

Develop predictive models that can determine given a particular compound whether it is active (1) or not (0).

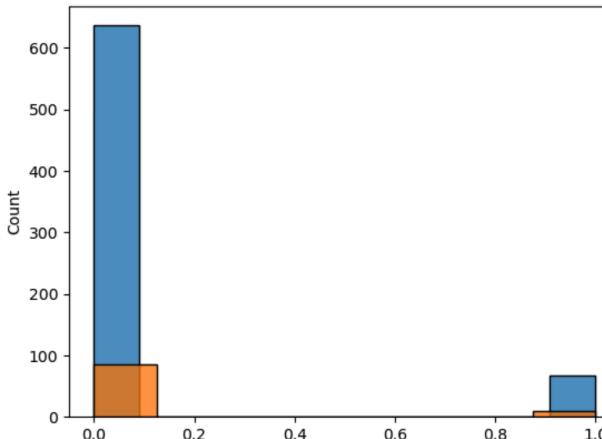
Drugs are typically small organic molecules that achieve their desired activity by binding to a target site on a receptor. The first step in the discovery of a new drug is usually to identify and isolate the receptor to which it should bind, followed by testing many small molecules for their ability to bind to the target site. This leaves researchers with the task of determining what separates the active (binding) compounds from the inactive (non-binding) ones. Such a determination can then be used in the design of new compounds that not only bind, but also have all the other properties required for a drug (solubility, oral absorption, lack of side effects, appropriate duration of action, toxicity, etc.).

The given dataset has the class labels as the first number in the rows in binary format 0 or 1. The features are represented in binary format as well, the rest of the numbers in a row represent the indices of the 1s in each record, and the indices represent the active compounds. I divided the class labels and indices into train Y and train X. The indices need to be mapped as 1 in each record in the corresponding index, in a compressed sparse row matrix for memory and speed optimization purposes. I created a class called “MatrixParser” to have two methods - The first method is called “calculateMatrixColumn”, this method will find out the max index in the records in the dataset and return the max index; the second method “parseIndexToMatrix” will take the returned value from the first method, and create an array of the max index length with only 0s for each data record and assign the indices that exist in each data record to 1 in the corresponding indices in the array. The parsed array for each data record will be appended to a “res” array and the “res” array will be transformed into a compressed sparse row matrix. Checking the sparse matrix for “res” of train X, we can see only the active compound, 1, in rows and columns. The same operation is applied to the test set.

Looking at the train X matrix and the characteristics of the dataset, the values are already normalized since they are already in binary format and they are not on different scales. However, some values can be dependent on others and redundancy can be eliminated to reduce the noise and dimensionality of the dataset for better accuracy. For my experiments with three types of classifiers, I applied “SelectKBest” with the chi2 score function or mutual information classification score function for feature selection [1]-[3]. With the chi2 test, the irrelevant features can be eliminated and I can choose k numbers of the top most relevant features, this is valid for this dataset because the features are non-negative. With mutual information classification, I can select the k numbers of the top most relevant features to the class label, this is valid for this dataset because the class label is discrete. For some experiments, I applied either chi2 SelectKBest or/and mutual information classification SelectKBest with different experimental k values to experiment with model accuracies based on F1 scores from cross-validation and from the validation set. The interpretations for experiments can be found in the last column of the tables.

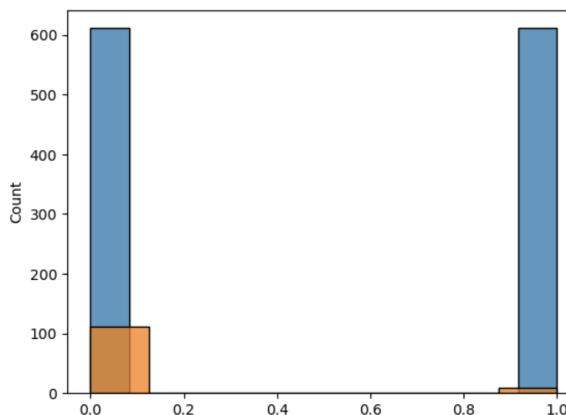
I divided train X to train X, Y, and validation X, Y. The reason is that I want to have a validation set to validate the model performance. Before dividing, train X and Y are shuffled to keep the randomness of the train set and validation set so that our model is not constantly learning the same data which could lead to overfitting. By comparing the cross-validation F1 score and validation set F1 score, I can actually see how the model performs on unseen data.

The dataset is very imbalanced as stated in the homework description and the graph drawn below. The blue represents train Y and the orange represents validation Y. From this graph, the difference between the total counts of labels 1 and 0 is really obviously significant. This is why we need to apply resampling.



I chose to apply random oversampling [4] to over-sample the underrepresented label in the train set. The reason is that the dataset is already small, it is helpful for model training and accuracy with more data, so we over-sample instead of under-sample. The new train, validation, and test ratio is roughly 72, 7, 21. See the screenshot below. Each experiment has slightly different values of train and validation class label amounts, details can be found in the screenshots column in the tables.

```
train/val/test ratio: 0.7222222222222222 0.07092198581560284 0.20685579196217493
Train label is 1: 611 Train label is 0: 611
Validation label is 1: 9 validation label is 0: 111
```



For cross-validation, I chose to use stratified k-fold cross-validation [5]. The reason is each fold can have equal instances of target class to prevent bias. For the Naive Bayes classifier, I choose to use Bernoulli Naive Bayes [6] because the dataset has discrete binary features and labels, for the same reason, Perceptron [7], [8] is chosen for the Neural Network classifier.

The random_state parameter is different for each experiment for reproducibility. The parameter is applied in the shuffle of train X and Y, train validation split, random oversampling, stratified k fold, decision tree classifier, and perceptron. The highlighted column is the final submission in Miner.

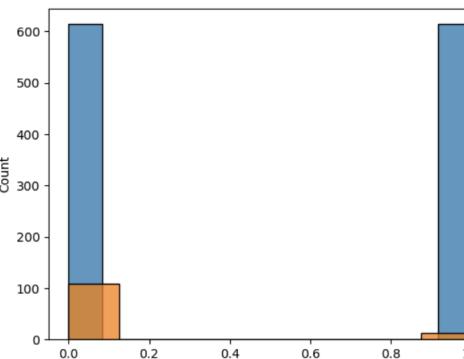
Stratified Cross Validation Experiments - Decision Tree						
random_state	params applied	CV n_splits	f1_scores	screenshots	interpretation	

4	max_dept h=8	10	average f1 score: 0.97 validation set f1 score: 0.62	<pre>train/val/test ratio: 0.7232037691401649 0.0706713780918728 0.2061248527679623 Train label is 1: 614 Train label is 0: 614 Validation label is 1: 12 validation label is 0: 108</pre> <pre> tree depth: 8 train label ratio: 0.5004524886877828 0.4995475113122172 test label ratio: 0.4959349593495935 0.5040650406504065 f1 score: 0.9838709677419354 confusion matrix: [[60 2] [0 61]] tree depth: 8 train label ratio: 0.5004524886877828 0.4995475113122172 test label ratio: 0.4959349593495935 0.5040650406504065 f1 score: 0.9354838709677418 confusion matrix: [[57 5] [3 58]] tree depth: 8 train label ratio: 0.5004524886877828 0.4995475113122172 test label ratio: 0.4959349593495935 0.5040650406504065 f1 score: 0.991869918699187 confusion matrix: [[61 1] [0 61]] tree depth: 8 train label ratio: 0.5004524886877828 0.4995475113122172 test label ratio: 0.4959349593495935 0.5040650406504065 f1 score: 0.975609756097561 confusion matrix: [[60 2] [1 60]] tree depth: 8 train label ratio: 0.4995475113122172 0.5004524886877828 test label ratio: 0.5040650406504065 0.4959349593495935 f1 score: 0.976 confusion matrix: [[59 2] [1 61]] tree depth: 8 train label ratio: 0.4995475113122172 0.5004524886877828 test label ratio: 0.5040650406504065 0.4959349593495935 f1 score: 0.96875 confusion matrix: [[57 4] [0 62]] tree depth: 8 train label ratio: 0.4995475113122172 0.5004524886877828 test label ratio: 0.5040650406504065 0.4959349593495935 f1 score: 0.9672131147540983 confusion matrix: [[60 1] [3 59]] tree depth: 8 train label ratio: 0.4995475113122172 0.5004524886877828 test label ratio: 0.5040650406504065 0.4959349593495935 f1 score: 0.9838709677419355 confusion matrix: [[60 1] [1 61]] tree depth: 8 train label ratio: 0.5 0.5 test label ratio: 0.5 0.5 f1 score: 0.9752066115702478 confusion matrix: [[60 1] [2 59]] tree depth: 8 train label ratio: 0.5 0.5 test label ratio: 0.5 0.5 f1 score: 0.9677419354838709 confusion matrix: [[58 3] [1 60]] average f1 score: 0.9677419354838709</pre>	<p>The max_depth is limited to prevent the model from learning training data too precisely.</p> <p>The validation score is closer to the average cross-validation score.</p> <p>Applying more parameters to see if we can improve accuracy.</p> <p>diff = 0.35</p>

min_samples_split and min_samples_leaf did not help with the accuracy difference. It could be because the minimum sample numbers are too big for such a small dataset.

I will lower the minimum split and leaf a bit and see what happens. I will increase the max_depth limit in case more depth is needed.

diff = 0.21

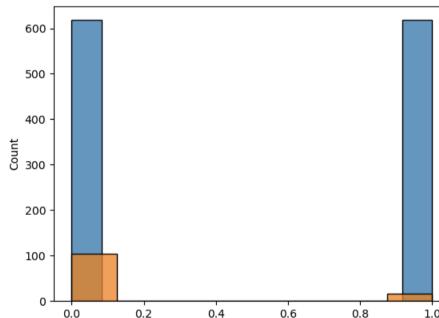
8	SelectKB est(chi2, k=500) max_dept h=15, min_sam ples_split =10, min_sam ples_leaf =4	10	average f1 score: 0.94 validation set f1 score: 0.69	<p>train/val/test ratio: 0.7232037691401649 0.0706713780918728 0.2061248527679623 Train label is 1: 614 Train label is 0: 614 Validation label is 1: 12 validation label is 0: 108</p>  <pre> tree depth: 15 train label ratio: 0.5004524886877828 0.4995475113122172 test label ratio: 0.4959349593495935 0.5040650406504065 f1 score: 0.9677419354838709 confusion matrix: [[160 1] [1 601]] tree depth: 15 train label ratio: 0.5004524886877828 0.4995475113122172 test label ratio: 0.4959349593495935 0.5040650406504065 f1 score: 0.9586774859504132 confusion matrix: [[160 2] [2 588]] tree depth: 15 train label ratio: 0.5004524886877828 0.4995475113122172 test label ratio: 0.4959349593495935 0.5040650406504065 f1 score: 0.959349593495935 confusion matrix: [[159 3] [3 59]] tree depth: 15 train label ratio: 0.5004524886877828 0.4995475113122172 test label ratio: 0.4959349593495935 0.5040650406504065 f1 score: 0.9752061135702478 confusion matrix: [[160 1] [1 601]] tree depth: 15 train label ratio: 0.4995475113122172 0.5004524886877828 test label ratio: 0.4959349593495935 0.4959349593495935 f1 score: 0.967213117540983 confusion matrix: [[160 1] [1 601]] tree depth: 15 train label ratio: 0.4995475113122172 0.5004524886877828 test label ratio: 0.4959349593495935 0.4959349593495935 f1 score: 0.975609756097561 confusion matrix: [[159 2] [2 598]] tree depth: 15 train label ratio: 0.4995475113122172 0.5004524886877828 test label ratio: 0.4959349593495935 0.4959349593495935 f1 score: 0.9682539682539683 confusion matrix: [[158 3] [3 59]] tree depth: 15 train label ratio: 0.4995475113122172 0.5004524886877828 test label ratio: 0.4959349593495935 0.4959349593495935 f1 score: 0.9677419354838709 confusion matrix: [[158 2] [2 596]] tree depth: 15 train label ratio: 0.4995475113122172 0.5004524886877828 test label ratio: 0.4959349593495935 0.4959349593495935 f1 score: 0.9411764705882353 confusion matrix: [[159 2] [5 56]] average f1 score: 0.9411764705882353 </pre> <p>f1 score validation set: 0.6896551724137931</p>
---	---	----	--	---

Stratified Cross Validation Experiments - Naive Bayes

			<pre>train/val/test ratio: 0.7238542890716804 0.07050528789659224 0.2056404230317274 Train label is 1: 616 Train label is 0: 616 Validation label is 1: 14 validation label is 0: 106</pre>	<p>let's apply both mutual info classfi select k best and chi2 select k best to eliminate dependent and irrelevant features.</p> <p>diff = 0.13</p>
--	--	--	---	---

18	force_alpha =True SelectKBes t(score_fun c=mutual_i nfo, k=600) SelectKBes t(chi2, k=300)	5	average f1 score: 0.88 validation set f1 score: 0.70	<pre>current train_x shape: (800, 100001) train_x shape after mutual info selection: (800, 600) current test_x shape: (350, 100001) test_x shape after mutual info selection: (350, 600) current train_x shape: (800, 600) train_x shape after chi2 selection: (800, 300) current test_x shape: (350, 600) test_x shape after chi2 selection: (350, 300) train/val/test ratio: 0.7241784037558685 0.07042253521126761 0.20539906103286384 Train label is 1: 617 Train label is 0: 617 Validation label is 1: 15 validation label is 0: 105</pre>	<p>The accuracy difference is bigger compared to just using mutual information classification.</p> <p>The next experiment will just use select k best for mutual information classification but make k bigger than what is in random_state 15.</p> <p>diff = 0.18</p>
----	---	---	--	--	---


```
train/val/test ratio: 0.7245017584994138 0.07033997655334115 0.20515826494724  
Train label is 1: 618 Train label is 0: 618  
Validation label is 1: 16 validation label is 0: 104
```

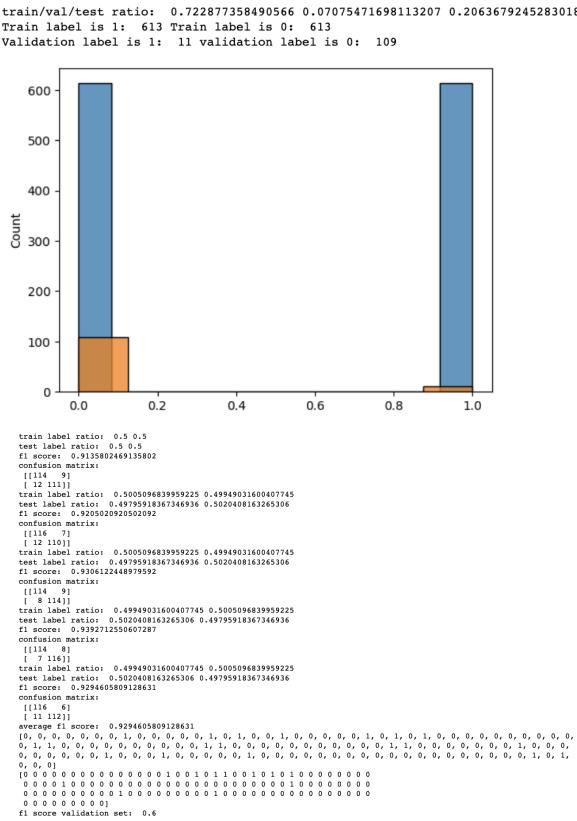


Stratified Cross Validation Experiments - Perceptron

			<pre>train/val/test ratio: 0.7225501770956316 0.07083825265643448 0.2066115702479339 Train label is 1: 612 Train label is 0: 612 Validation label is 1: 10 validation label is 0: 110</pre>	the next experiment. diff = 0.43
22	SelectKBest(score_func=mutual_info, k=500) eta0=0.1	5	<p>average f1 score: 0.96 validation set f1 score: 0.53</p> <pre>current train_x shape: (800, 100001) train_x shape after mutual info selection: (800, 500) current test_x shape: (350, 100001) test_x shape after mutual info selection: (350, 500) train/val/test ratio: 0.7238542890716804 0.07050528789659224 0.2056404230317274 Train label is 1: 616 Train label is 0: 616 Validation label is 1: 14 validation label is 0: 106</pre>	<p>apply a smaller learning rate than the default 1.</p> <p>The difference between the accuracies is still pretty big.</p> <p>Will increase the learning rate and see what happens.</p>

				<pre>train/val/test ratio: 0.7222222222222222 0.07092198581560284 0.20685579196217493 Train label is 1: 611 Train label is 0: 611 Validation label is 1: 9 validation label is 0: 111</pre>	max_iter. diff = 0.40
26	SelectKBest(score_func=mutual_info, k=500) eta0=0.5 max_iter=100	5	average f1 score: 0.93 validation set f1 score: 0.59	<pre>current train_x shape: (800, 100001) train_x shape after chi2 selection: (800, 500) current test_x shape: (350, 100001) test_x shape after chi2 selection: (350, 500)</pre> <hr/> <pre>train/val/test ratio: 0.7090301003344481 0.09587513935340022 0.19509476031215162 Train label is 1: 636 Train label is 0: 636 Validation label is 1: 86 validation label is 0: 86</pre>	The default max_iter is 1000, this experiment we set it to 100. The accuracy difference decreased a little bit but not significantly. Next will try a bigger value for max iteration. diff = 0.34

			<pre> train label ratio: 0.5004916420845624 0.4995083579154376 test label ratio: 0.4980392156862745 0.5019607843137255 f1 score: 0.9402985074626865 confusion matrix: [[113 15] [1 126]] train label ratio: 0.4995083579154376 0.5004916420845624 test label ratio: 0.5019607843137255 0.4980392156862745 f1 score: 0.9144981412639406 confusion matrix: [[109 18] [5 123]] train label ratio: 0.5 0.5 test label ratio: 0.5 0.5 f1 score: 0.9618320610687023 confusion matrix: [[118 9] [1 126]] train label ratio: 0.5 0.5 test label ratio: 0.5 0.5 f1 score: 0.9077490774907748 confusion matrix: [[106 21] [4 123]] train label ratio: 0.5 0.5 test label ratio: 0.5 0.5 f1 score: 0.9402985074626865 confusion matrix: [[112 15] [1 126]] average f1 score: 0.9402985074626865 </pre>	
27	SelectKBest(score_func=mutual_info, k=500) eta0=0.5 max_iter=500	5	<p>average f1 score: 0.94 validation set f1 score: 0.63</p> <pre> current train_x shape: (800, 100001) train_x shape after mutual info selection: (800, 500) current test_x shape: (350, 100001) test_x shape after mutual info selection: (350, 500) train/val/test ratio: 0.72387738490566 0.0705471698113207 0.206367924532830188 Train label is 1: 613 Train label is 0: 613 Validation label is 1: 11 validation label is 0: 109 Count 0.0 0.2 0.4 0.6 0.8 1.0 600 500 400 300 200 100 0 </pre> <pre> train label ratio: 0.5 0.5 test label ratio: 0.5 0.5 f1 score: 0.9402985074626865 confusion matrix: [[109 20] [1 120]] train label ratio: 0.4995083579154376 0.49949031600407745 test label ratio: 0.49795918367346936 0.5020408163265306 f1 score: 0.9618320610687023 confusion matrix: [[113 10] [1 126]] train label ratio: 0.500596839959225 0.49949031600407745 test label ratio: 0.49795918367346936 0.5020408163265306 f1 score: 0.9761165048543889 confusion matrix: [[118 11] [1 121]] train label ratio: 0.49949031600407745 0.500596839959225 test label ratio: 0.49795918367346936 0.49795918367346936 f1 score: 0.9374999999999999 confusion matrix: [[109 13] [1 120]] average f1 score: 0.9374999999999999 </pre>	<p>It looks like max_iter has minimal impact on our model. Either it is applied or not, the accuracy difference between the average cross-validation score and validation set score are all around 0.30. The problem might not be because of max epochs, it could be something else.</p> <p>Even when the learning rate is adjusted the difference between accuracies is still similar to when it is not specified.</p> <p>Next experiment we will adjust the number of k selected</p>



References

- [1] “sklearn.feature_selection.SelectKBest,” *Scikit-learn*.
https://scikit-learn.org/stable/modules/generated/sklearn.feature_selection.SelectKBest.html#sklearn.feature_selection.SelectKBest [accessed Oct. 25. 2023.]
 - [2] “Sklearn.feature_selection.chi2,” *Scikit-learn*.
https://scikit-learn.org/stable/modules/generated/sklearn.feature_selection.chi2.html#sklearn.feature_selection.chi2 [accessed Oct. 25. 2023.]
 - [3] “sklearn.feature_selection.mutual_info_classif,” *Scikit-learn*.
https://scikit-learn.org/stable/modules/generated/sklearn.feature_selection.mutual_info_classif.html#sklearn.feature_selection.mutual_info_classif [accessed Oct. 25. 2023.]

- [4] “RandomOverSampler — Version 0.11.0.”
https://imbalanced-learn.org/stable/references/generated/imblearn.over_sampling.RandomOverSampler.html#imblearn.over_sampling.RandomOverSampler
[accessed Oct. 25. 2023.]
- [5] “sklearn.model_selection.StratifiedKFold,” *Scikit-learn*.
https://scikit-learn.org/stable/modules/generated/sklearn.model_selection.StratifiedKFold.html [accessed Oct. 25. 2023.]
- [6] https://scikit-learn.org/stable/modules/generated/sklearn.naive_bayes.BernoulliNB.html [accessed Oct. 25. 2023.]
- [7] “sklearn.linear_model.Perceptron,” *Scikit-learn*.
https://scikit-learn.org/stable/modules/generated/sklearn.linear_model.Perceptron.html [accessed Oct. 28. 2023.]
- [8] J. Brownlee, “Perceptron algorithm for classification in Python,”
MachineLearningMastery.com, Aug. 05, 2020.
<https://machinelearningmastery.com/perceptron-algorithm-for-classification-in-python/> [accessed Oct. 28. 2023.]