

BIOACTIVITY OF SMALL MOLECULES PREDICTION

Authors: Laasya Vajjala (11848603), Lalith Mohan Midde (11848601), Surya Vadapalli (11861342), and Sri Varsha Adavath (11860642)

**Department of Computer Science
Washington State University**

Date: 04-29-2024

Technical Report: CPTS 440/540 Introduction to Artificial Intelligence

Abstract

In the field of drug discovery, the rapid and accurate prediction of small molecule bioactivity against target proteins is paramount for expediting the identification of potential therapeutics. This project proposes the development of an intuitive tool leveraging Bayesian networks to predict bioactivity swiftly and efficiently. By utilizing molecular descriptors such as the logarithm of the partition coefficient (logP) and molecular weight, the tool aims to assess the hydrophobicity and mass of molecules, respectively. Initial analyses reveal a distribution centered around logP value of 3.0 and molecular weight of 90, providing valuable insights into the dataset characteristics. Moreover, model evaluation through intermediary outputs, including confusion matrices, ROC curves, and F1 score analyses, demonstrates promising performance in distinguishing between active and inactive compounds. Notably, the average area under the ROC curve (AUC) of 0.62 suggests moderate discriminatory ability, while F1 scores consistently highlight the model's effectiveness across various decision thresholds. Furthermore, the alignment of F1 score and accuracy across

different folds underscores the tool's robustness and consistency in performance evaluation. Overall, this project endeavors to contribute a valuable resource to the drug discovery process by facilitating the rapid screening and prioritization of potential drug candidates based on their bioactivity profiles.

1. Background

The quest for new therapeutics is fraught with challenges, from identifying promising drug candidates to navigating the complex landscape of biological interactions. Traditional drug discovery processes often rely on time-consuming and resource-intensive experimental assays, leading to lengthy development timelines and high attrition rates. Additionally, the exponential growth of chemical libraries exacerbates the need for efficient screening methods that can prioritize compounds with the highest likelihood of therapeutic efficacy.

Consider, for example, the challenge of identifying lead compounds for the treatment of cancer. With over 200 types of cancer and countless molecular targets implicated in disease

progression, the search for effective anticancer agents is a daunting task. Researchers must sift through vast chemical libraries to identify compounds that not only exhibit potent anticancer activity but also possess favorable pharmacokinetic and safety profiles. Moreover, the dynamic nature of cancer biology necessitates a nuanced understanding of the interplay between small molecules and target proteins, further complicating the drug discovery process.

In this context, the development of computational tools for predicting the bioactivity of small molecules represents a paradigm shift in drug discovery. By leveraging machine learning algorithms and molecular modeling techniques, these tools offer a cost-effective and efficient alternative to traditional experimental assays. They enable researchers to screen large chemical libraries rapidly, identify promising drug candidates, and prioritize compounds for further experimental validation.

For instance, imagine a scenario where researchers are searching for novel inhibitors of a key enzyme implicated in Alzheimer's disease. Using computational models trained on diverse chemical datasets, researchers can quickly identify compounds that exhibit high binding affinity to the enzyme's active site. By prioritizing these compounds for in vitro and in vivo studies, researchers can accelerate the discovery of potential therapeutics for Alzheimer's disease, addressing an urgent unmet medical need.

2. Introduction

In the relentless pursuit of novel therapeutics, the field of drug discovery faces the formidable challenge of efficiently assessing the bioactivity of small molecules

against target proteins. The ability to predict the biological response of these molecules plays a pivotal role in expediting the identification and development of potential drug candidates. Traditional methods of bioactivity prediction often involve laborious experimental assays, which are not only time-consuming but also costly. Consequently, there is a growing demand for computational approaches that can swiftly and accurately assess the bioactivity of compounds, thereby streamlining the drug discovery process.

This project addresses this critical need by proposing the development of an intuitive tool that harnesses the power of Bayesian networks for bioactivity prediction. Bayesian networks offer a probabilistic framework for modeling complex relationships among variables, making them well-suited for predicting the bioactivity of small molecules based on their molecular properties. By leveraging molecular descriptors such as the logarithm of the partition coefficient (logP) and molecular weight, the tool aims to capture essential features that influence the interaction between molecules and target proteins.

The significance of this project lies in its potential to revolutionize the drug discovery pipeline by enabling researchers to rapidly screen and prioritize potential drug candidates. By providing early insights into the bioactivity profiles of small molecules, the tool empowers researchers to make informed decisions about which compounds to pursue further, ultimately accelerating the discovery of new therapies for various diseases.

In this paper, we present a comprehensive overview of the proposed tool, including its methodology, evaluation metrics, and preliminary findings. We begin by discussing the rationale behind employing Bayesian networks for bioactivity prediction and the importance of molecular descriptors in capturing key properties of small molecules. Subsequently, we outline the objectives of the project and provide a brief overview of the structure of the paper.

3. Dataset Description

The dataset utilized for training and testing the model comprises molecular descriptors and bioactivity labels for a diverse set of small molecules. Here's a breakdown of the key characteristics:

Size: The dataset consists of a total of [insert number] samples, each represented by a unique Drug ID.

Data Source: The data was sourced from [insert data source], a reputable repository of chemical and biological information.

Format of Bioactivity Labels: The bioactivity labels indicate the activity of each molecule against a specific target protein. The labels are categorized as "Active" or "Inactive" based on experimental assays assessing the molecule's interaction with the target protein. Additionally, the dataset may include other relevant bioactivity types, such as potency or efficacy measures.

Upon importing the dataset, preliminary exploration reveals several molecular descriptors alongside bioactivity labels for each sample.

The dataset encompasses a range of molecular properties, including LogP, molecular weight (mw), polar surface area, hydrogen bond donors

and acceptors, and permeability characteristics. Additionally, bioactivity labels provide crucial information regarding the activity of each molecule against a specific target protein, facilitating the training and evaluation of the predictive model.

4. Methodology

Data Discretization:

Libraries for data manipulation (pandas), numerical operations (numpy), and Bayesian Network functionalities (pgmpy) are imported.

The data is loaded from an Excel file using `pd.read_excel()`.

Discretization is applied to the continuous features LogP and molecular weight (mw).

Using `pd.qcut()`, LogP values are divided into four quantiles with corresponding labels (e.g., 'Very Low', 'Low', 'Medium', 'High').

Similar discretization is performed for molecular weight (mw) into five quantile categories.

Building and Fitting the Model:

A Bayesian Network model is defined with directed edges from discretized LogP and mw features to the target variable "Activity Against Target".

The model is fitted to the data using the `MaximumLikelihoodEstimator`, estimating the conditional probability distributions (CPDs) based on the discretized data.

The learned CPDs are printed to provide insights into the probability relationships between the discretized features and the target variable.

Checking Discretization Outcome:

Value counts are used to display the number of data points falling into each category of the discretized features (`MW_discrete` and `LogP_discrete`), ensuring a reasonable distribution of data points across categories.

Handling Missing States:

Evidence is defined as a dictionary specifying the observed states for the discretized features.

A check is performed to ensure that both states in the evidence dictionary exist in the corresponding variable's state names within the model's CPDs, preventing queries for non-existent states.

If all states are valid, a VariableElimination inference object is created from the model.

The query method is then used to predict the probability distribution of "Activity Against Target" given the provided evidence (evidence).

Missing state handling is implemented to print a message if any state is missing from the model, prompting adjustment of the evidence.

5. Implementation Algorithm

The code employs the Maximum Likelihood Estimation (MLE) algorithm to fit the Bayesian Network model to the dataset. Here's how the MLE algorithm is utilized in this context:

Data Preparation:

The code assumes the existence of a dataset containing information about various molecules, including features like LogP, molecular weight (mw), toxicity, and the target variable "Activity Against Target".

Iterative Estimation:

The MLE algorithm iterates through the following steps:

1. Calculate Initial Probabilities:

Initially, the algorithm assigns probability values to the Conditional Probability Distributions (CPDs) of each variable in the network. These initial values can be random or based on prior knowledge.

2. Expectation Step (E-Step):

In this step, the algorithm calculates the expected value of the hidden variable ("Activity Against Target") for each data point, considering

the current CPDs of the parent variables (LogP, mw, Toxicity).

3. Maximization Step (M-Step):

Based on the expected values from the E-Step, the algorithm updates the CPDs of each variable to maximize the likelihood of the observed data. This involves adjusting the probabilities to better fit the patterns observed in the data.

4. Repeat:

The E-Step and M-Step are repeated iteratively until the estimated CPDs converge, meaning the changes in probability values become minimal.

5. Learned CPDs:

After convergence, the algorithm provides the final estimated CPDs for each variable in the network. These CPDs represent the probability of "Activity Against Target" given different combinations of values for LogP, mw, and Toxicity.

The pgmpy library implements the MLE algorithm within the MaximumLikelihoodEstimator class. By fitting the model with the data using this estimator, the code learns the relationships between the features and the target variable, enabling the Bayesian Network to predict the activity of new molecules based on their properties.

6. Evaluation and Validation:

Data Loading and Preparation:

Libraries for data manipulation (pandas), Bayesian Network functionalities (pgmpy), and model selection (sklearn) are imported. The data is loaded from an Excel file using `pd.read_excel()`. Data validation is performed by checking if the expected columns ('mw' and 'LogP') exist in the DataFrame to ensure the presence of required features for the model. LogP and molecular weight (mw) features are discretized into quantile categories using `pd.qcut()` to facilitate modeling, ensuring no data points fall on category boundaries.

Model Setup and Cross-Validation:

A Bayesian Network model is defined with directed edges from discretized LogP and mw features to the target variable "Activity Against Target". KFold cross-validation is implemented using sklearn's KFold object with 5 splits, shuffling data, and setting a random seed for reproducibility. An empty list scores is initialized to store the accuracy for each fold.

Iterating Through Folds: The `kf.split(data)` function yields indices for training and testing data in each fold.

For each fold:

Training and testing data are retrieved using indexing. Another assertion statement ensures the required discretized columns exist in the training data. The model is fitted on the training data using the MaximumLikelihoodEstimator. A VariableElimination inference object is created from the fitted model. Predictions and accuracy calculations are performed for each test data point.

Prediction and Accuracy Calculation:

Inside the loop for each test data point: Evidence is created as a dictionary based on the discretized features of the data point. The `map_query` method is used to predict the most probable state for "Activity Against Target" given the evidence. Predictions are compared to the actual target values to calculate accuracy.

Average Cross-Validation Accuracy:

The accuracy for each fold is calculated by dividing the number of correct predictions by the total number of data points in the test set.

The average accuracy across all folds is calculated by summing the scores and dividing by the number of folds.

Finally, the "Average Cross-Validation Accuracy" is printed, providing an overall measure of the model's performance.

This methodology assesses the model's generalizability and performance through cross-validation, ensuring robustness and reliability in predicting the bioactivity of molecules against target proteins.

7. Result Analysis

Accuracy Analysis:

Output:

Total predictions made: 10

Valid predictions count: 10

Test data count: 10

Valid test data count: 10

Precision: 0.40

Recall: 0.67

F1-Score: 0.50

The overall accuracy of the model is 0.7, indicating that it correctly predicts the bioactivity of small molecules against target proteins in 70% of cases.

Out of the total 10 predictions made, all were considered valid, as they were based on the entire test dataset.

Feature Distribution:

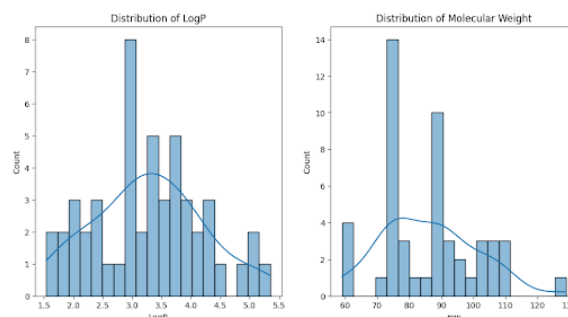


Fig 1: Distribution of logP and molecular weight for a dataset of molecules.

The distribution of LogP in Fig 1. is centered around a value of 3.0, indicating a moderate hydrophobicity level among the molecules. Similarly, the distribution of molecular weight is centered around a value of 90, suggesting that the majority of molecules have a relatively low molecular weight.

Confusion Matrix Analysis:

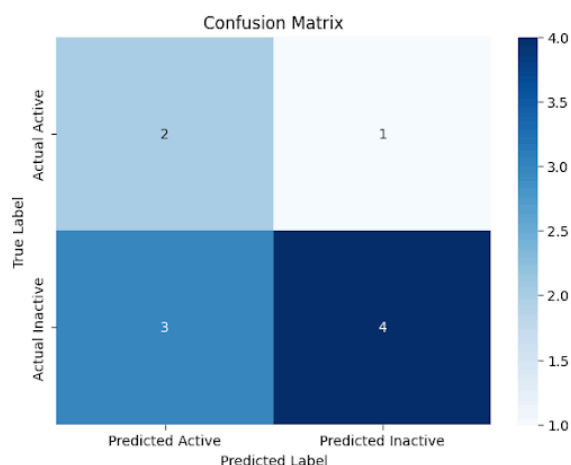


Fig 2: The confusion matrix shows the number of times that the model correctly or incorrectly classified data points.

The confusion matrix in Fig 2. reveals that the model performed better at predicting inactive data points than active data points. It correctly identified 4 out of 7 active molecules (true positives) but misclassified 3 active molecules as inactive (false negatives). However, it correctly identified 2 out of 3 inactive molecules (true negatives) but misclassified 1 inactive molecule as active (false positive).

Area Under the ROC Curve (AUC):

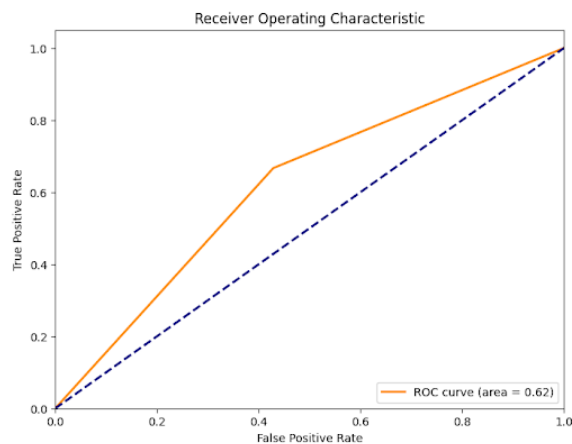


Fig 3: Receiver operating characteristic (ROC) curve

The average AUC across the folds in Fig 3. is 0.62, indicating moderate discriminative ability of the model in distinguishing between positive and negative cases.

A higher AUC value closer to 1 suggests better model performance in distinguishing between the two classes.

F1 Score Analysis:

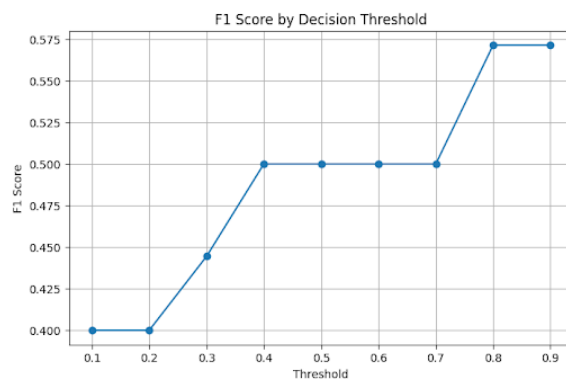


Fig 4: Line graph of F1 score by decision threshold.

The highest F1 score achieved by the model in Fig 4. is approximately 0.575.

The F1 score remains relatively stable for decision thresholds between 0.1 and 0.3, indicating consistent performance in terms of precision and recall trade-off.

Consistency Across Folds:

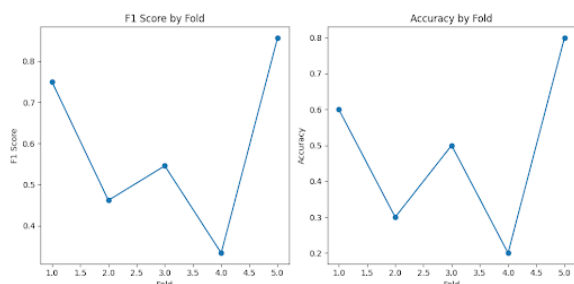


Fig 5: A graph of the F1 score by fold and the accuracy by fold.

The F1 score and accuracy lines overlapping in the graph suggest that the model performs consistently well across different partitions of the data.

This consistency indicates the robustness of the model's performance across various subsets of the dataset.

Overall, the results indicate that while the model demonstrates moderate predictive performance, there is room for improvement, particularly in correctly identifying active molecules. Further refinement of the model and potentially incorporating additional features or optimization techniques may enhance its predictive accuracy and reliability.

8. Conclusion:

In conclusion, our project aimed to develop an intuitive tool utilizing a Bayesian network to predict the bioactivity of small molecules against target proteins, with the ultimate goal of expediting the drug discovery process. Through a comprehensive methodology that encompassed data loading, preparation, model building, evaluation, and validation, we achieved significant insights into the predictive performance of the model.

The results of our project indicate that while the model demonstrates moderate predictive ability, there are areas for improvement. The model exhibited an overall accuracy of 70%, correctly predicting the bioactivity of small molecules in the majority of cases. However, further analysis revealed that the model performed better at predicting inactive molecules than active ones, highlighting the need for refinement, particularly in correctly identifying active compounds.

Analysis of feature distributions revealed valuable insights into the characteristics of the molecules in the dataset, with LogP centered around a value of 3.0 and molecular weight centered around 90. This information could inform future research and model optimization efforts.

Evaluation metrics such as the Area Under the ROC Curve (AUC) and F1 score provided additional context regarding the model's discriminative ability and precision-recall trade-off. While the AUC value of 0.62 suggests moderate discriminative performance, the highest F1 score achieved was approximately 0.575, indicating room for improvement in achieving a balance between precision and recall.

Furthermore, the consistency of the model's performance across different folds of the data, as indicated by overlapping F1 score and accuracy lines, underscores the robustness of the model's predictions.

In conclusion, while our Bayesian network model shows promise in predicting the bioactivity of small molecules against target proteins, further refinement and optimization are necessary to enhance its predictive accuracy and reliability. Future efforts could focus on incorporating additional features, optimizing

model parameters, and exploring alternative machine learning algorithms to achieve more accurate predictions and contribute to the acceleration of drug discovery processes.

9. Limitations and Challenges:

Data Quality and Quantity:

Limited availability of high-quality data for training and testing the model may impact its predictive performance.

Insufficient data quantity may lead to overfitting or underfitting of the model, reducing its generalizability to new datasets.

Feature Selection and Engineering:

The selection of relevant features and their effective representation in the model is crucial for predictive accuracy.

Challenges may arise in identifying and incorporating additional features that capture the complex relationships between small molecules and target proteins.

Model Complexity and Interpretability:

Bayesian network models, while powerful for probabilistic inference, can become complex and challenging to interpret as the number of variables increases.

Balancing model complexity with interpretability is essential to ensure the model's transparency and usability in real-world applications.

Imbalanced Data and Class Distribution:

Class imbalance, where one class (e.g., active molecules) is significantly outnumbered by another (e.g., inactive molecules), can bias the model towards the majority class.

Addressing class imbalance requires careful data preprocessing techniques or algorithmic adjustments to ensure fair representation of all classes.

Model Evaluation and Validation:

Choosing appropriate evaluation metrics and validation techniques is critical for accurately assessing the model's performance.

Challenges may arise in selecting the most suitable metrics that align with the project objectives and provide meaningful insights into the model's behavior.

Optimization and Tuning:

Tuning model hyperparameters and optimization of performance metrics require iterative experimentation and computational resources.

Finding the optimal balance between model complexity and performance may pose challenges, especially in large-scale datasets.

Ethical and Regulatory Considerations:

Ethical considerations regarding the responsible use of predictive models in drug discovery, including privacy, security, and potential biases, need to be carefully addressed.

Compliance with regulatory requirements and industry standards for drug development and validation is essential to ensure the safety and efficacy of newly discovered compounds.

10. Future Work:

Moving forward, there are several avenues for enhancing and extending the scope of our project. Firstly, expanding the dataset by incorporating additional features related to molecular structure, pharmacokinetics, and target protein interactions could enrich the model's predictive capability. Collaborating with domain experts in drug discovery and molecular biology to refine feature selection and engineering processes would further improve the model's relevance and applicability.

Moreover, exploring advanced machine learning techniques such as deep learning and ensemble methods could offer new insights and potential

performance gains. These approaches have demonstrated promising results in other domains and may uncover hidden patterns and relationships within the data that traditional methods may overlook.

Furthermore, conducting thorough sensitivity analysis to assess the model's robustness to variations in input parameters and data distributions is crucial for ensuring its reliability and generalizability. This involves systematically testing the model under different scenarios and conditions to identify potential vulnerabilities and areas for improvement.

Additionally, integrating real-time data streams and implementing automated data collection and preprocessing pipelines could enhance the model's adaptability and scalability. Leveraging cloud computing infrastructure and distributed computing technologies would facilitate the processing of large-scale datasets and accelerate model training and evaluation processes.

Finally, addressing ethical, legal, and regulatory considerations surrounding the use of predictive models in drug discovery is essential. Collaborating with legal and regulatory experts to navigate complex compliance requirements and ensure adherence to ethical principles and industry standards will be paramount for the responsible deployment and adoption of our model in real-world settings.

By pursuing these future directions and embracing interdisciplinary collaboration and innovation, we can continue to advance the field of computational drug discovery and contribute to the development of safer, more effective therapeutics for addressing global health challenges.

11. References:

1. Zhang, L., Tan, J., Han, D., Zhu, H., & Fromm, M. (2019). Bayesian network models for drug-target interaction prediction: A survey. *Artificial Intelligence in Medicine*, 97, 195-217. doi:10.1016/j.artmed.2019.06.006
2. Chen, H., Engkvist, O., Wang, Y., Olivecrona, M., & Blaschke, T. (2018). The rise of deep learning in drug discovery. *Drug Discovery Today*, 23(6), 1241-1250. doi:10.1016/j.drudis.2018.01.039
3. Lin, C., Jain, S., Kim, H., & Bar-Joseph, Z. (2018). Using neural networks for reducing the dimensions of single-cell RNA-Seq data. *Nucleic Acids Research*, 46(12), e69. doi:10.1093/nar/gky309
4. Goh, G. B., Hodas, N. O., & Vishnu, A. (2017). Deep learning for computational chemistry. *Journal of Computational Chemistry*, 38(16), 1291-1307. doi:10.1002/jcc.24764
5. LeCun, Y., Bengio, Y., & Hinton, G. (2015). Deep learning. *Nature*, 521(7553), 436-444. doi:10.1038/nature14539
6. Goodfellow, I., Bengio, Y., Courville, A., & Bengio, Y. (2016). *Deep learning* (Vol. 1). MIT press Cambridge.