



BUTLER

COLLEGE of PHARMACY
and HEALTH SCIENCES

Use of LogBB data to model and predict small molecule Blood Brain Barrier Penetrance

Caiden Lukan, Da Shi, Ph.D., Sarangan Ravichandran, Ph.D., and Amanda Paulson, Ph.D.

¹Butler University, ²ATOM Consortium
2021 ATOM Consortium Summer Internship Program
West Dundee, IL



Abstract

Modelling blood brain barrier (BBB) drug penetrance is essential for the discovery of small molecules, especially since cancer drugs typically do not cross the blood brain barrier. Additionally, an associated issue is the expense that comes along with testing *in vitro*, which is usually necessary for understanding which drugs can cross the BBB. The literature that has already been established contains data with LogBB values (the log of the ratio of concentrations of a drug in the brain to the body at steady state) as well as a categorical value of penetrance. While the data from the literature are used en masse to predict whether or not a molecule crosses the BBB, there are often conflicting values between sources of how a molecule was judged to cross the BBB. The data gathered from the literature were further curated to ensure BBB penetrance was assessed in a consistent manner. After this data curation, a regression and classification model were made. The best models made were neural networks with a validation regression R² score of 0.61 and a classification ROC AUC score of 0.94. Additionally, feature importance was analyzed for the best models. Features of a molecule that indicate polarity are a helpful measurement in determining the permeability of that molecule. Using artificial intelligence and machine learning techniques, models with high accuracy can be made and used for prediction of BBB penetrance in the future.

Introduction

- The Blood Brain Barrier (BBB) is a barrier between the Central Nervous System and the rest of the body that prevents molecules from causing CNS damage.
 - The metric that is used as a parameter to measure whether or not a molecule passes through the BBB is LogBB.
- $$\log(BB) = \log\left(\frac{[Brain]}{[Blood]}\right)$$
- LogBB is best used to measure concentrations of a drug at equilibrium.
 - LogBB < 0 means that a molecule does not pass the BBB.
 - LogBB ≥ 0 means that a molecule passes the BBB.

Materials and Methods

- Seven journal articles were used as sources of LogBB data.
- The assays that are used to calculate logBB are Parallel Artificial Membrane Permeability Assay (PAMPA) and Immobilized Artificial Membrane technique.
- The data from the journal article was the curated in order to improve the accuracy of the models.
- From the curated data, a Regression model and Classification model were made.
- The Regression model is useful for predicting a LogBB value.
- The Classification model is useful to predict whether or not a molecule will pass through the BBB.
- Molecules were then featurized and trained with an Xgboost, Neural Network, and Random Forest model.
- Parameters for the models were analyzed and improved to improve the score of the valid set of data.
- Finally, results are analyzed and the best model can be used as a standard for predicting LogBB values or whether or not a molecule will pass the BBB.

Data Curation

- Collected 2590 compounds from 7 different sources.
- Converted Compounds to base_rdkit_smiles to look for duplicate compounds.
- Removed duplicate compounds to have a final total of 1181 compounds.
- Prior to moving to the modelling step, LogBB data was analyzed by a histogram (**Figure 1**) and Penetration data was analyzed with a bar graph (**Figure 2**) to ensure a proper balance of data.
- The data of the compounds structure was also analyzed with a heat map (**Figure 3**) and tanimoto distance plot (**Figure 4**) to ensure a proper diversification of SMILES strings.

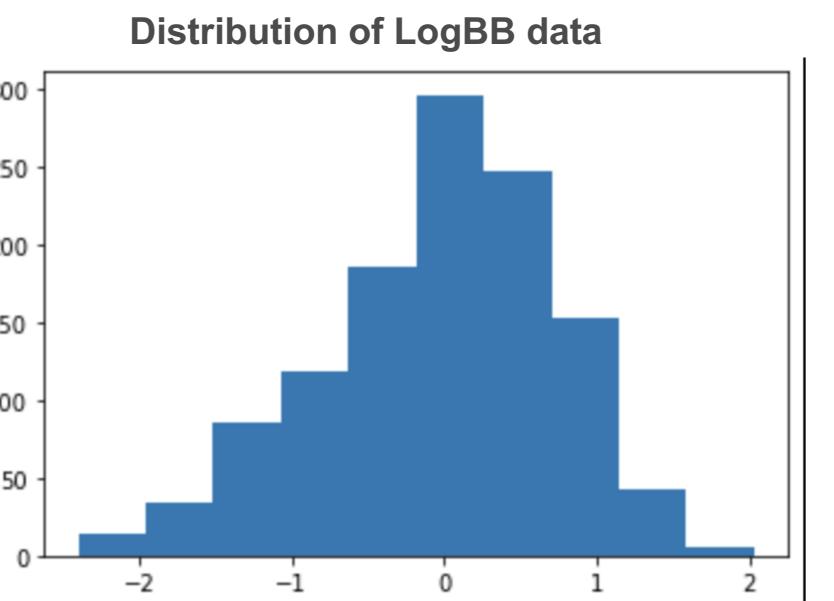


Figure 1. Distribution of LogBB data shows a relatively even distribution of values. The minimum value is -2.4 and the maximum value is 2.03.

Penetration Plot

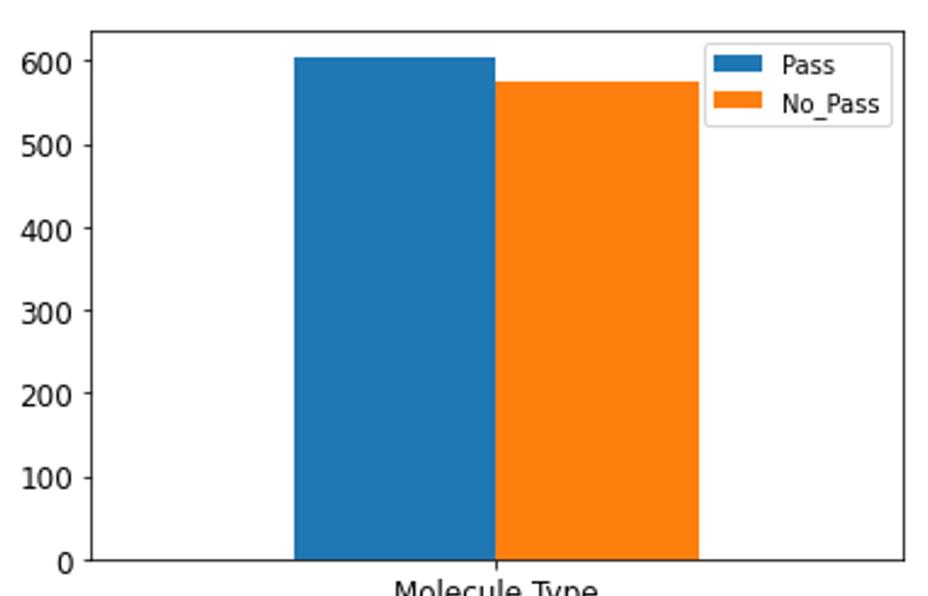


Figure 2. A graphical representation of the molecule that pass (605) and do not pass (576) the blood brain barrier

Data Curation

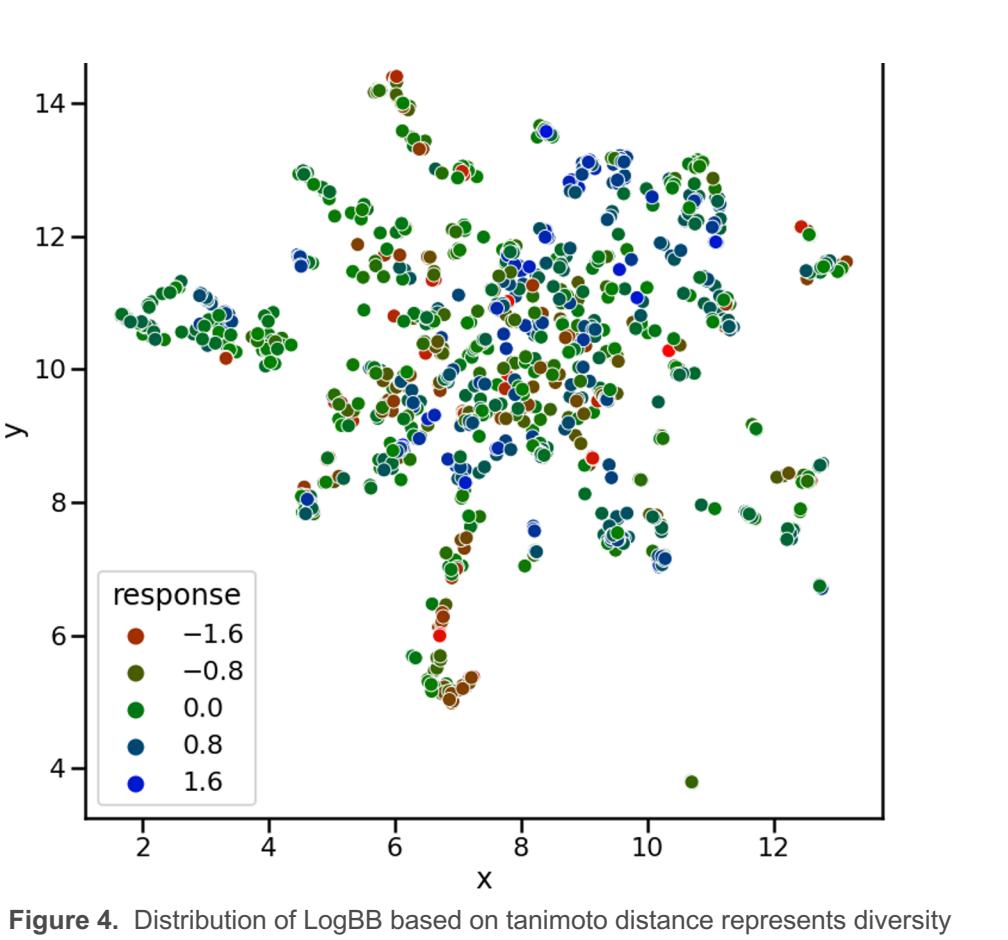


Figure 4. Distribution of LogBB based on tanimoto distance represents diversity of compounds

Modelling

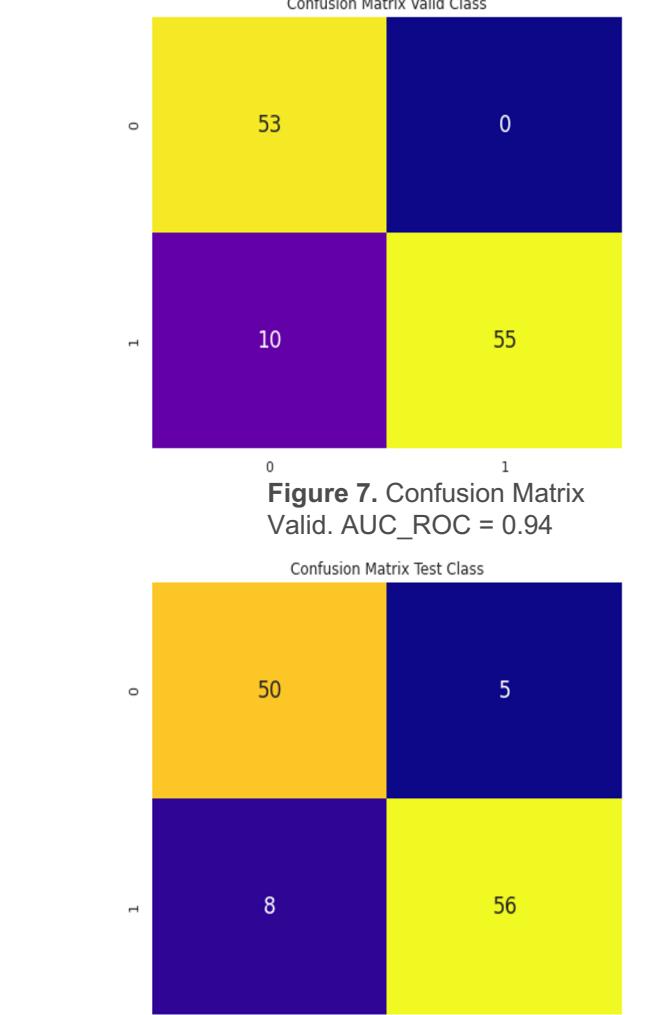


Figure 7. Confusion Matrix Valid. AUC_ROC = 0.94

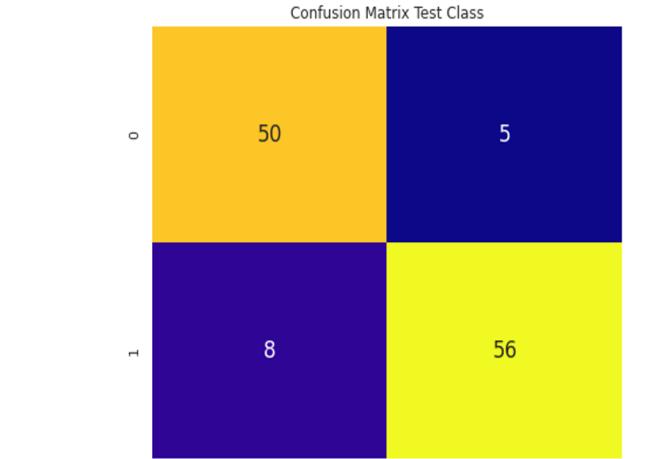


Figure 8. Confusion Matrix Test. AUC_ROC = 0.83

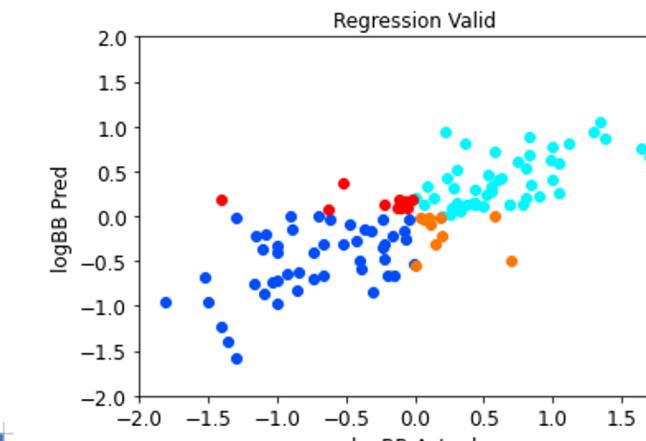


Figure 9. Regression Valid. r² = 0.61. ROC_AUC based on threshold of 0 = .83

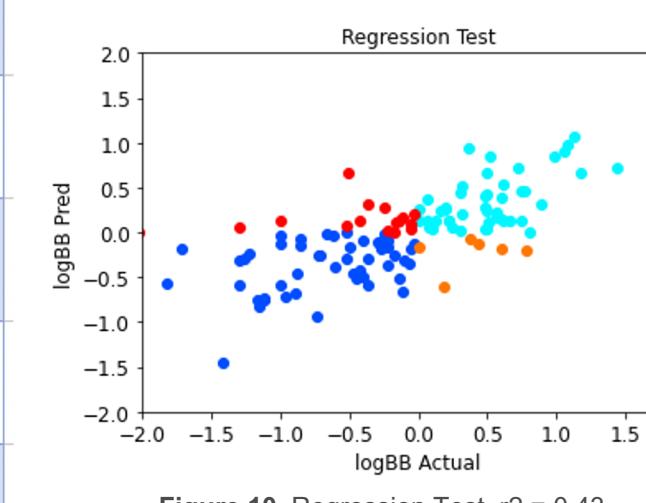


Figure 10. Regression Test. r² = 0.43. ROC_AUC based on threshold of 0 = .82

Data Exploration

- Comparing the regression and classification model with the shuffling method (**Figure 13**), the most important features were Qed and PEOE_VSA6 respectively.
- A change in PEOE (Partial Equalization of Orbital Electronegativities) would change the polarity of the molecule. This means that a change in polarity has an effect on the molecules ability to penetrate the BBB.
- A change in Qed (Quantitative estimation of drug likeness) is similar to changing parameters in Lipinski's rule of 5. This means that a change in those properties has a great effect on the molecules ability to penetrate the BBB.

Conclusion

- Molecules that can cross the BBB have properties such as being non-polar based on a feature importance algorithm.
- Neural networks show the best results for all feature types.
- BBB penetrance modelling is important for future drug development regardless of whether or not the purpose of the drug is to cross.

Future Aims

-

The best regression and classification models can be used to predict molecules that will pass through the blood brain barrier.

The models can be improved when more assays are completed with small molecules that have unknown LogBB values.

Compare effectiveness with molecular dynamic models.

References

Radchenko EV, Dyabina AS, Palyulin VA. Towards Deep Neural Network Models for the Prediction of the Blood-Brain Barrier Permeability for Diverse Organic Compounds. *Molecules*. 2020; 25(24):5901. <https://doi.org/10.3390/molecules25245901>

Singh S, Zeng H, Li H, et al. Development of blood brain barrier permeation prediction models for organic and inorganic biocidal active substances. *Chemosphere*. 2021; 277: 130330. <https://doi.org/10.1016/j.chemosphere.2021.130330>

Vilar S, Cherkarati M, Costanzo S. Prediction of passive blood-brain partitioning: Straightforward and effective classification model based on *in silico* derived physicochemical descriptors. *J. Mol. Graph. Model.* 2010; 28(8): 89. <https://doi.org/10.1016/j.jmgm.2010.03.008>

Zhang J, Zhang H, Li H, et al. Targeted QSAR modeling of the blood-brain barrier permeability for diverse organic compounds. *Pharm Res*. 2008;25(8): 1902-1914. doi:10.1007/s10908-006-9069-0

Plisson F, Piggott AM. Predicting Blood-Brain Barrier Permeability of Marine-Derived Kinase Inhibitors Using Ensemble Classifiers Reveals Potential Hits for Neurodegenerative Disorders. *Marine Drugs*. 2019; 17(2):81. <https://doi.org/10.3390/mdd17020081>

Zhu, L., Zhao, J., Zhang, Y., et al. ADME properties evaluation in drug discovery: *in silico* prediction of blood-brain partitioning. *Mol Divers* 22, 979–990 (2018). <https://doi.org/10.1007/s10021-018-0986-8>

Ciriani L. Feature Importance with Neural Network. Medium. <https://towardsdatascience.com/feature-importance-with-neural-network-346eb6205743>. Published January 21, 2020. Accessed July 13, 2021.

This project has been funded in whole or in part with federal funds from the National Cancer Institute, National Institutes of Health, under contract HHSN26120080001E. The content of this publication does not necessarily reflect the views or policies of the Department of Health and Human Services, nor does mention of trade names, commercial products, or organizations imply endorsement by the U.S. Government.

All animals used in this research project were cared for and used humanely according to the following policies: the U.S. Public Health Service Policy on Humane Care and Use of Animals (2000); the Guide for the Care and Use of Laboratory Animals (1985); and the U.S. Government Principles for Utilization and Care of Vertebrate Animals Used in Testing, Research, and Training (1985). All Frederick National Laboratory animal facilities and the animal program are accredited by the Association for Assessment and Accreditation of Laboratory Animal Care International.

I would like to thank Amanda Paulson for directing me through this project and her excellent teaching skills. Next, I would like to thank Da Shi for helping me with resolving coding troubles that I encountered during the time he was with ATOM.

Finally, I would like to thank Sarangan Ravichandran (Ravi) for answering additional questions of mine about docking and other steps of the process of modelling.

This project has been funded in whole or in part with federal funds from the National Cancer Institute, National Institutes of Health, under contract HHSN26120080001E. The content of this publication does not necessarily reflect the views or policies of the Department of Health and Human Services, nor does mention of trade names, commercial products, or organizations imply endorsement by the U.S. Government.

All animals used in this research project were cared for and used humanely according to the following policies: the U.S. Public Health Service Policy on Humane Care and Use of Animals (2000); the Guide for the Care and Use of Laboratory Animals (1985); and the U.S. Government Principles for Utilization and Care of Vertebrate Animals Used in Testing, Research, and Training (1985). All Frederick National Laboratory animal facilities and the animal program are accredited by the Association for Assessment and Accreditation of Laboratory Animal Care International.

Figure 5. Parameters chosen for Hyperparameter Optimization of Classification Model

	RF	xgboost
Rdkit	rfe_choice = [63,32,64,28,26]	xgb_choice = [0,5,0,1,0,15,0,2]
ECFP	rfe_choice = [12,8,19,25,26]	xgb_choice = [0,1,0,2,0,3]
moe	rfe_choice = [12,8,19,25,26]	xgb_choice = [0,1,0,2,0,3]
Mordred	rfe_choice = [12,8,19,25,26]	xgb_choice = [0,1,0,2,0,3]
GCN	rfe_choice = [10,0,1,0,0,0,0,0,0,0]	xgb_choice = [0,1,0,2,0,25]

Figure 6. Parameters chosen for Hyperparameter Optimization of Regression Model

Data Exploration

- Feature importance of the classification and regression model is important in understanding what part of the molecule has the greatest impact on the outcome.
- The way that feature importance can be accomplished is via shuffling the data. "If we, with our shuffle, break a strong relationship will compromise what our model has learned during training, resulting in higher errors (high error = high importance)".

X_1	X_2	X_...	y
X_1a	X_2a	:	y_1
X_1b	X_2b	:	y_2
X_1c	X_2c	:	y_3
:	:	:	:
X_1z	X_2z	X....	y_n?

Figure 13. Steps on how feature shuffling is done. The first dataframe is unchanged. The second dataframe shows the first feature being shuffled while the second columns remains as it was in the initial dataframe and the result column changing. The third dataframe shows the second feature being shuffled while the first feature remains as it was in the initial dataframe and the result column changing. This process is repeated for as many features are in the featurizer.

Acknowledgements

I would like to thank Amanda Paulson for directing me through this project and her excellent teaching skills. Next, I would like to thank Da Shi for helping me with resolving coding troubles that I encountered during the time he was with ATOM.

Finally, I would like to thank Sarangan Ravichandran (Ravi) for answering additional questions of mine about docking and other steps of the process of modelling.

This project has been funded in whole or in part with federal funds from the National Cancer Institute, National Institutes of Health, under contract HHSN26120080001E. The content of this publication does not necessarily reflect the views or policies of the Department of Health and Human Services, nor does mention of trade names, commercial products, or organizations imply endorsement by the U.S. Government.

All animals used in this research project were cared for and used humanely according to the following policies: the U.S. Public Health Service Policy on Humane Care and Use of Animals (2000); the Guide for the Care and Use of Laboratory Animals (1985); and the U.S. Government Principles for Utilization and Care of Vertebrate Animals Used in Testing, Research, and Training (1985). All Frederick National Laboratory animal facilities and the animal program are accredited by the Association for Assessment and Accreditation of Laboratory Animal Care International.