
Graph Neural Network for chemistry

Léon Ter*

University of Nantes, France
leon.ter@etu.univ-nantes.fr

Duc Mai Chu*

University of Nantes, France
duc.chu@etu.univ-nantes.fr

1 Sharing of the work

Léon Ter: Analysis of the MoSS results, experiment GCN, training and analysis on Sirtuin6

Duc Mai Chu: Test on MoSS with the three different input files, find the Sirtuin6 dataset

2 Frequent subgraph mining for molecules

To test out the MoSS (Molecular Substructure miner) program, we run 3 different examples of input files in the SMILES format.

- example1.dat: 6 entries
- example2.dat: 3 entries
- steroids.dat: 17 entries

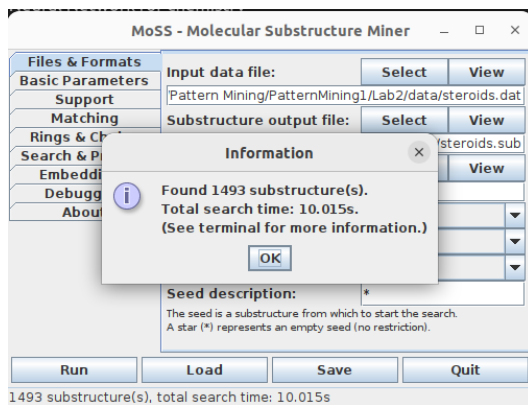


Figure 1: Information tab after finding the substructures of the steroids.dat file

*Equal contribution.

works. We then work on the Sirtuin6 Small Molecules dataset, for classification, which includes 100 molecules with descriptors to determine the candidate inhibitors of a target protein. The molecules are grouped based on low- and high-BFE which we use for the classification. It can be accessed with the following link: <https://archive.ics.uci.edu/dataset/748/Sirtuin6+small+molecules-1>

We perform evaluation on the regular measures for models with the function `classification_report` from `sklearn.metrics`. It covers precision, recall, f1-score and support for the classes.

- Precision : measures the accuracy of positive predictions.
- Recall : measures the ability of the model to find all positive samples.
- f1-score : mean of precision and recall.
- support : number of true samples in each class.

	precision	recall	f1-score	support
0	0.70	0.88	0.78	8
1	0.90	0.75	0.82	12
accuracy			0.80	20
macro avg	0.80	0.81	0.80	20
weighted avg	0.82	0.80	0.80	20

Figure 5: Result of classification report

The accuracy is evidence that the model correctly classifies with 80% of the samples.

Observing precision, recall and the `f1_score` allows us to individually assess the quality of the classification over each of the possible values predicted. We notice High_BFE (class 1) has a higher precision rate than Low_BFE (class 0) while inversely High_BFE has a lower recall rate. The precision suggests that the model incorrectly classifies low_BFE as high_BFE more often. On the contrary the recall suggests there are fewer false negatives for low_BFE.

We can also print the confusion matrix which compares the predicted labels with the true labels.

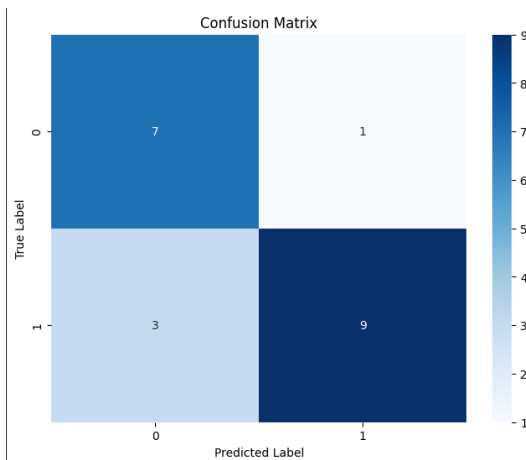


Figure 6: Result of the confusion matrix

To read the confusion matrix in this case, we could look at the false positives (1). The false positives are especially important in this context of drug molecules where we want to minimize the number of times we wrongfully classify a molecule as a class it doesn't belong.