

# **DRUG CHARACTER PREDICTION USING BLOOD BRAIN BARRIER PENETRATION DATASET**

A project Report  
submitted in the fulfilment of the requirements for the course  
CSCI B-565 Data Mining

By

*Durga Sai Sailesh Chodabattula*

*Sai Jagan Reddy Lakku*

*Raja Simha Reddy Allampati*

*Dhanusha Duraiyan*

Under the supervision of

**Dr. Yuzhen Ye**

Associate Professor  
Luddy School of Informatics, Computing and Engineering, Indiana  
University Bloomington

## ABSTRACT

The blood-brain barrier (BBB) has proven to be a significant obstacle to drug delivery to the brain. The BBB in a healthy brain serves as a diffusion barrier that prevents most substances from passing from the bloodstream to the brain, allowing only tiny molecules to pass through. The BBB is interrupted in specific pathological states of diseases like stroke, diabetes, seizures, multiple sclerosis, Parkinson's disease, and Alzheimer's disease. As a result, breaking through the barrier has long been a problem in the development of medications that target the central nervous system. Using blood brain barrier penetration (BBBP) dataset which includes binary labels for over 2000 compounds on their permeability properties, the permeability of the drug for the test set were checked. This project proposes deep learning methods and neural networks to predict the BBB permeability based on the clinical phenotypes data. Different models like linear, lasso, Elastic Net, SVM, and Naive Bayes were used to train the data and to predict the performance of the test data. The results indicated that neural network approaches may considerably increase drug BBB penetration prediction accuracy, allowing researchers to minimize clinical trials and identify novel CNS medicines.

*Keywords: BBB, BBBP, Deep learning methods, Neural networks, Gradient boosting, correlation*

## TABLE OF CONTENTS

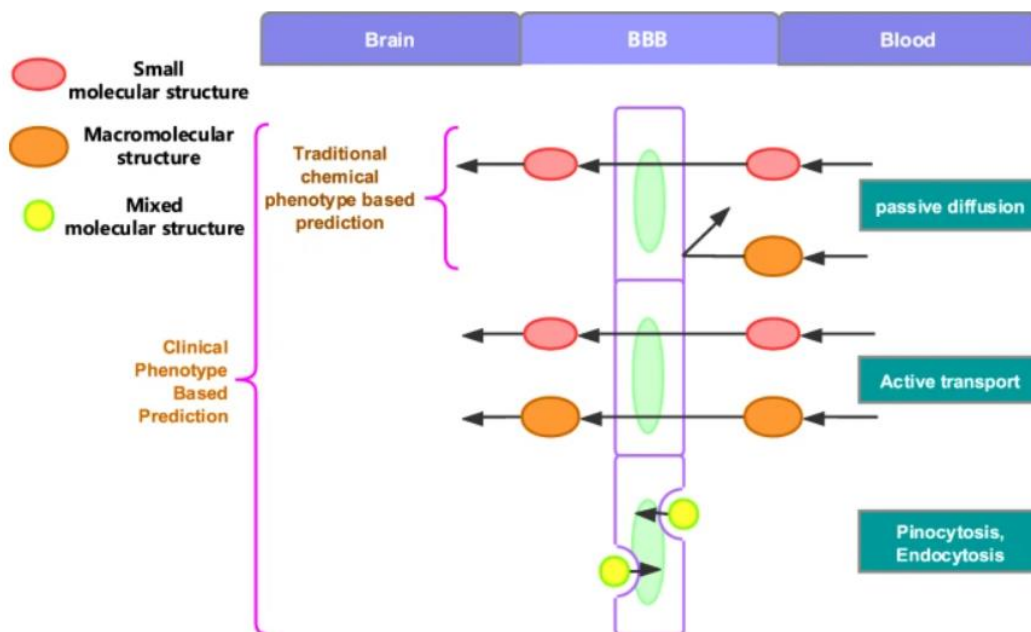
Introduction	4
Experimental Investigations	5
Results	9
Discussion	15
References	16

## 1. INTRODUCTION

Brain disorders, such as central nervous system (CNS) issues and brain tumors, are among the most frequent, fatal, and undertreated ailments. Because the number of seniors and patients with CNS issues is increasing, global drug development for brain diseases will need to increase significantly during the next 20 years. However, as compared to other therapeutic sectors, the discovery of medications for brain ailments has the lowest success rates. The time it takes to develop CNS drugs is frequently significantly longer than the time it takes to develop non-CNS treatments. Clinical trials of CNS drugs are challenging due to the brain's complexity, adverse effects, and the blood-brain barrier's (BBB) impermeability.

CNS drug development is limited by a lack of adequate technologies for transporting drugs over the BBB, in addition to the complications of brain disorders. Small and macromolecules are being investigated as possible therapeutic agents for a wide range of brain diseases. The BBB can only be crossed by minuscule lipid-soluble molecules with a molecular weight of less than 400 Da; most macromolecules cannot pass through the brain endothelium. The BBB is a physiological barrier that stops 95% of chemicals from becoming medicines. The BBB filters the blood, including the solutes found in the blood going to the brain.

To overcome all these complications in the process of drug discovery and obstacles in caused by the BBB to the drugs, we use the Blood Brain Barrier Penetration Dataset (BBBP) which gives the information of whether a drug can get through the BBB's impermeability.



**Figure 1.** Mechanisms of drug transport through the BBB and the applicability of prediction methods (Miao et al., 2019).

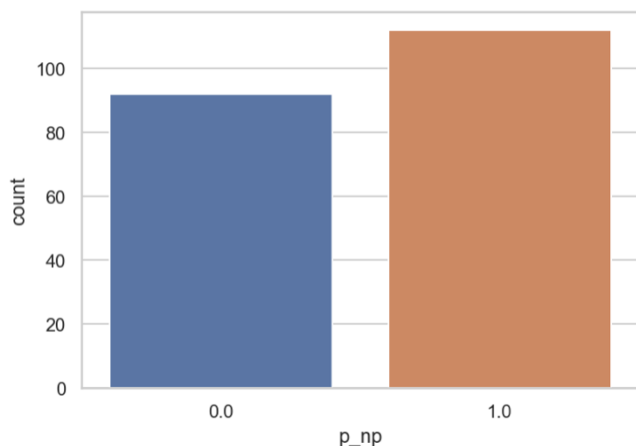
## 2. EXPERIMENTAL INVESTIGATIONS:

### 2.1 Dataset

DeepChem provides a high-quality open-source toolchain that democratizes the application of deep-learning in drug discovery, materials science, quantum chemistry, and biology.

The BBBP dataset from the deepchem module was used. BBBP dataset contains 2053 items with four attributes :

- the index number from 1 to 2053 (“num”)
- the name of the compound (“name”)
- the penetrating or nonpenetrating properties (“p\_np”)
- the SMILES string of the compound (“smiles”)



**Figure 2.** Class balance of p\_np

### 2.1.1 Class balance

As the count of p\_np for values of 0 and 1 are close to each other, we could say that the classes are balanced.

### 2.1.2 Data splitting and cross validation

The deepchem module primarily uses the scaffold split by default to perform the function. Scaffold split algorithm in MoleculeNet splits the given dataset into training, validation, and test data. It creates an unbalanced split which makes the prediction harder.

Initially, the molecular compounds are grouped into scaffold sets on the basis of the skeletal ring structures termed scaffolds then the compounds and scaffold sets are sorted in reverse order (largest to the smallest). Finally the dataset are split into training, validation, and test data in an 8:1:1 ratio from the top.

In scaffold splitting, the sets are sorted from the largest to the smallest. The test data will consist of compounds that are less related to others. This leads to the overfitting of the data which ends up in poor accuracy in both qualitative and quantitative prediction as the split is biased.

So to overcome this problem we used k – fold cross validation technique along with shuffling is done to increase the randomness in data for the train and test data. In k-fold, the training set is subdivided into k smaller subsets (other approaches are described below, but generally follow the same principles). For each of the k "folds," the following approach is used. A model is trained using the k-1 folds as training data, and the resulting model is verified using the remaining data

(i.e., it is used as a test set to compute a performance measure such as accuracy). The values obtained are then averaged to give the performance metric of the different models we use in our project.

### *2.1.3 Conversion of the data to SMILES*

The simplified molecular-input line-entry system (SMILES) is a specification in form of a line notation for describing the structure of chemical species using short ASCII strings. There are almost 1800 characteristics of each smile which needs to be examined to determine which characteristics influence the penetration of a drug.

When using the raw .csv file, the raw data has to be converted to SMILES but while using the deepchem module this is simplified as the dataset already has the SMILES in it.

### *2.1.4 Conversion of SMILES to fingerprints*

Deep learning models nearly always require numerical arrays as inputs. We need to represent each molecule as one or more arrays of numbers if we wish to process molecules with them. Many (but not all) types of models need fixed-size inputs. This can be difficult for molecules because the amount of atoms in each molecule varies. If we wish to employ these models, we must somehow represent variable-sized molecules with fixed-sized arrays.

Fingerprints are intended to address these issues. A fingerprint is a fixed-length array in which distinct entries represent the existence of various properties in the molecule. If two molecules have similar fingerprints, it means they have many of the same traits and, as a result, will most likely have comparable chemistry.

DeepChem recognizes a type of fingerprint known as a "Extended Connectivity Fingerprint," or "ECFP" for short. They are also known as "circular fingerprints" at times. The ECFP method starts by categorizing atoms only based on their direct characteristics and bonds. Each distinct pattern is a distinguishing trait. For example, "carbon atom coupled to two hydrogens and two heavy atoms" is a feature, and for each molecule that includes that feature, a specific element of the fingerprint is set to 1. It then looks at bigger circular regions to identify new traits repeatedly. A higher level feature is formed when one specific feature is bound to two additional specific features, and the associated element is set for every molecule that includes it.

## ***2.2 Models***

The following models were used in the project:

- Linear Regression
- Lasso classifier
- SVC Model
- Naïve Bayes Classifier
- Neural Network
- Gradient Boosting on Linear regression and SVC model

### *2.2.1 Linear Regression*

Linear regression attempts to model the relationship between two variables by fitting a linear equation to observed data. One variable is considered to be an explanatory variable, and the other is considered to be a dependent variable.

### *2.2.2 Lasso Classifier*

Lasso (least absolute shrinkage and selection operator) is a regression analysis method that performs variable selection as well as regularization to improve the prediction accuracy and interpretability of the resulting statistical model. Lasso is a linear model with L1 prior as regularizer.

### *2.2.3 Support Vector Classifier (SVC)*

The Linear Support Vector Classifier (SVC) method applies a linear kernel function to perform classification and it performs well with a large number of samples. If we compare it with the SVC model, the Linear SVC has additional parameters such as penalty normalization which applies 'L1' or 'L2' and loss function. The kernel method cannot be changed in linear SVC, because it is based on the kernel linear method.

### *2.2.4 Naïve Bayes Classifier*

It is a classification technique based on Bayes Theorem with an assumption of independence among predictors. In simple terms, a Naive Bayes classifier assumes that the presence of a particular feature in a class is unrelated to the presence of any other feature.



### 2.2.5 Neural network

Artificial neural network models are behind many of the most complex applications of machine learning. Classification, regression problems, and sentiment analysis are some of the ways artificial neural networks are being leveraged today. As an emerging field, there are many different types of artificial neural networks. They differ for a variety of reasons, including: B. Complexity, network architecture, density, and data flow. However, different types share a common goal of modelling neuronal behaviour and replicating neuronal behaviour to improve machine learning.

### 2.2.6 Gradient Boosting on Linear regression and SVC model

Gradient boosting is a machine learning technique used especially for regression and classification tasks. This produces a predictive model in the form of an ensemble of weak predictive models, which is usually a decision tree. If the decision tree is a weak learner, the resulting algorithm is called a gradient boosting tree. Usually better than Random Forest. The gradient boosting tree model is built in stages like any other boosting method but generalizes the other methods by allowing optimization of the differentiable loss function.

## 3. RESULTS

### 3.1 Linear Regression

	precision	recall	f1-score	support
0.0	0.77	1.00	0.87	71
1.0	1.00	0.84	0.91	133
accuracy			0.90	204
macro avg	0.89	0.92	0.89	204
weighted avg	0.92	0.90	0.90	204

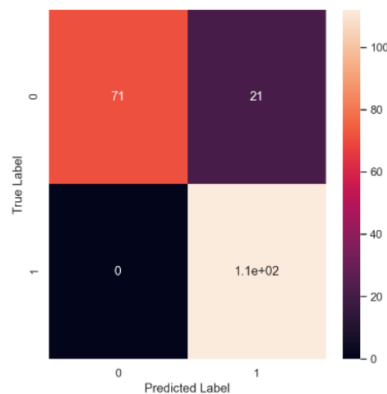
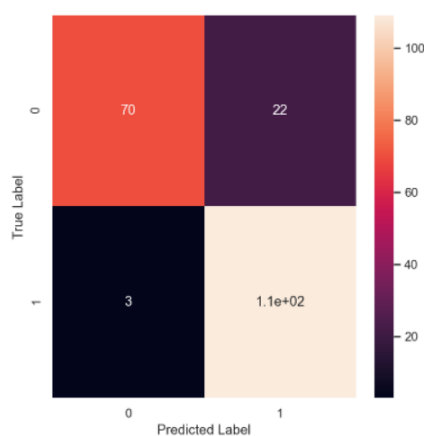


Figure 3

Overall, the accuracy for the model is 0.90. From the confusion matrix we can see that false negatives are not found.

### 3.2 Lasso Classifier

	precision	recall	f1-score	support
0.0	0.76	0.96	0.85	73
1.0	0.97	0.83	0.90	131
accuracy			0.88	204
macro avg	0.87	0.90	0.87	204
weighted avg	0.90	0.88	0.88	204



**Figure 4**

From the scores, we can see that precision for class 1.0 is high and yet the recall for it lower than class 0 but the F1 score is higher. Overall, the accuracy for the model is 0.88. From the confusion matrix we can see that false negatives are low as 3.

### 3.3 Support Vector Classifier (SVC)

	precision	recall	f1-score	support
0.0	0.72	1.00	0.84	66
1.0	1.00	0.81	0.90	138
accuracy			0.87	204
macro avg	0.86	0.91	0.87	204
weighted avg	0.91	0.87	0.88	204

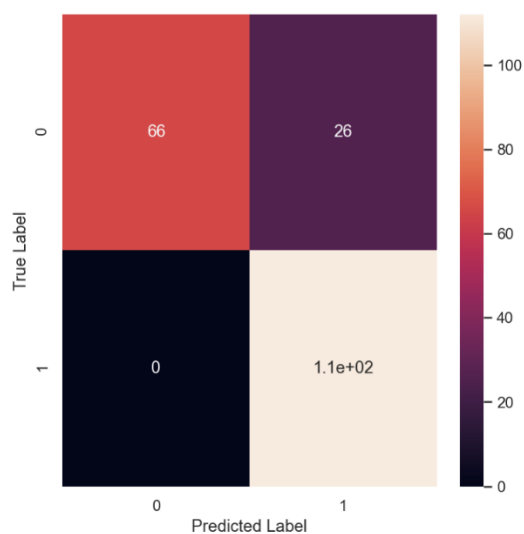


Figure 5

The confusion matrix shows that the model had no false negative results and the overall accuracy for the model is 0.87.

### 3.4 Naïve Bayes Classifier

	precision	recall	f1-score	support
0.0	0.39	0.44	0.41	82
1.0	0.59	0.54	0.56	122
accuracy			0.50	204
macro avg	0.49	0.49	0.49	204
weighted avg	0.51	0.50	0.50	204

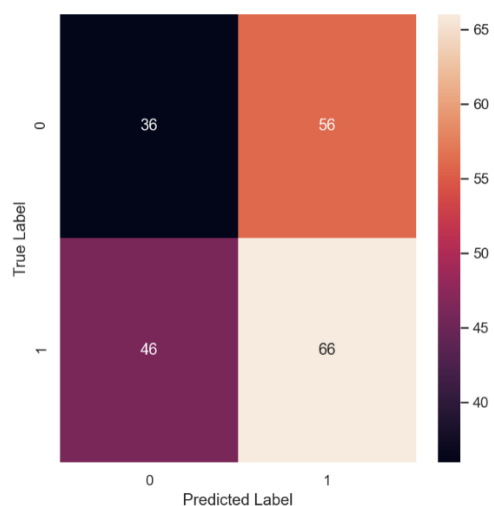


Figure 6

The naïve bayes classifier doesn't have a good performance as the accuracy comes only upto 50% and it has high counts of both false positives and false negatives which is also reflected in the poor F1 scores.

### 3.5 Neural network

```
In [83]: pred = bbbp_svc_gb_calib.predict(X_test)
         f1_score(y_test,pred)
```

```
Out[83]: 0.9226006191950465
```

```
In [84]: pred = bbbp_svc_gb_calib.predict_proba(X_test)
         roc_auc_score(y_test,pred[:,1])
```

```
Out[84]: 0.9038461538461539
```

```
In [85]: pred = pred[:,1]
         pred_svc_gb = np.copy(pred)
         pred[pred<=threshold] = 0
         pred[pred>threshold] = 1
         svc_gb_score = f1_score(y_test,pred)
         print(svc_gb_score)
```

```
0.9221556886227544
```

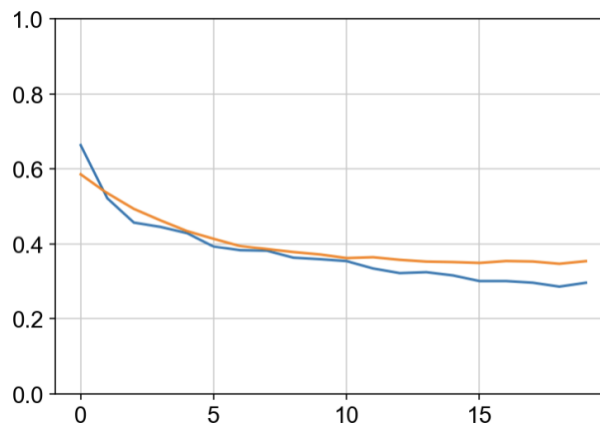


Figure 7

The F1 score for the neural networks is 0.92 which is pretty good. The test and train loss are depicted in the graph.

### 3.6 Gradient Boosting on Linear regression and SVC model

```
In [98]: pred = bbbp_xgb_gb_calib.predict(X_test)
         f1_score(y_test,pred)
```

```
Out[98]: 0.9259259259259259
```

```
In [99]: pred = bbbp_xgb_gb_calib.predict_proba(X_test)
         roc_auc_score(y_test,pred[:,1])
```

```
Out[99]: 0.9180021367521367
```

```
In [100]: pred = pred[:,1]
          pred_xgb_gb = np.copy(pred)
          pred[pred<=threshold] = 0
          pred[pred>threshold] = 1
          xgb_gb_score = f1_score(y_test,pred)
          print(xgb_gb_score)
```

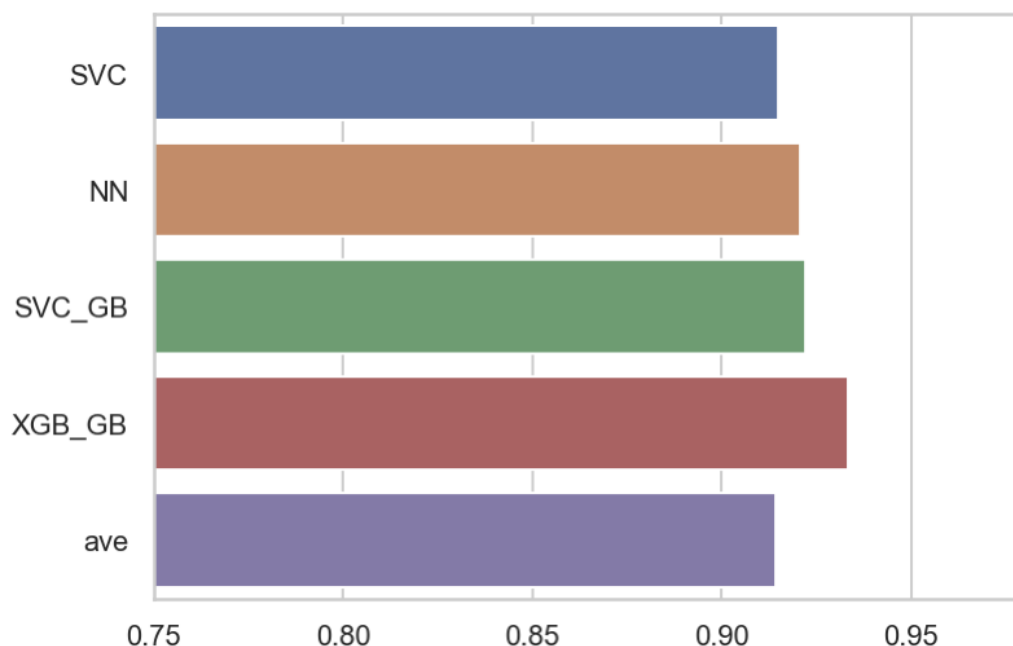
```
0.9333333333333333
```

**Figure 8**

Gradient boosting on linear regression and SVC models have elevated F1 score and accuracy. Also they have high roc\_auc score which is better.

### 3.7 Overall performance of all the models

	Model	Validation fold 1 Score	Validation fold 2 Score	Validation fold 3 Score	Validation fold 4 Score	Validation fold 5 Score	Mean AUPR Score
0	Logistic Regression	0.912470	0.922513	0.928506	0.937749	0.920286	0.924305
1	Lasso Classifier (L1)	0.928900	0.906990	0.920275	0.913917	0.916375	0.917291
2	Support Vector Classifier	0.900022	0.913080	0.909301	0.895114	0.928481	0.909200
3	Gaussian Naive Bayes Classifier	0.843203	0.884258	0.878238	0.826347	0.845667	0.855543



**Figure 9**

The overall performance of the all the models show that the gradient boosting performed better than all the other models tested followed by neural networks.

### 3.8 Correlation

```
n_np is highly correlated with ATSC1c    0.50154
Name: p_np, dtype: float64
Top 10 features that have decisive role on permeability are
MATS1are    0.310021
Name: p_np, dtype: float64
MATS1pe    0.315559
Name: p_np, dtype: float64
AATSC1c    0.333928
Name: p_np, dtype: float64
Lipinski    0.334814
Name: p_np, dtype: float64
SsssCH    0.360272
Name: p_np, dtype: float64
SdssC    0.385012
Name: p_np, dtype: float64
ATSC1pe    0.40103
Name: p_np, dtype: float64
ATSC1se    0.410719
Name: p_np, dtype: float64
ATSC1are    0.4139
Name: p_np, dtype: float64
ATSC1c    0.50154
Name: p_np, dtype: float64
```

**Figure 10**

## 4. DISCUSSION

### *4.1 Comparison of different models and neural networks*

We used the train dataset to train the 5 ML models and incorporated a 5-fold cross validation (5-CV) in the process. We conducted the model predictions on the unseen Test data and displayed the results after training models and evaluated performance with 5-CV.

All models except Naïve bayes have very low False Negative score, but almost all models have a little high False Positive score. The classification reports show that even though accuracy of logistic regression are very high, these models have trouble classifying non-penetratable targets. This could be because there is more penetratable class data than non-p class data.

One of the reasons why Logistic Regression outperforms other models is because we did not include any regularization values -  $C=1$  for logistic. So, by doing so, the model considers all features as part of its learning process and does not completely shrink/remove any feature weights, which is important in this biological problem statement where each feature has a meaningful contribution to finding the target values and thus including all of them helps the model get better results.

Furthermore, the main reason Gaussian Naive Bayes performs poorly (low accuracy, very low recall, precision, and high FN, FP) is due to its main idea of considering features to be independent of each other - this assumption will not work in many cases, particularly in biological problems where each molecular description and chemical feature is important to understanding the pattern in data.

Usually, the neural networks have better results than gradient boosted models. Here from figure 9, we can see that gradient boosted models work better for the dataset than the neural networks. This could be because of the number of datapoints available. The datapoints might not be sufficient for the neural network to work at its best.

### *4.2 Significance of the study*

It is estimated that more than 98 percent of small organic molecules (drugs) will not penetrate the BBB. The time taken to develop drugs for CNS related diseases can be reduced significantly, as we can eliminate the possibility of a chemical being useful as a drug which can cross the blood

brain barrier. This minimizes the number of clinical trials to be conducted which saves resources, time and is cost effective.

#### *4.3 Cost of false positives and false negatives*

When a drug is predicted that it can penetrate the barrier but actually it doesn't, it is a false positive. This costs us a lot of time and clinical trials because when a chemical is considered as predictable, further research and work is done in developing the drug and when we figure out that it is actually not penetrable, efforts are gone for no use.

Similarly when a drug is predicted negative and is actually positive, it is false negative. False negatives doesn't cost more because as a result we just lose one chemical which could have potentially become a drug.

Efforts and investment wise, false positives are costly. At the expense of lives being lost due to a disease, false negatives are costly.

#### *4.4 Correlation*

When correlation was calculated between the variables it was found that  $n_{np}$  is highly correlated with ATSC1c with the value of 0.50154. Also, the top 10 features which are decisive on permeability was found as seen in figure 10. This could infer that these features could influence a chemical being able to penetrate or not penetrate BBB.

### **5. REFERENCES**

- Abbott, N. Joan, N. Joan Abbott, Adjanie A. K. Patabendige, Diana E. M. Dolman, Siti R. Yusof, and David J. Begley. 2010. "Structure and Function of the Blood-brain Barrier." *Neurobiology of Disease*. <https://doi.org/10.1016/j.nbd.2009.07.030>.
- Banks, William A. 2009. "Characteristics of Compounds That Cross the Blood-Brain Barrier." *BMC Neurology* 9 Suppl 1 (June): S3.



- Geldenhuys, Werner J., Afroz S. Mohammad, Chris E. Adkins, and Paul R. Lockman. 2015. "Molecular Determinants of Blood–brain Barrier Permeation." *Therapeutic Delivery* 6 (8): 961–71.
- Guerra, Angela, Juan A. Pérez, and Nuria E. Campillo. 2008. "Artificial Neural Networks in ADMET Modeling: Prediction of Blood-Brain Barrier Permeation." *QSAR & Combinatorial Science* 27 (5): 586–94.
- Miao, Rui, Liang-Yong Xia, Hao-Heng Chen, Hai-Hui Huang, and Yong Liang. 2019. "Improved Classification of Blood-Brain-Barrier Drugs Using Deep Learning." *Scientific Reports* 9 (1): 8802.
- MoleculeNet. Available online: <https://moleculenet.org/> (accessed on 29 March 2022).
- Sakiyama, Hiroshi, Motohisa Fukuda, and Takashi Okuno. 2021. "Prediction of Blood-Brain Barrier Penetration (BBBP) Based on Molecular Descriptors of the Free-Form and In-Blood-Form Datasets." *Molecules* 26 (24). <https://doi.org/10.3390/molecules26247428>.
- Weininger, David. 1988. "SMILES, a Chemical Language and Information System. 1. Introduction to Methodology and Encoding Rules." *Journal of Chemical Information and Computer Sciences* 28 (1): 31–36.
- Wu, Zhenqin, Bharath Ramsundar, Evan N. Feinberg, Joseph Gomes, Caleb Geniesse, Aneesh S. Pappu, Karl Leswing, and Vijay Pande. 2018. "MoleculeNet: A Benchmark for Molecular Machine Learning." *Chemical Science* 9 (2): 513–30.
- Along with these, other sites were referred for missing information and definitions that was required to carry on with the work.