantity of a drug needed to inhibit a biological process by half but also are important in determining the toxicity of a drug and how it impacts patients. Using regression-based machine learning models and bioactivity data of compounds and target proteins, these values can be predicted and outputted. Upon completion of this project, we developed multiple regression models for the target protein of the SARS coronavirus. After statistical analysis, the best model was chosen: DecisionTreeRegressor. This model had a root mean squared error score of 0.34 and an R-squared score of 0.9019. This implies that this model fits the situation and makes accurate predictions. We conclude that drug discovery can become quicker, more accurate, and cost-effective through the implementation of machine learning algorithms.

Keywords: Drug discovery, IC₅₀, pIC₅₀, machine learning, regression, bioactivity, r-squared score, root mean squared error score

Introduction/Research

In this expanding world where illness and disease are increasingly common, efficient, accurate, and cost-effective approaches to drug discovery are becoming more necessary. Inaccuracies and other sources of error can greatly delay drug development which can have disastrous impacts on individuals around the world. It was calculated that there were 29 life-years lost in North America alone per one hour of delay in drug approval (Helwick, 2015) and 90 percent of clinical drug development fails (Sun et al., 2022). A major step in drug development is determining the toxicity of the drug, its potency, and overall, its safety. Current methods require careful experimentation and individual assessment of large amounts of potential drug candidates by hand (Blanco-González et al., 2023). Although current methods are good, they are prone to errors, are costly, and are time-consuming (Blanco-González et al., 2023). Consequently, together, the high likelihood of failure combined with the inherent problems with current methods create a system that fails to perform effectively and efficiently. This is detrimental for both developers— as they lose time and money— and patients who do not receive proper treatment at the correct time. To overcome this, a new and innovative system is required: one that can provide information about the potency and toxicity of a certain molecule or drug and how it absorbs into the body in a fast, cost-effective, and accurate manner. The proposed method that this project will employ to solve this problem is the use of regression-based machine learning algorithms in Python with preprocessed data to give information on drug toxicity and potency by predicting IC_{50} and pIC_{50} values. To build the model, any target protein can be utilized (the only difference between models of different target proteins is what they are trained and tested on). For this project, the SARS coronavirus 3C-like proteinase target will be used. The model's success will be determined by statistical analysis.

This type of project will greatly impact the scientific community as it will not only provide a model to predict IC₅₀ and pIC₅₀ values but pave the way for the use of machine learning and related technologies for drug discovery. This project will build on past research on this problem that utilized different types of models and architectures.

Modern-day drug discovery is a process in which new drugs and therapeutic entities are created based on computational, experimental, translational, and clinical models (Zhou et al, 2017). The drug creation process can be broken up into six main steps: target identification, target validation, “hit” identification, lead discovery, “hit-to-lead” (H2L) stage, and lead optimization (Alvarellos, 2023). Target identification is the stage where the target protein for the disease or illness that the drug will target is chosen (Alvarellos, 2023). During target validation, the selected biological molecule is analyzed to confirm whether or not it plays a role in disease development and whether it needs to be targeted (Alvarellos, 2023). During “hit” identification, molecules that are active, inactive, or intermediate at the target biological molecule (Alvarellos, 2023). In other words, at this stage, it is determined which molecules have impacts on the biological molecule of the disease. In lead discovery, each molecule is screened for its properties pertaining to safety, efficacy, and potency (Alvarellos, 2023). In the H2L stage, scientists optimize the characteristics of the molecule they want based on results from the lead discovery stage and test to ensure they have reached their goals (Alvarellos, 2023). Finally, during the lead optimization stage, the most promising candidates get further optimized (Alvarellos, 2023). The problem with these conventional methods is that they are very difficult, lengthy, costly, and inaccurate (Zhou et al, 2017). This problem is exemplified when considering that each clinical trial can cost anywhere from 30 million dollars to 310 million dollars and fewer than ten percent of these trials succeed (Alvarellos, 2023). Imperfect representations by humans, errors in

experiments to determine the properties of drugs, and false positives all contribute to the failure of many clinical trials (Alvarellos, 2023). Machine learning and data-driven drug discovery have shown promise and could be a possible solution due to their characteristics: unbiased data, data standardization through preprocessing, and overall because of consistency and high performance if done right (Alvarellos, 2023). In this project, for the most part, the steps of H2L and lead discovery will be addressed and improved through machine learning as these steps are the most critical and most prone to error. A newer, more effective method called drug design that employs machine learning is a possible solution to the discussed problem and is the avenue this project will be taking. Drug design employs machine learning algorithms, computational methods, and bioinformatic methods to analyze drugs and molecules based on the knowledge of a biological target (Zhou et al, 2017). During the lead discovery and H2L stages, many factors are taken into consideration, but in this project, the factors of potency and drug toxicity will be taken into account. These metrics directly provide information about how well a patient will respond to the drug in question and what the drug's effect will be (Zhou et al, 2017). This project aims to be able to take a molecule, analyze it, and output values that indicate potency and toxicity effectively, accurately, and on time.

Drug toxicity and potency can be described by IC₅₀ values and pIC₅₀ values. IC₅₀ stands for inhibitory concentration at fifty percent (Visikol, 2022). pIC₅₀ values are the same values in concept, except they have had a negative logarithmic function applied (Navre, 2019). This is because the IC₅₀ values follow a logarithmic phenomenon and pattern like pH and pOH values (Navre, 2019). In addition, once converted to pIC₅₀, the data becomes easier to read and interpret. These values can be indicative of how much of a drug is needed to create an impact, what the drug's potency is, and its systemic toxicity (Visikol, 2022). For example, a lower IC₅₀

value can indicate higher drug potency as less of the drug is required to inhibit a biological process and is less toxic because less of it is needed (Berrouet et al., 2022). There are many methods of calculating IC₅₀ values such as using raw data from various assays and using curve fitting software and point-by-point analysis (isirv Antiviral group, n.d.). For the purposes of this project, we will be using a different methodology. IC₅₀ and pIC₅₀ values can be calculated through the SMILES notation of a molecule or drug and other bioactivity data (whether the drug or molecule is active or inactive at the target biological organism) (Liu et al., 2019). SMILES is the Simplified Molecular Input Line Entry System that converts chemical formulas into a notation composed of strings that are easily readable by a computer (US EPA, 2012). The data used in this project will be sourced from ChEMBL. ChEMBL is a curated database of molecules with drug-like properties (ChEMBL, n.d.). Since this machine learning project will be tested with SARS coronavirus data, the bioactivity data will contain bioactivity data and SMILES notation of certain molecules that have been tested with the target protein of SARS coronavirus. This will be the data that will be processed and will be used by the machine learning algorithm during training and testing.

Machine learning will be leveraged in this project with the data aforementioned. Machine learning is a branch of artificial intelligence and computer science that utilizes data to create algorithms that learn, predict, and improve (IBM, 2023). Within the field of machine learning, there is supervised learning and unsupervised learning. The difference between these types of learning for a machine learning model is that supervised learning uses labeled data outputs and inputs while unsupervised learning does not (Seldon, 2021). This project will employ supervised learning as the data is labeled by SMILES notation and bioactivity and use regression-based machine learning models. Regression-based machine learning models take in data and use

mathematical functions to establish relationships between certain inputs and outputs to eventually create a classifier and be able to predict outcomes (Kumara, 2023). This project will test many different machine learning regression models using LazyRegressor and then determine the best one to model drug design for the scenario presented within the project. LazyRegressor is a library within Python that can automate machine learning training, testing, and analysis of models (Patil, 2021).

This project is fully dependent on the use of the programming language Python and the data sourced from ChEMBL. As a result, this project employs the use of many libraries and applications. This project will be programmed on Google Colab rather than conventional IDEs such as Pycharm or Visual Studio Code. Google Colab is an online cloud-based Jupyter Notebook environment (Sharma, 2024). There are many benefits to using this as opposed to other IDEs. The benefit of utilizing Google Colab is that it is cloud-based, so all of the machine learning training, testing, and simple running of code is run on GPUs and TPUs hosted by Google (Sharma, 2020). This allows for the model to be trained quickly and efficiently while also not being hindered by the performance of a personal device. Next, this project requires the use of many Python libraries. The most prominent libraries used in this project were Pandas, Numpy, RDkit, Matplotlib, PaDEL, Seaborn, and LazyRegressor.

The main objective of this computer science project is to use regression-based machine learning models to use bioactivity data and SMILES notation of drugs and molecules to predict the IC₅₀ and pIC₅₀ values. This will give valuable insight into the properties of the drug such as toxicity and potency in a cost-effective, accurate, and timely manner. The model will be tested for how well it can predict and fit the situation. This will undoubtedly allow for quicker drug development for scientists and allow for the saving of millions of lives.

Experimental Design

Location

This experiment will be executed indoors at home or school with access to the internet. It will take place on a 16-inch Apple Macbook Pro. This laptop is equipped with an Apple Silicon M2 Pro chip, 16 gigabytes of RAM, and 1 terabyte of SSD storage.

Constraints

Since this experiment is solely based on programming, there are no independent variables, dependent variables, or controls to consider. Instead of this, the project will be tested on accuracy metrics. One limitation to consider is the amount of data available. Many programs related to machine learning—such as the regression models in this project—require a training dataset to allow the program to learn how to identify what is being shown. A testing set is also required to calculate the accuracy of this program and to ensure that the program is functioning correctly. This creates a great demand for large datasets. Since datasets for bioactivity are not widely available (and when they are, they are very small), the program will not be able to use the most optimal amount of data, thereby unable to reach its full potential.

Safety

This entire project revolves around using a computer to program a machine learning model. The testing and training will also occur on a laptop. Therefore, there are no safety concerns to consider for the execution of this experiment.

Materials

For this experiment, a laptop or a type of personal computer able to run Google Colab is required. In this experiment, a 16-inch Apple Macbook Pro equipped with an Apple Silicon M2

Pro chip, 16 gigabytes of RAM, and 1 terabyte of SSD storage will be used. Since this experiment is performed using Google Colab, most computers will work. Next, the SARS Coronavirus 3C-like proteinase data needs to be downloaded from the ChEMBL Database. A Google account will be necessary to utilize Google Colab, but no installation is required as it can be accessed without any downloads. Finally, the necessary libraries need to be installed: Pandas, Numpy, RDkit, Matplotlib, PaDEL, Seaborn, SkLearn, ChEMBL web source, and LazyRegressor.

Procedure 1: Get Bioactivity Data

Once the ChEMBL web source is installed, it can be accessed within the code editor, and SARS Coronavirus 3C-like proteinase data can be downloaded for future use. Using the `target.search` function the data can be found for SARS Coronavirus 3C-like proteinase and saved as a CSV file. The dataset may have some missing data or other problems. To filter these out, the `data frame drop` function can be used. Next, create a new data frame with only the molecule ChEMBL ID, canonical smiles, and standard value data by using the selection function. The rest of the data is not necessary for this project. Next, we must identify each molecule's bioactivity as inactive, active, or intermediate. If the standard value for a molecule is above 10000 it can be considered active, if it is under 1000 it can be considered inactive, and any values in between are intermediates. Now save the new data as a CSV file that can be accessed in the next steps.

Procedure 2: Exploratory Data Analysis/Descriptor (Lipinski) Calculation

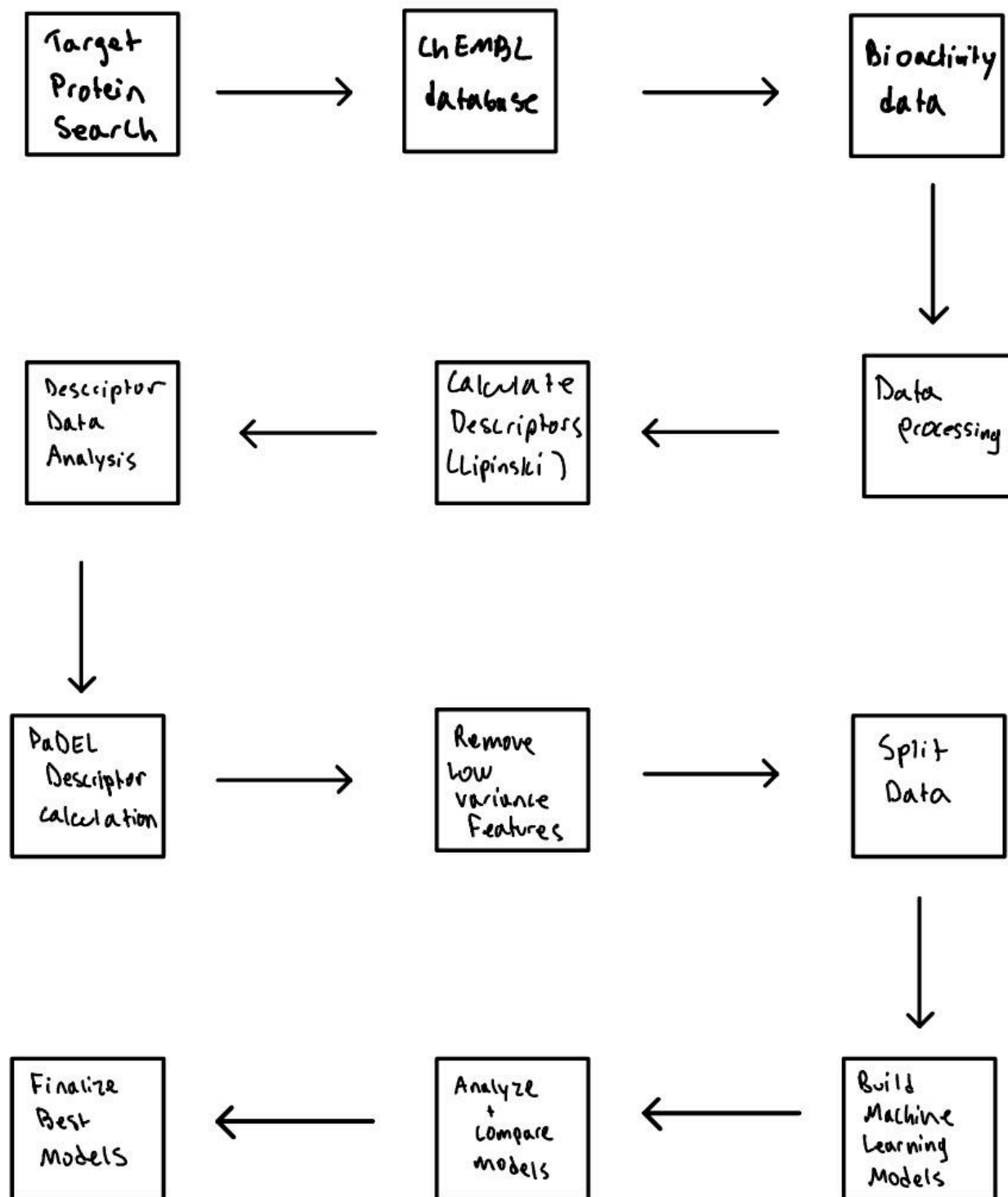
Once Rdkit has been installed and the previous steps have been completed import `chem` and `lipinski` from Rdkit. Next, use the Lipinski descriptor calculator outlined by Codeocean.com to calculate the descriptors (Codeocean.com, n.d.). This code can also be found in Appendix 2. Combine the outputted data with the CSV file from procedure 1. This dataset can be found in

Appendix 1. Next, to perform statistical analysis to determine which descriptors are indicative of active versus inactive, a Mann-Whitney U Test will be performed for statistical significance (Laerd, 2012). The alpha value will be 0.05. The code for this test can be found in Appendix 3. Perform the data analysis. Finally, the combined data will be converted into descriptor fingerprints using PaDEL for the regression model read as x and y values. Use the bash function to do this.

Procedure 3: Regression Models Creation

First import the created dataset from procedure 2 to use in the model creation. Next, remove the low variance data points by using the variance threshold function. Split the dataset into training and testing sets with 80 percent of the data going into training and the rest into testing using the split function. Define the x and y for training and testing and then use the LazyRegressor function to train and test models with the dataset. Output the results and choose the best model to model the pIC50 and IC50 values for molecules for SARS Coronavirus 3C-like proteinase.

Experiment Visualized:



Data Collection/Statistical Tests:

Descriptor data analysis

It is essential to know which of the descriptors calculated influence whether or not a molecule or drug is active or inactive. A Mann-Whitney U test is one method to determine this. It will determine the significance of the descriptor by calculating a p-value and seeing if it is smaller than the alpha value (alpha=0.05) that will be used. If the p-value is smaller than the alpha value, there is a statistically significant difference. Example data table:

| | Descriptor | Statistics | p | alpha | Interpretation |
|---|------------|------------|---|-------|--|
| 0 | LogP | - | - | 0.05 | Fail to reject null hypothesis or reject null hypothesis |

Machine learning regression model data analysis:

R-squared values and root mean squared error (RMSE) scores will be used to evaluate the model's performance. The R-squared value will describe how well the regression model predicts the data on a scale from 0 to 1 where 1 is perfect. In other words, it shows how well the model fits the data. The RMSE score displays the average difference between the model's prediction and the actual value. The goal of this project is to achieve an R-squared value of at least 0.9 and a maximum RMSE value of 0.5.

| Model | R-Squared Score | RMSE Score |
|-------|-----------------|------------|
| - | - | - |

Data Analysis/Results

Figure 1a. Amount Active and Inactive Molecules in dataset

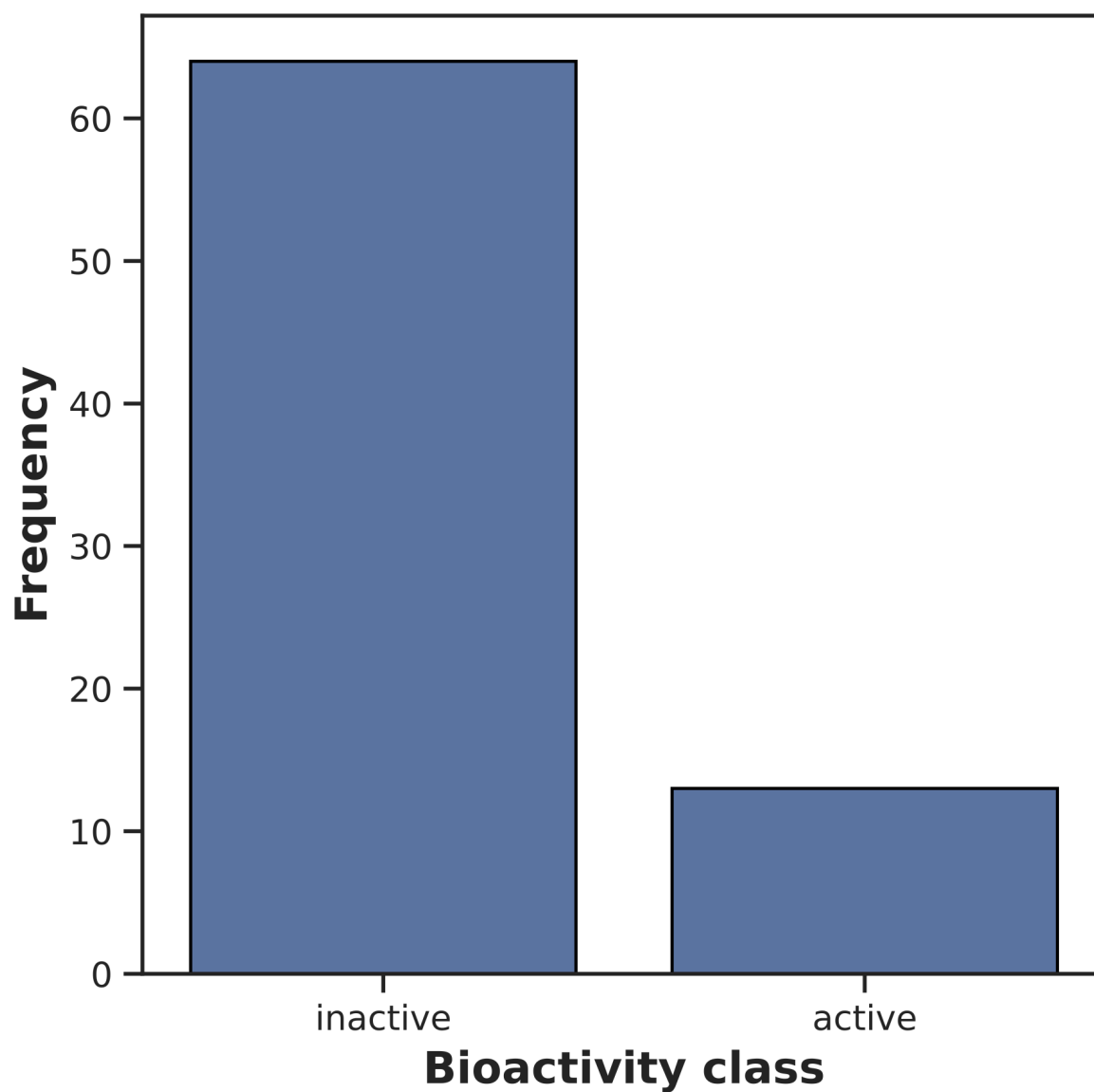


Figure 1a displays the quantity of inactive and active molecules within the dataset. 13 molecules are active at the target protein and 65 are inactive molecules. There was a total of 78 molecules.

Table 1a. Statistical Significance of pIC50

| | Descriptor | Statistics | p | alpha | Interpretation |
|---|------------|------------|----------------------|-------|------------------------------------|
| 0 | pIC50 | 832.0 | 1.58809447522824E-08 | 0.05 | Different distribution (reject H0) |

This table displays the result of the Mann-Whitney U test for pIC50. There was a statistically significant difference between the pIC50 value for active and inactive molecules and drugs. This can be used as a distinguisher in the machine learning regression model.

Table 1b. Statistical Significance for Number of Hydrogen Bond Donors

| | Descriptor | Statistics | p | alpha | Interpretation |
|---|------------|------------|-----------------------|-------|------------------------------------|
| 0 | NumHDonors | 178.5 | 0.0006202477871294170 | 0.05 | Different distribution (reject H0) |

This table displays the result of the Mann-Whitney U test for the number of hydrogen bond donors. There was a statistically significant difference between the number of hydrogen bond donors for active and inactive molecules and drugs. This can be used as a distinguisher in the machine learning regression model.

Table 1c. Statistical Significance for Number of Hydrogen Bond Acceptors

| | Descriptor | Statistics | p | alpha | Interpretation |
|---|----------------|------------|----------------------|-------|------------------------------------|
| 0 | NumHAacceptors | 248.0 | 0.020963972338247300 | 0.05 | Different distribution (reject H0) |

This table displays the result of the Mann-Whitney U test for the number of hydrogen bond acceptors. There was a statistically significant difference between the number of hydrogen bond

acceptors for active and inactive molecules and drugs. This can be used as a distinguisher in the machine learning regression model.

Table 1d. Statistical Significance of Molecular Weight

| | Descriptor | Statistics | p | alpha | Interpretation |
|---|------------|------------|----------------------|-------|------------------------------------|
| 0 | MW | 249.0 | 0.023566956483549300 | 0.05 | Different distribution (reject H0) |

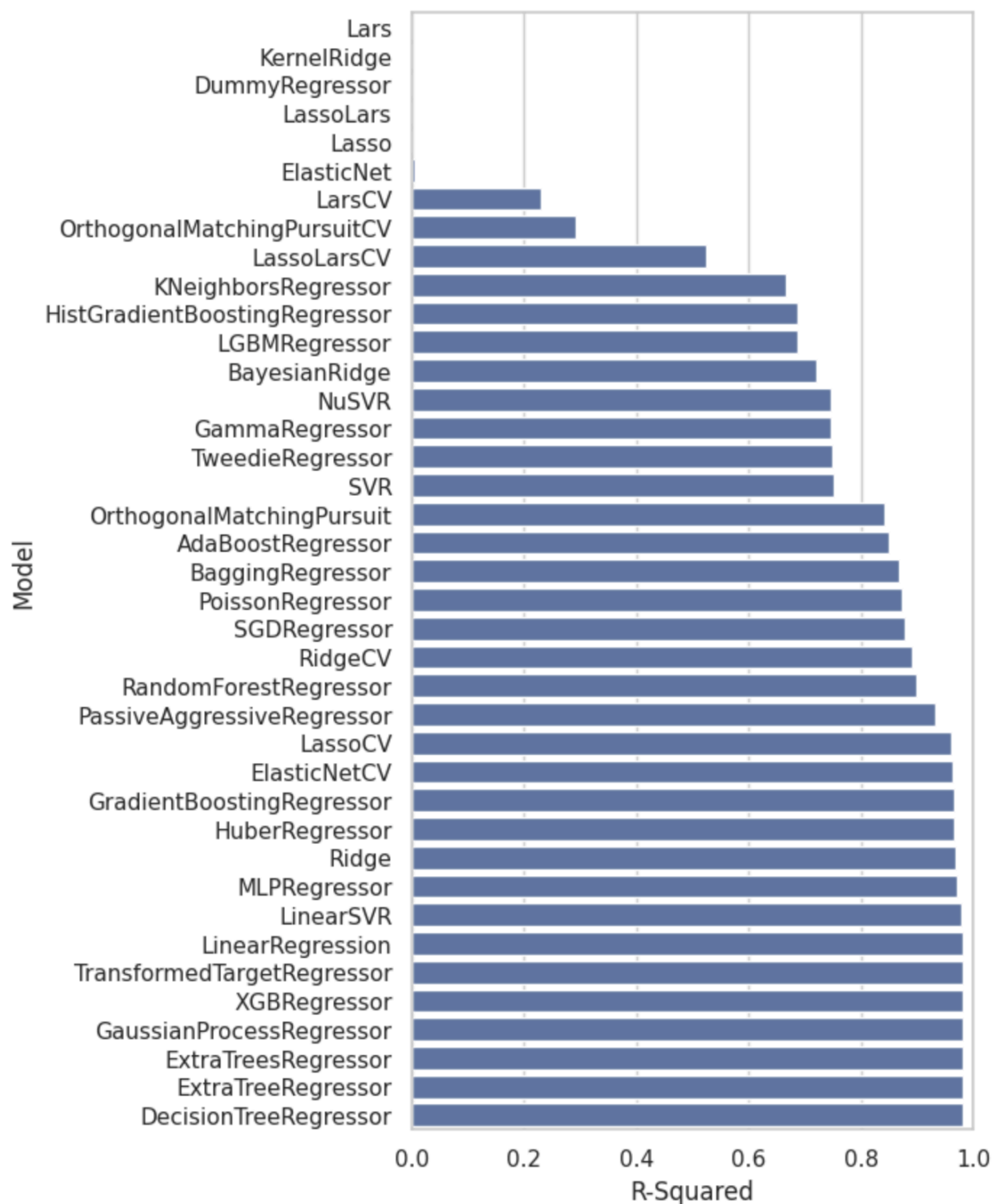
This table displays the result of the Mann-Whitney U test for molecular weight. There was a statistically significant difference between the molecular weight of active and inactive molecules and drugs. This can be used as a distinguisher in the machine learning regression model.

Table 1e. Statistical Significance of LogP

| | Descriptor | Statistics | p | alpha | Interpretation |
|---|------------|------------|--------------------|-------|---------------------------------------|
| 0 | LogP | 425.0 | 0.9079804356663140 | 0.05 | Same distribution (fail to reject H0) |

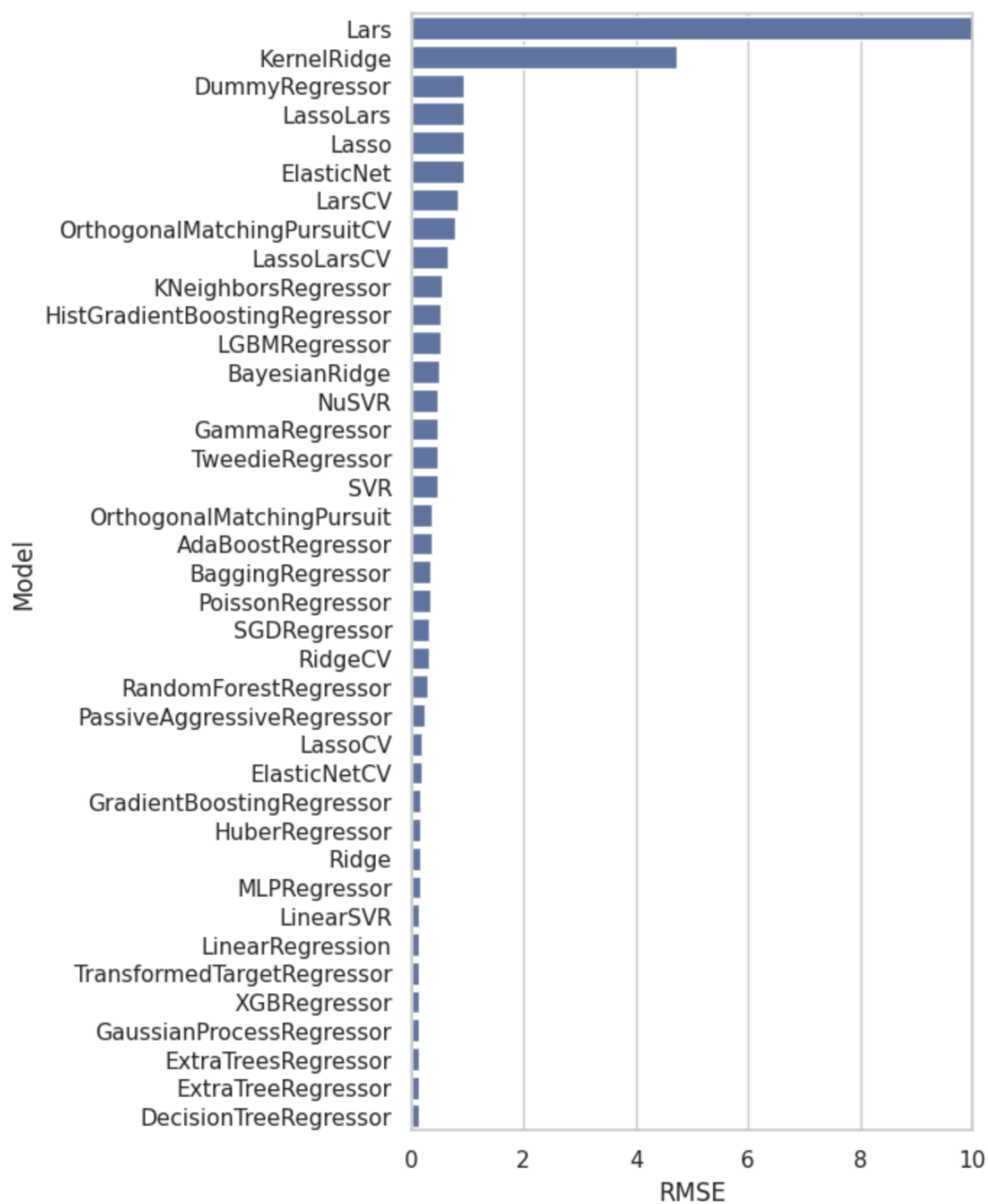
This table displays the result of the Mann-Whitney U test for druglikeness (LogP). There was not a statistically significant difference between the LogP of active and inactive molecules and drugs. This cannot be used as a distinguisher in the machine learning regression model.

Figure 2a. R-Squared values of the machine learning models during training



DecisionTreeRegressor performed the best on this data.

Figure 2b. RMSE values of the machine learning models during training

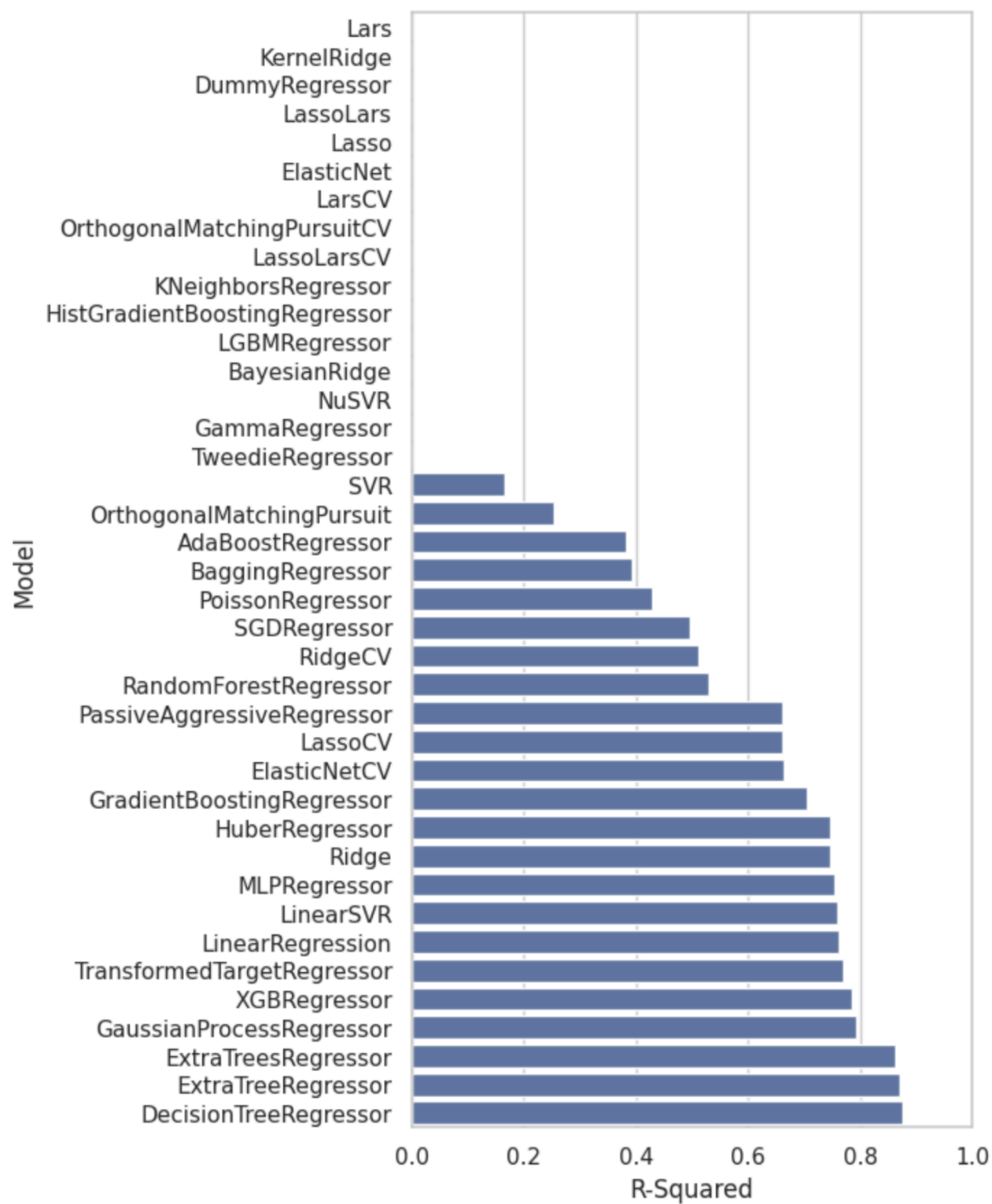


DecisionTreeRegressor performed the best on this data.

Table 2a. R-squared and RMSE values of the best model during training

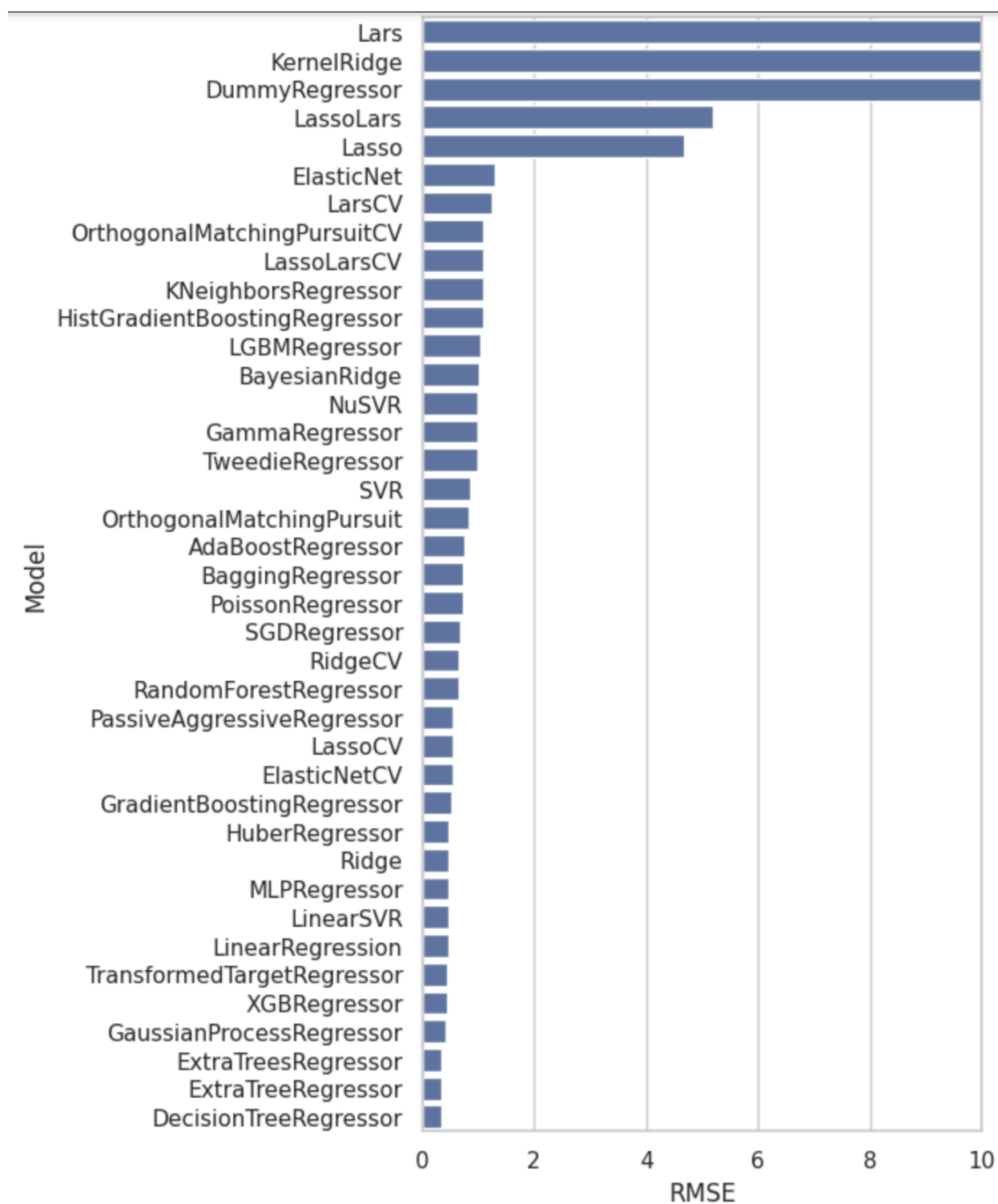
| Model | R-Squared Score | RMSE Score |
|-----------------------|-----------------|------------|
| DecisionTreeRegressor | 0.98 | 0.13 |

Figure 3a. R-Squared values of the machine learning models during testing



DecisionTreeRegressor performed the best on this data.

Figure 3b. RMSE values of the machine learning models during testing



DecisionTreeRegressor performed the best on this data.

Table 3a. R-squared and RMSE values of the best model during testing

| Model | R-Squared Score | RMSE Score |
|-----------------------|-----------------|------------|
| DecisionTreeRegressor | 0.9019 | 0.34 |

Discussion

As scientists face more difficulties in drug development whether it be during the analysis of properties or simply with time and cost of their research, this research aims to solve those problems by creating a regression-based machine learning model that utilizes bioactivity data and the molecule or drug's SMILES string to predict pIC50 and IC50 values. These pIC50 and IC50 values can be used to interpret a drug's potency and thus toxicity. This information is crucial in the H2L and lead discovery stages. This project will also save money and time in these stages of drug discovery.

Based on the results of the research, it can be said that models such as this regression-based machine learning model for SARS Coronavirus 3C-like proteinase are a plausible route to take in the journey of improving conventional drug discovery practices. This model was trained on chemical footprints determined by the descriptors. The statistical tests show that all descriptors except for druglikeness (LogP) contribute (they create a statistically significant difference as their p-values are less than 0.05) to whether or not a molecule will be active or inactive and thereby have an impact on pIC50 and IC50 values. Next, the dataset was split into training and testing sets to begin the creation, training, and testing of the machine learning models. Using LazyRegressor and the dataset, many models were created and analyzed by their RMSE and R-squared values. The best model after testing was determined to be the

DecisionTreeRegressor machine learning model. This model had the lowest RMSE value (0.34) and the highest R-squared value (0.9019). Having a large R-squared value is favored as that implies that the model fits the situation well and the predictions are accurate because the predictions are closer to the line of best fit (the actual values). RMSE is the average distance that all of the predictions are from the actual values. Having lower RMSE values is better because it would imply that the predictions are not very far off from the actual values. Together the low RMSE value and high R-squared value convey that the developed model is accurate and fast whilst also being cost-effective.

One limitation that greatly hindered the performance of this model was the amount of data available for the target protein of SARS Coronavirus 3C-like proteinase. Machine learning models and artificial intelligence models in general require immense amounts of data to go off to create better predictions and learn properly during training. This is especially true if there is significant variance within the data. The database utilized in this project had 78 data points, which is far too few to reach accuracies above 95 to 97 percent. Hence, this model did not reach higher accuracies. Overall, at this point, since the developed model has not reached an accuracy of higher than 95 percent, the models developed in this project can be used as more of an aid and supplement to scientists to make preliminary judgments and to confirm that the data they receive is accurate.

Implications and Applications

An important idea to note is that the regression model developed within this project is not only limited to molecules for the SARS Coronavirus 3C-like proteinase target protein. These models are very versatile in their uses. The same process can be applied to any target protein and

drug data. The only change that would be necessary is changing the dataset. By doing this, high-performance models for all drugs and target proteins can be developed.

Future Research

In the future, similar regression-based machine learning models can be expanded to other fields and even areas of drug discovery rather than focusing on the target protein SARS Coronavirus 3C-like proteinase. Work can also be done towards optimizing the current model or even testing new models with the same situation to look for ways to improve accuracy and decrease prediction error. The main action that can be taken to improve the accuracy of these types of models is gathering more data. Machine learning algorithms rely heavily on training and testing data. Without ample data, it is not possible to create models that reach accuracies of 97 percent higher. Finally, to improve the usability of these models, web applications can be formed to streamline the process of entering data and outputting a prediction.

Conclusion

In the world we live in today, the field of drug discovery is riddled with difficulties whether it be with the property analysis of drugs or the immense cost of clinical trials. The ultimate goal of this project was to develop and test regression-based machine learning models that could make predictions about drugs and molecules designed for SARS Coronavirus 3C-like proteinase while also being accurate, fast, and cost-effective. The models would predict the potency and toxicity by outputting an IC₅₀ value and pIC₅₀ value for the molecule. These models made predictions based on SMILES notation and bioactivity data of the molecules and drugs from the database ChEMBL. This data was first processed to remove any unnecessary data

and to ensure it was readable by the machine learning models. Then, the data was split into training and testing data and imported into LazyRegressor to start training the various types of regression models. In the end, the best regression-based machine learning model, DecisionTreeRegressor achieved an RMSE value of 0.34 and a R-squared value of 0.9019. This shows that the model was fairly accurate while also being cost-effective and quick. The results could have been improved with the inclusion of more data about SARS Coronavirus 3C-like proteinase. Despite the lack of data and average accuracy of less than 95 percent, we deem this project to be a success. This is because although it did not reach a very high accuracy it was still fairly accurate considering the constraints. In addition, this project proved that using these regression-based machine learning models is a possible avenue in improving the field of drug discovery. At the very least, the model developed in this project can be used as more of an aid and supplement to scientists to make preliminary judgments and to confirm that the data they receive is accurate. By continuing to build on the work conducted in this project we can continue to improve the field of drug discovery and save millions of lives.

Appendix

Appendix 1

Preprocessed dataset used for machine learning regression model training and testing

| molecule_ch embl_id | canonical_smiles | class | MW | LogP | NumHD onors | NumHAcc eptors | pIC50 |
|--------------------------|---|--------------|------------------------|------------------------|----------------|-------------------|------------------------|
| CHEMBL187 579 | <chem>Cc1noc(C)c1CN1C(=O)C(=O)c2cc(C#N)ccc21</chem> | intermediate | 281.271 | 1.89262 | 0.0 | 5.0 | 5.142667503 568730 |
| CHEMBL188 487 | <chem>O=C1C(=O)N(Cc2ccc(F)cc2Cl)c2ccc(I)cc21</chem> | intermediate | 415.589 | 3.8132000000 00000 | 0.0 | 2.0 | 5.026872146 400300 |
| CHEMBL185 698 | <chem>O=C1C(=O)N(CC2COc3ccccc3O2)c2ccc(I)cc21</chem> | inactive | 421.1900000 0000000 | 2.6605000000 000000 | 0.0 | 4.0 | 4.869666231 5049900 |
| CHEMBL426 082 | <chem>O=C1C(=O)N(Cc2cc3ccccc3s2)c2ccccc21</chem> | inactive | 293.3470000 0000000 | 3.6308000000 000000 | 0.0 | 3.0 | 4.882397308 3099200 |
| CHEMBL187 717 | <chem>O=C1C(=O)N(Cc2cc3ccccc3s2)c2c1ccc2[N+](=O)[O-]</chem> | intermediate | 338.3440000 000000 | 3.5390000000 000000 | 0.0 | 5.0 | 5.698970004 336020 |

| | | | | | | | |
|--------------------------|--|--------------|------------------------|--------------------------|-----|-----|------------------------|
| CHEMBL365 134 | <chem>O=C1C(=O)N(Cc2cc3ccccc3s2)c2c(Br)ccc21</chem> | active | 372.2430000 0000000 | 4.3933000000 00000 | 0.0 | 3.0 | 6.008773924 307510 |
| CHEMBL187 598 | <chem>O=C1C(=O)N(Cc2cc3ccccc3s2)c2ccc(F)cc21</chem> | intermediate | 311.3370000 0000000 | 3.7699000000 000000 | 0.0 | 3.0 | 5.316952961 76115 |
| CHEMBL190 743 | <chem>O=C1C(=O)N(Cc2cc3ccccc3s2)c2ccc(I)cc21</chem> | active | 419.2430000 0000000 | 4.2354000000 00000 | 0.0 | 3.0 | 6.022276394 711150 |
| CHEMBL365 469 | <chem>O=C1C(=O)N(Cc2cc3ccccc3s2)c2ccc(Cl)cc21</chem> | inactive | 327.7920000 000000 | 4.2842000000 00000 | 0.0 | 3.0 | 4.950781977 329820 |
| CHEMBL188 983 | <chem>O=C1C(=O)N(C/C=C/c2cc3ccccc3s2)c2ccc(I)cc21</chem> | inactive | 445.2810000 0000000 | 4.7486000000 00000 | 0.0 | 3.0 | 4.628932137 728260 |
| CHEMBL191 575 | <chem>O=C(Nc1ccc(Cl)cc1)c1ccc(CN2C(=O)C(=O)c3ccc(I)ccc32)s1</chem> | inactive | 522.7510000 000000 | 4.9879000000 00000 | 1.0 | 4.0 | 4.900664722 314040 |
| CHEMBL370 923 | <chem>O=C1C(=O)N(Cc2ccc(C(=O)N3CCCCC3)s2)c2ccc(I)cc21</chem> | inactive | 480.327 | 3.7083000000 00000 | 0.0 | 4.0 | 4.756961951 313710 |
| CHEMBL194 398 | <chem>CCOC(=O)/C=C/[C@H](C[C@@H]1CCNC1=O)NC(=O)[C@@H](CC(=O)[C@@H](NC(=O)c1cc(C)on1)C(C)C)Cc1ccccc1</chem> | inactive | 580.6820000 000000 | 2.6858200000 000000 | 3.0 | 8.0 | 4.346787486 224660 |
| CHEMBL196 635 | <chem>CCOC(=O)/C=C/[C@H](C[C@@H]1CCNC1=O)NC(=O)[C@@H](CC=C(C)C)CC(=O)[C@@H](NC(=O)c1cc(C)on1)C(C)C</chem> | inactive | 558.6760000 000000 | 2.7994200000 000000 | 3.0 | 8.0 | 4.154901959 985740 |
| CHEMBL209 287 | <chem>CCCCN1C(=O)C(=O)c2cc(I)ccc21</chem> | inactive | 329.137 | 2.6206000000 000000 | 0.0 | 2.0 | 4.180456064 458130 |
| CHEMBL358 279 | <chem>NC(=O)c1ccc2c(c1)C(=O)C(=O)N2Cc1cc2ccccc2c1</chem> | active | 330.3430000 0000000 | 2.6682000000 000000 | 1.0 | 3.0 | 6.431798275 933010 |
| CHEMBL348 660 | <chem>NC(=O)c1ccc2c(c1)C(=O)C(=O)N2Cc1ccccc1</chem> | inactive | 280.2830000 0000000 | 1.5150000000 000000 | 1.0 | 3.0 | 4.903089986 991940 |
| CHEMBL379 727 | <chem>CCCCN1C(=O)C(=O)c2cc(C(N)=O)ccc21</chem> | inactive | 246.2660000 0000000 | 1.1149000000 000000 | 1.0 | 3.0 | 4.721246399 047170 |
| CHEMBL210 525 | <chem>CCCN1C(=O)C(=O)c2cc(C(N)=O)ccc21</chem> | inactive | 232.2390000 0000000 | 0.7248 | 1.0 | 3.0 | 4.602059991 3279600 |
| CHEMBL148 483 | <chem>CN1C(=O)C(=O)c2cc(C(N)=O)ccc21</chem> | inactive | 204.1850000 0000000 | -0.0554000000 0000030 | 1.0 | 3.0 | 4.148741651 2809200 |
| CHEMBL383 725 | <chem>O=C1C(=O)N(Cc2ccc3ccccc3c2)c2ccc(I)cc21</chem> | intermediate | 413.2140000 0000000 | 4.1739000000 00000 | 0.0 | 2.0 | 5.958607314 841780 |
| CHEMBL118 596 | <chem>O=C1C(=O)N(Cc2ccccc2)c2ccc(I)cc21</chem> | inactive | 363.1540000 0000000 | 3.0207000000 000000 | 0.0 | 2.0 | 4.301029995 663980 |
| CHEMBL208 732 | <chem>O=C(CSc1nccc(-c2csc(-c3ccccc3)n2)n1)Nc1cc(Cl)cc(Cl)c1</chem> | intermediate | 473.4100000 000000 | 6.3047000000 00000 | 1.0 | 6.0 | 5.520000247 4083600 |
| CHEMBL210 146 | <chem>COc1cccc(-c2nc(SCC(=O)Nc3ccc([N+](=O)[O-])cc3)nc(O)c2C#N)c1</chem> | inactive | 437.4370000 000000 | 3.3684800000 000000 | 2.0 | 9.0 | 3.520000002 9339800 |
| CHEMBL207 458 | <chem>CCOC(=O)c1cnc(SCC(=O)Nc2ccc([N+](=O)[O-])cc2)nc1N</chem> | inactive | 377.3820000 000000 | 1.8745000000 000000 | 2.0 | 9.0 | 3.600000005 4478700 |
| CHEMBL207 484 | <chem>COC(OC)c1cc(O)nc(SCC(=O)Nc2ccc(C(F)(F)F)cc2)n1</chem> | inactive | 403.3820000 0000000 | 3.2232000000 00000 | 2.0 | 7.0 | 3.690000009 501580 |
| CHEMBL207 207 | <chem>Cc1cc(O)nc(SCC(=O)Nc2cc(Cl)ccc2Oc2ccccc2)n1</chem> | inactive | 401.8750000 000000 | 4.6671200000 00000 | 2.0 | 6.0 | 4.0 |
| CHEMBL210 487 | <chem>CCC(Sc1nc(O)c(C#N)c(-c2cccc(OC)c2)n1)C(=O)Nc1ccc(C(C)=O)cc1</chem> | inactive | 462.5310000 0000000 | 4.4414800000 00000 | 2.0 | 8.0 | 4.219999989 963100 |

| | | | | | | | |
|--------------------------|--|--------------|------------------------|------------------------|-----|-----|------------------------|
| CHEMBL380 470 | <chem>COc1ccc(NC(=O)CSc2nc(O)cc(-c3ccccc3)n2)cc1OC</chem> | inactive | 397.4560000 0000000 | 3.5972000000 000000 | 2.0 | 7.0 | 4.339999990 132750 |
| CHEMBL210 612 | <chem>CC(C)c1ccc(NC(=O)CSc2nccc(-c3ccccc3)n2)cc1</chem> | inactive | 369.5150000 0000000 | 5.0593000000 00000 | 1.0 | 5.0 | 4.389999976 337710 |
| CHEMBL209 667 | <chem>O=C(CSc1ncccc(-c2cc(-c3ccccc3Cl)no2)n1)Nc1ccc(Cl)cc1</chem> | inactive | 457.3420000 000000 | 5.8362000000 000000 | 1.0 | 6.0 | 4.820000071 285180 |
| CHEMBL210 097 | <chem>O=C(CSc1ncccc(-c2cc(-c3ccc(Cl)cc3Cl)no2)n1)Nc1ccc(C(F)(F)F)cc1</chem> | inactive | 525.3390000 000000 | 6.8550000000 00000 | 1.0 | 6.0 | 4.820000071 285180 |
| CHEMBL378 674 | <chem>CSc1sc(-c2nc(C)cs2)c(C)c1-c1ccnc(SCC(=O)Nc2ccccc2Cl)n1</chem> | inactive | 519.1420000 000000 | 7.0515400000 00010 | 1.0 | 8.0 | 4.920000156 997060 |
| CHEMBL210 216 | <chem>Cc1nc(SCC(=O)Nc2ccc([N+](=O)[O-])cc2)nc(O)c1C</chem> | inactive | 334.3570000 000000 | 2.43804 | 2.0 | 7.0 | 3.0 |
| CHEMBL210 195 | <chem>Cc1cc(C(F)(F)F)nc(SCC(=O)Nc2ccc(Cl)cc2F)n1</chem> | inactive | 379.766 | 4.3271200000 00000 | 1.0 | 4.0 | 3.300000003 143150 |
| CHEMBL210 437 | <chem>COc1ccc(-c2ccnc(SCC(=O)Nc3ccc(C(C)C)cc3)n2)cc1</chem> | inactive | 393.5120000 0000000 | 5.0064000000 00000 | 1.0 | 5.0 | 3.389999997 659040 |
| CHEMBL378 677 | <chem>CC(C)(C)c1ccc(NC(=O)CSc2nc(O)c(C#N)c(C3CCCC3)n2)cc1</chem> | inactive | 424.5700000 000000 | 5.1298800000 00000 | 2.0 | 6.0 | 3.449999999 06189 |
| CHEMBL210 972 | <chem>CC(Sc1nc(O)c(C#N)c(-c2ccccc2)n1)C(=O)Nc1ccc(Cl)cc1</chem> | inactive | 410.8860000 000000 | 4.4935800000 00000 | 2.0 | 6.0 | 3.520000002 9339800 |
| CHEMBL210 145 | <chem>CCOc1ccc(N2C(=O)CC(Sc3nc(C)cc(C)n3)C2=O)cc1</chem> | inactive | 357.4350000 0000000 | 2.9163400000 000000 | 0.0 | 6.0 | 3.520000002 9339800 |
| CHEMBL377 225 | <chem>N#Cc1c(O)nc(SCC(=O)Nc2ccc(C(F)(F)F)c2)nc1-c1ccccc1</chem> | inactive | 430.4110000 0000000 | 4.4704800000 00000 | 2.0 | 6.0 | 3.690000009 501580 |
| CHEMBL210 823 | <chem>CCCCc1ccc(NC(=O)CSc2nc(O)c(C#N)c(C3CCCC3)n2)c(C)c1</chem> | inactive | 438.5970000 000000 | 5.4834000000 00010 | 2.0 | 6.0 | 3.690000009 501580 |
| CHEMBL207 381 | <chem>CC(C)c1ccc(NC(=O)CSc2nccc(-c3csc(COc4cccc4Cl)n3)n2)cc1</chem> | inactive | 511.0720000 0000000 | 6.6867000000 00010 | 1.0 | 7.0 | 3.690000009 501580 |
| CHEMBL210 511 | <chem>Cc1cc(O)nc(SCC(=O)Nc2ccc(Oc3ccc(Cl)cc3)cc2)n1</chem> | inactive | 401.8750000 0000000 | 4.6671200000 00000 | 2.0 | 6.0 | 3.690000009 501580 |
| CHEMBL210 632 | <chem>COc1cccc(-c2nc(SCC(=O)Nc3ccc(C(C)=O)cc3)nc(O)c2C#N)c1</chem> | inactive | 434.4770000 000000 | 3.6628800000 00000 | 2.0 | 8.0 | 4.219999989 963100 |
| CHEMBL210 497 | <chem>COc1cccc(-c2nc(SCC(=O)Nc3ccc(S(N)(=O)=O)cc3)nc(O)c2C#N)c1</chem> | inactive | 471.5200000 000000 | 2.10768 | 3.0 | 9.0 | 4.389999976 337710 |
| CHEMBL208 584 | <chem>CCCc1cc(O)nc(SCC(=O)Nc2ccc(Cl)cc2)n1</chem> | inactive | 337.8320000 0000000 | 3.5189000000 00000 | 2.0 | 5.0 | 4.519999959 791460 |
| CHEMBL208 763 | <chem>O=C(CSc1ncccc(-c2cc(-c3ccccc3)no2)n1)Nc1ccc(Cl)cc1</chem> | inactive | 422.8970000 0000000 | 5.1828000000 00000 | 1.0 | 6.0 | 4.820000071 285180 |
| CHEMBL209 227 | <chem>Cc1nc(-c2nc(-c3ccnc(SCC(=O)Nc4ccc(Cl)cc4)n3)cs2)cs1</chem> | inactive | 460.0090000 0000000 | 5.4162200000 00000 | 1.0 | 8.0 | 4.849999859 991130 |
| CHEMBL210 092 | <chem>CSc1sc(-c2nc(C)cs2)c(C)c1-c1ccnc(SCC(=O)Nc2ccc(Cl)cc2)n1</chem> | inactive | 519.1420000 000000 | 7.0515400000 00010 | 1.0 | 8.0 | 4.950000175 844590 |
| CHEMBL377 150 | <chem>Cn1nc(-c2ccc(-c3ccnc(SCC(=O)Nc4ccc(Cl)cc4)n3)s2)cc1C(F)(F)F</chem> | inactive | 509.966 | 6.0086000000 00000 | 1.0 | 7.0 | 5.0 |
| CHEMBL212 454 | <chem>O=C(Oc1ccc(S(=O)(=O)c2ccc(OC(=O)C(Cl)=C(Cl)Cl)cc2)cc1)C(Cl)=C(Cl)Cl</chem> | active | 565.041 | 6.1012000000 00000 | 0.0 | 6.0 | 6.045757490 560680 |
| CHEMBL213 581 | <chem>Nc1ncc(S(=O)(=O)c2ccc(Cl)cc2)c(N)n1</chem> | intermediate | 284.728 | 1.1272000000 000000 | 2.0 | 6.0 | 5.221848749 616360 |

| | | | | | | | |
|--------------------------|---|--------------|------------------------|------------------------|-----|------|------------------------|
| CHEMBL380 403 | <chem>O=[N+](O)c1ccc(C(F)(F)F)ccc1S(=O)(=O)c1ccc(Cl)cc1</chem> | inactive | 365.716 | 4.0998000000 00000 | 0.0 | 4.0 | 4.920818753 952380 |
| CHEMBL212 504 | <chem>Cc1nc(S(=O)(=O)c2ccccc2)c(C#N)c(C)c1[N+](=O)[O-]</chem> | inactive | 317.3260000 000000 | 2.3111200000 000000 | 0.0 | 6.0 | 4.886056647 693160 |
| CHEMBL211 969 | <chem>Cc1ccc(S(=O)(=O)c2nc(C)c([N+](=O)[O-])c(C)c2C#N)cc1</chem> | inactive | 331.3530000 0000000 | 2.6195400000 000000 | 0.0 | 6.0 | 4.886056647 693160 |
| CHEMBL384 739 | <chem>O=[N+](O)c1ccc(S(=O)(=O)c2ccc(Cl)c2)[n+](O)c1</chem> | inactive | 314.706 | 1.7144000000 000000 | 0.0 | 5.0 | 4.823908740 944320 |
| CHEMBL215 732 | <chem>CCOC(=O)/C(C#N)=C/Nc1ccc(S(=O)(=O)c2ccc(N=C/C(C#N)=C(O)OCC)cc2)c1</chem> | inactive | 494.5290000 000000 | 3.9338600000 00000 | 2.0 | 10.0 | 4.795880017 344080 |
| CHEMBL212 240 | <chem>O=C(O)c1ccc(S(=O)(=O)c2cc(Br)c(O)c(Br)c2)cc1</chem> | inactive | 436.0770000 0000000 | 3.4482000000 00000 | 2.0 | 4.0 | 4.795880017 344080 |
| CHEMBL377 253 | <chem>CC(=O)c1cccc1S(=O)(=O)c1cccc1C(=O)O</chem> | inactive | 304.3230000 0000000 | 2.4202000000 000000 | 1.0 | 4.0 | 4.795880017 344080 |
| CHEMBL215 397 | <chem>O=[N+](O)c1ccc(S(=O)(=O)c2ccc([N+](=O)[O-])cc2)cc1</chem> | inactive | 308.271 | 2.3358 | 0.0 | 6.0 | 4.602059991 3279600 |
| CHEMBL378 342 | <chem>CCOC(=O)C(=CNc1ccc(S(=O)(=O)c2cc(N=C/C(C(=O)OCC)=C(O)OCC)cc2)c1)C(=O)OCC</chem> | inactive | 588.6350000 000000 | 4.0129000000 00000 | 2.0 | 12.0 | 4.494850021 680090 |
| CHEMBL379 642 | <chem>O=C(Sc1nnc(C(F)(F)F)[nH]1)c1ccc(C#Cc2ccccc2)o1</chem> | intermediate | 363.3200000 000000 | 3.7488000000 00000 | 1.0 | 5.0 | 5.522878745 280340 |
| CHEMBL427 404 | <chem>Cc1noc(NC(=O)c2ccc(-c3cc(C(F)(F)F)n3C)s2)c1[N+](=O)[O-]</chem> | intermediate | 401.3260000 000000 | 3.6243200000 000000 | 1.0 | 8.0 | 5.301029995 663980 |
| CHEMBL212 190 | <chem>Cc1oc(C(C)(C)C)cc1-c1cc(NS(=O)(=O)c2cccs2)[nH]n1</chem> | inactive | 365.48 | 4.1379200000 00000 | 2.0 | 5.0 | 5.0 |
| CHEMBL378 700 | <chem>CSc1[nH]nc(NC(=O)c2cccs2)c1S(=O)(=O)c1ccccc1</chem> | inactive | 379.4880000 000000 | 3.2782000000 00000 | 2.0 | 6.0 | 4.823908740 944320 |
| CHEMBL212 019 | <chem>CC1(C)CC(=O)c2c(NCc3ccco3)sc(C#N)c2C1</chem> | inactive | 300.3830000 000000 | 3.9799800000 00000 | 1.0 | 5.0 | 4.795880017 344080 |
| CHEMBL212 399 | <chem>CSc1nn(-c2c([N+](=O)[O-])c(C)nn2C)c(-c2cccs2)c1C#N</chem> | inactive | 360.4240000 0000000 | 3.1445000000 000000 | 0.0 | 9.0 | 4.744727494 8966900 |
| CHEMBL215 733 | <chem>O=S(=O)(Cc1[nH]c(-c2ccc(Cl)s2)c[s+1]c1cccs1.[Br-])</chem> | inactive | 442.8340000 0000000 | 1.7786 | 1.0 | 4.0 | 4.744727494 8966900 |
| CHEMBL375 130 | <chem>Cc1nn(C)c(NCc2ccc(-c3cccs3)s2)c1[N+](=O)[O-]</chem> | inactive | 334.4260000 000000 | 4.0388200000 00000 | 1.0 | 7.0 | 4.698970004 336020 |
| CHEMBL214 372 | <chem>O=C(Cc1nccs1)c1nccs1</chem> | inactive | 210.2830000 0000000 | 2.025 | 0.0 | 5.0 | 4.397940008 6720400 |
| CHEMBL212 218 | <chem>Cc1cc(S(=O)(=O)c2c([N+](=O)[O-])cc(C(F)(F)F)cc2[N+](=O)[O-])c(Cl)cc1Cl</chem> | active | 459.1850000 000000 | 4.9698200000 00000 | 0.0 | 6.0 | 6.522878745 280340 |
| CHEMBL222 840 | <chem>O=C(Oc1cncc(Cl)c1)c1ccco1</chem> | active | 223.615 | 2.5472 | 0.0 | 4.0 | 7.221848749 616360 |
| CHEMBL222 769 | <chem>O=C(Oc1cncc(Cl)c1)c1ccc(-c2ccc(Cl)cc2)o1</chem> | active | 334.158 | 4.8676000000 00000 | 0.0 | 4.0 | 7.200659450 546420 |
| CHEMBL222 735 | <chem>COc1cccc(C(=O)Oc2cncc(Cl)c2)c1</chem> | active | 263.68 | 2.9628000000 000000 | 0.0 | 4.0 | 6.468521082 957750 |
| CHEMBL222 628 | <chem>O=C(Oc1cncc(Cl)c1)c1cscn1</chem> | active | 240.6710000 0000000 | 2.4107000000 000000 | 0.0 | 5.0 | 6.568636235 841010 |

| | | | | | | | |
|---------------------------|---|----------|------------------------|------------------------|-----|-----|------------------------|
| CHEMBL222 893 | <chem>O=C(Oc1cncc(Cl)c1)c1cc2ccccc2s1</chem> | active | 289.7430000 0000000 | 4.1689000000 00000 | 0.0 | 4.0 | 7.022276394 711150 |
| CHEMBL225 515 | <chem>O=C(Oc1cncc(Cl)c1)c1cc2ccccc2[nH]1</chem> | active | 272.6910000 0000000 | 3.4355000000 00000 | 1.0 | 3.0 | 7.187086643 357140 |
| CHEMBL222 234 | <chem>O=C(Oc1cncc(Br)c1)c1ccco1</chem> | active | 268.066 | 2.6563000000 000000 | 0.0 | 4.0 | 7.301029995 663980 |
| CHEMBL426 898 | <chem>O=C(Oc1cncc(Cl)c1)c1cc2ccccc2o1</chem> | active | 273.6750000 0000000 | 3.7004000000 00000 | 0.0 | 4.0 | 6.769551078 621730 |
| CHEMBL187 266 | <chem>Cc1coc2c1C(=O)C(=O)c1c-2ccc2c1CC CC2(C)C</chem> | inactive | 294.35 | 4.2479200000 00000 | 0.0 | 3.0 | 4.050122295 963130 |
| CHEMBL215 254 | <chem>Cc1coc2c1C(=O)C(=O)c1c-2ccc2c1CC CC2(C)CO</chem> | inactive | 310.3490000 0000000 | 3.2203200000 000000 | 1.0 | 4.0 | 4.605548319 173780 |
| CHEMBL214 6517 | <chem>COC(=O)[C@@]1(C)CCCc2c1ccc1c2C(=O) C(=O)c2c(C)coc2-1</chem> | inactive | 338.3590000 000000 | 3.4010200000 000000 | 0.0 | 5.0 | 4.675717544 702310 |
| CHEMBL187 460 | <chem>C[C@H]1COC2=C1C(=O)C(=O)c1c2ccc 2c1CCCC2(C)C</chem> | inactive | 296.366 | 3.4433000000 000000 | 0.0 | 3.0 | 3.644548479 8734800 |
| CHEMBL363 535 | <chem>Cc1coc2c1C(=O)C(=O)c1c-2ccc2c(C)cc cc12</chem> | inactive | 276.291 | 4.0956400000 00000 | 0.0 | 3.0 | 4.412289034 981090 |
| CHEMBL227 075 | <chem>Cc1cccc2c3c(ccc12)C1=C(C(=O)C3=O)[C@@H](C)CO1</chem> | inactive | 278.307 | 3.2910200000 000000 | 0.0 | 3.0 | 4.841637507 904750 |
| CHEMBL458 30 | <chem>CC(C)C1=Cc2ccc3c(c2C(=O)C1=O)CC CC3(C)C</chem> | inactive | 282.3830000 000000 | 4.1053000000 00000 | 0.0 | 2.0 | 4.675717544 702310 |

Appendix 2

Code for Lipinski descriptor calculation

```
def lipinski(smiles, verbose=False):

    moldata= []
    for elem in smiles:
        mol=Chem.MolFromSmiles(elem)
        moldata.append(mol)
    baseData= np.arange(1,1)
    i=0
    for mol in moldata:
        desc_MolWt = Descriptors.MolWt(mol)
        desc_MolLogP = Descriptors.MolLogP(mol)
        desc_NumHDonors = Lipinski.NumHDonors(mol)
        desc_NumHAacceptors = Lipinski.NumHAacceptors(mol)
        row = np.array([desc_MolWt,
            desc_MolLogP,
            desc_NumHDonors,
            desc_NumHAacceptors])
        if(i==0):
```

```
baseData=row
else:
baseData=np.vstack([baseData, row])
i=i+1
columnNames=["MW", "LogP", "NumHDonors", "NumHAcceptors"]
descriptors = pd.DataFrame(data=baseData, columns=columnNames)
return descriptors
```

Appendix 3

Mann-Whitney U Test for statistical significance code. For analysis use the function

“mannwhitney(“”).” For example for pIC50 the function would be mannwhitney(“pIC50”).

```
def mannwhitney(descriptor, verbose=False):

from numpy.random import seed
from numpy.random import randn
from scipy.stats import mannwhitneyu

# seed the random number generator
seed(1)

# actives and inactives
selection = [descriptor, 'bioactivity_class']
df = df_2class[selection]
active = df[df.bioactivity_class == 'active']
active = active[descriptor]

selection = [descriptor, 'bioactivity_class']
df = df_2class[selection]
inactive = df[df.bioactivity_class == 'inactive']
inactive = inactive[descriptor]

# compare samples
stat, p = mannwhitneyu(active, inactive)
#print('Statistics=%.3f, p=%.3f' % (stat, p))

# interpret
alpha = 0.05
if p > alpha:
interpretation = 'Same distribution (fail to reject H0)'
```

```

else:
interpretation = 'Different distribution (reject H0)'
results = pd.DataFrame({'Descriptor':descriptor,
'Statistics':stat,
'p':p,
'alpha':alpha,
'Interpretation':interpretation}, index=[0])
filename = 'mannwhitneyu_' + descriptor + '.csv'
results.to_csv(filename)

return results

```

References

- Alexandre Blanco-González, Cabezón, A., Seco-González, A., Conde-Torres, D., Antelo-Riveiro, P., Ángel Piñeiro, & Garcia-Fandino, R. (2023). The Role of AI in Drug Discovery: Challenges, Opportunities, and Strategies. *The Role of AI in Drug Discovery: Challenges, Opportunities, and Strategies*, 16(6), 891–891. <https://doi.org/10.3390/ph16060891>
- Alvarellos, M. (2023, October 5). *What are the current challenges of drug discovery?* Wwww.lifebit.ai. <https://www.lifebit.ai/blog/current-challenges-of-drug-discovery>
- Aykul, S., & Martinez-Hackert, E. (2016). Determination of half-maximal inhibitory concentration using biosensor-based protein interaction analysis. *Analytical Biochemistry*, 508(1), 97–103. <https://doi.org/10.1016/j.ab.2016.06.025>
- Berrouet, C., Dorilas, N., Rejniak, K. A., & Tuncer, N. (2020). Comparison of Drug Inhibitory Effects (IC50) in Monolayer and Spheroid Cultures. *Bulletin of Mathematical Biology*, 82(6). <https://doi.org/10.1007/s11538-020-00746-7>
- ChEMBL. (n.d.). *ChEMBL Database*. Wwww.ebi.ac.uk. <https://www.ebi.ac.uk/chembl/>
- Code Ocean. (n.d.). *Code Ocean*. Codeocean.com. Retrieved February 1, 2024, from <https://codeocean.com/explore/capsules?query=tag:data-curation>

IBM. (2023). *What is Machine Learning?* IBM.

<https://www.ibm.com/topics/machine-learning>

isirv Antiviral group. (n.d.). *Analysis of IC50 data*. Wwww.isirv.org. Retrieved January 31, 2024, from

<https://www.isirv.org/site/index.php/methodology/analysis-of-ic50-data#:~:text=IC50%20values%20can%20be>

Kurama, V. (2019, September 4). *Regression in Machine Learning: What it is and Examples of Different Models*. Built In.

<https://builtin.com/data-science/regression-machine-learning>

Laerd Statistics . (2013). *Mann-Whitney U Test in SPSS Statistics*. Laerd Statistics.

<https://statistics.laerd.com/spss-tutorials/mann-whitney-u-test-using-spss-statistics.php>

Liu, P., Li, H., Li, S., & Leung, K.-S. (2019). Improving prediction of phenotypic drug response on cancer cell lines using deep convolutional network. *BMC Bioinformatics*, 20(1). <https://doi.org/10.1186/s12859-019-2910-6>

Navre, M. (2019, November 21). *Why using pIC50 instead of IC50 will change your life*. Wwww.collaborativedrug.com.

<https://www.collaborativedrug.com/cdd-blog/why-using-pic50-instead-of-ic50-will-change-your-life#:~:text=A%20closely%20related%20concept%20is>

October 25, C. H., & 2015. (2015, October 15). *Delays in Drug Approval Are Deadly, Highlighting the Need for Improved Regulatory Efficiency - The ASCO Post*. Ascopost.com.

<https://ascopost.com/issues/october-25-2015/delays-in-drug-approval-are-deadly->

highlighting-the-need-for-improved-regulatory-efficiency/#:~:text=Worldwide%2C%201%20life%20year%20was

Patil, R. (2021, May 6). *Lazy Predict - Best Suitable Model for You*. Analytics Vidhya.

<https://www.analyticsvidhya.com/blog/2021/05/lazy-predict-best-suitable-model-for-you/>

Seldon. (2021, October 16). *Supervised vs Unsupervised Learning Explained*. Seldon.

<https://www.seldon.io/supervised-vs-unsupervised-learning-explained>

Sharma, A. (2020, March 22). *A Comprehensive Guide to Google Colab*. Analytics Vidhya.

<https://www.analyticsvidhya.com/blog/2020/03/google-colab-machine-learning-deep-learning/#:~:text=Benefits%20of%20Google%20Colab>

Sun, D., Gao, W., Hu, H., & Zhou, S. (2022). Why 90% of clinical drug development fails and how to improve it? *Acta Pharmaceutica Sinica B*, 12(7).

<https://doi.org/10.1016/j.apsb.2022.02.002>

US EPA. (2012). *Sustainable Futures / P2 Framework Manual 2012 EPA-748-B12-001 Appendix F. SMILES Notation Tutorial*.

<https://www.epa.gov/sites/default/files/2015-05/documents/appendf.pdf>

Visikol. (2022). *The Importance of IC50 Determination* | Visikol. Visikol.

<https://visikol.com/blog/2022/06/07/the-importance-of-ic50-determination/>

Zhou, S.-F., & Zhong, W.-Z. (2017). Drug Design and Discovery: Principles and

Applications. *Molecules*, 22(2), 279. <https://doi.org/10.3390/molecules22020279>