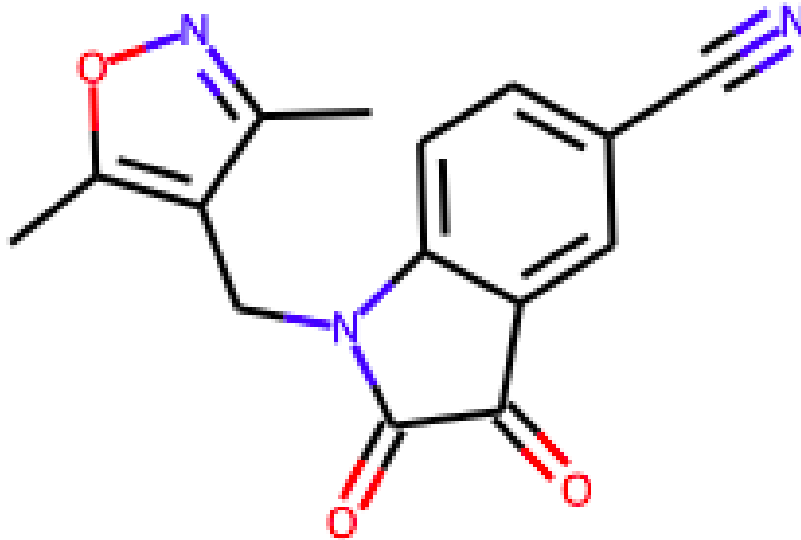


# ID2214 – Programming for Data Science

## Assignment 4a - Project report

Group 3: Francesco Gelain, Borja Javierre, Jingyi Hu



## Table of Contents

---

Objective .....	3
Methodology .....	3
Features collection .....	4
Data preparation.....	6
Modelling.....	7
Introduction.....	7
Random Forest Classifier (RFC) .....	8
Multilayer Perceptron Classifier .....	9
Multinomial Naïve Bayes .....	10
Comparison .....	11
Results.....	<b>Error! Bookmark not defined.</b>
Future work.....	13
References .....	14

## Objective

---

The objective of this project is to predict the activity of chemical compounds from a given data set. This is going to be achieved thanks to an optimized model developed in Python.

## Methodology

---

The flow diagram presented below in Figure 1 summarizes the project's methodology.

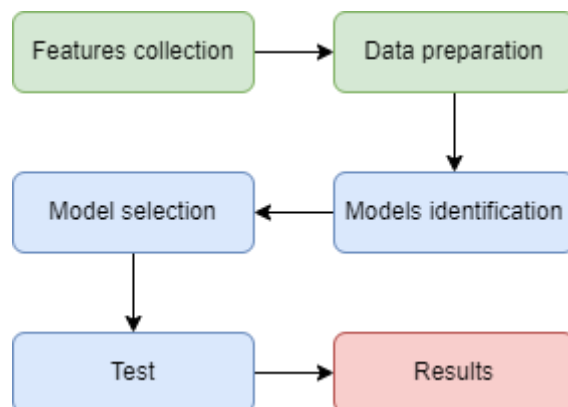


Figure 1 - Methodology

## Features collection

Features extraction and selection was achieved thanks to the open-source cheminformatics toolkit "RDkit" [1].

### Library

```
m = Chem.MolFromSmiles(df_copy['SMILES'][i]) # convert SMILES to molecule. SMILES string and m will be assigned an object representing
#chemical compound, for which various properties might be derived
numatoms = m.GetNumAtoms() # numbers of atoms
molwt = d.CalcExactMolWt(m) # molecule's exact molecular weight
heavyatom = l.HeavyAtomCount(m) # Number of heavy atoms a molecule
aroring = d.CalcNumAromaticRings(m) # number of aromatic rings for a molecule, they are very stable and do not break apart easily
amidebond = d.CalcNumAmideBonds(m) # number of amide bonds in a molecule
rotabond = d.CalcNumRotatableBonds(m) # number of rotatable bonds for a molecule
saturatedring = d.CalcNumSaturatedRings(m) # returns the number of saturated rings for a molecule
alcoo = f.fr_AL_COO(m) # Number of aliphatic carboxylic acids
benzene = f.fr_benzene(m) # Number of benzene rings
```

Figure 2.1 – "Basic" features

We selected these features because we thought they were the most representative ones for chemistry elements. Those are:

- Number of atoms (numatoms): the total number of atoms that a molecule contains
- Molecular weight (molwt): total weight of the molecule
- Heavy atoms (heavyatom): number of heavy atoms that a molecule has. These atoms are all of those which are not hydrogen atoms.
- Aromatic rings (aroring): number of aromatic rings for a molecule. They are very stable, are not easy to break and contain benzene or another related ring structure.
- Amide bonds (amidebond): number of amide bonds in a molecule. This one is syntetized when the carboxyl of the other amino acid molecule reacts with the other amino acid molecule.
- Rotatable bonds (rotabond): number of rotatable bonds for a molecule. It is any single non-rigid bond which are attached to a non-terminal, non-hydrogen atom. Amide C-N bonds are not counted because of their high barrier to rotation.
- Saturated rings (saturatedring): number of saturated rings for a molecule. A saturated molecule contains just single bonds and no rings. Also, a saturated molecule has a maximum number of hydrogen atoms possible to be an acyclic alkane.
- Aliphatic carboxylic acids (alcoo): number of aliphatic carboxylic acids. These carbon acids perform different types of industrial functions and others are important in processes of nutrition and many are intermediates in biochemical processes.
- Benzene rings (benzene): number of Benzene rings. The Benzene is an organic chemical compound composed of 6 atoms of carbon (C) and 6 of hydrogen (H).

['INDEX', 'SMILES', 'NumAtoms', 'MolWt', 'HeavyAtom', 'AroRing', 'AmideBond', 'RotatableBond', 'SaturatedRing', 'AL_COO', 'Benzene', 'ACTIVE']													
	INDEX	SMILES	NumAtoms	MolWt	HeavyAtom	AroRing	AmideBond	RotatableBond	SaturatedRing	AL_COO	Benzene	ACTIVE	
	0	1	CC(C)N1CC(=O)C(c2nc3cccc3[nH]2)=C1N	19.0	256.132411	19.0	2.0	0.0	2.0	0.0	0.0	1.0	0.0
	1	2	Cc1ccc(-c2ccc3c(N)c(C(=O)c4ccc(OC)c(OC)c4)sc3...	30.0	420.114378	30.0	4.0	0.0	6.0	0.0	0.0	2.0	0.0
	2	3	Cc1ccc(C(=O)COC(=O)CCc2nc(=O)c3cccc3[nH]2)cc1	27.0	364.142307	27.0	3.0	0.0	7.0	0.0	0.0	2.0	0.0
	3	4	O=C(CN1CCOCC1)Nc1ccc(S(=O)(=O)N2CCCCC2)cc1	26.0	381.172227	26.0	1.0	1.0	5.0	2.0	0.0	1.0	0.0
	4	5	C=CC(Nc1cccc1)c1ccc(OC)c(OC)c1	21.0	283.157229	21.0	2.0	0.0	7.0	0.0	0.0	2.0	0.0
	...	...	...	...	...	...	...	...	...	...	...	...	...
156253	156254	O=C(N/N=C/c1ccc(F)cc1)c1ccc(Cn2cc(Br)c([N+])(=O)...	27.0	434.997844	27.0	3.0	1.0	6.0	0.0	0.0	0.0	1.0	0.0
156254	156255	COc1ccc(NS(=O)(=O)c2cc(NC(=O)C=C/c3cc(OC)ccc3...	39.0	551.209007	39.0	3.0	1.0	10.0	1.0	0.0	3.0	0.0	0.0
156255	156256	O=c1nc(N2CCN(Cc3cccc3)CC2)[nH]c2c1CCC2	23.0	310.179361	23.0	2.0	0.0	3.0	1.0	0.0	1.0	0.0	0.0
156256	156257	Cc1onc(-c2cccc2)c1C(=O)N/N=C/c1ccco1	22.0	295.095691	22.0	3.0	1.0	4.0	0.0	0.0	0.0	1.0	0.0
156257	156258	CC(C)CCNC(=O)c1ccc2nc[nH]c2c1	17.0	231.137162	17.0	2.0	1.0	4.0	0.0	0.0	0.0	1.0	0.0

156258 rows x 12 columns

Figure 2.2: Result for the features collection using original library features

## Fingerprint

A way to vectorizing the features that represents the molecules. By applying fingerprints to the DataFrame we can know and select which is the best model according to its performance. The fingerprints are based on the Morgan algorithm, so they are similar to the ECFP and FCFP fingerprints [10].

	INDEX	ACTIVE	rd_form	fcpc0	fcpc1	fcpc2	fcpc3	fcpc4	fcpc5	fcpc6	...	fcpc114	fcpc115	fcpc116	fcpc117	fcpc118	fcpc119	fcpc120	fcpc121	fcpc122	fcpc123
0	1	0.0	<rdkit.Chem.rdchem.Mol object at 0x121707ca0>	1	0	1	0	1	1	1	...	0	0	1	1	0	0	0	0	1	0
1	2	0.0	<rdkit.Chem.rdchem.Mol object at 0x121707d10>	1	1	1	1	1	0	1	...	1	0	0	0	0	0	0	0	0	0
2	3	0.0	<rdkit.Chem.rdchem.Mol object at 0x121707d80>	1	0	1	0	1	1	1	...	0	0	0	0	0	0	1	0	0	0
3	4	0.0	<rdkit.Chem.rdchem.Mol object at 0x121707df0>	1	1	1	1	1	0	0	...	1	0	0	0	0	1	0	0	0	1
4	5	0.0	<rdkit.Chem.rdchem.Mol object at 0x121707e60>	1	0	1	1	1	0	0	...	0	0	0	0	1	0	0	0	0	0
...	...	...	...	...	...	...	...	...	...	...	...	...	...	...	...	...	...	...	...	...	...
156253	156254	0.0	<rdkit.Chem.rdchem.Mol object at 0x122816110>	1	1	1	1	1	0	1	...	0	0	1	0	1	0	0	0	0	1
156254	156255	0.0	<rdkit.Chem.rdchem.Mol object at 0x122816180>	1	1	1	1	1	1	0	...	1	0	0	0	0	0	0	1	1	1
156255	156256	0.0	<rdkit.Chem.rdchem.Mol object at 0x1228161f0>	1	0	1	0	1	1	1	...	0	0	1	0	0	0	0	1	0	0
156256	156257	0.0	<rdkit.Chem.rdchem.Mol object at 0x122816260>	1	1	1	0	1	0	1	...	0	0	0	0	0	0	0	0	0	1
156257	156258	0.0	<rdkit.Chem.rdchem.Mol object at 0x1228162d0>	1	1	1	0	1	1	1	...	0	0	0	0	0	0	0	0	1	1

156258 rows x 127 columns

## Library + Fingerprint

Combination of both the original “basic” library of features and the fingerprint features. We get a bigger table as we add the columns representing the fingerprints together with the columns of the features from the original library in the same DataFrame.

	INDEX	NumAtoms	MolWt	HeavyAtom	AroRing	AmideBond	RotatableBond	SaturatedRing	AL_COO	Benzene	...	fcpc114	fcpc115	fcpc116	fcpc117	fcpc118	fcpc119	fcpc120	fcpc121	fcpc122	fcpc123
0	1	19.0	256.132411	19.0	2.0	0.0	2.0	0.0	0.0	1.0	...	0	0	1	1	0	0	0	0	1	0
1	2	30.0	420.114378	30.0	4.0	0.0	6.0	0.0	0.0	2.0	...	1	0	0	0	0	0	0	0	0	0
2	3	27.0	364.142307	27.0	3.0	0.0	7.0	0.0	0.0	2.0	...	0	0	0	0	0	0	1	0	0	0
3	4	26.0	381.172227	26.0	1.0	1.0	5.0	2.0	0.0	1.0	...	1	0	0	0	0	1	0	0	0	1
4	5	21.0	283.157229	21.0	2.0	0.0	7.0	0.0	0.0	2.0	...	0	0	0	0	1	0	0	0	0	0
...	...	...	...	...	...	...	...	...	...	...	...	...	...	...	...	...	...	...	...	...	...
156253	156254	27.0	434.997844	27.0	3.0	1.0	6.0	0.0	0.0	1.0	...	0	0	1	0	1	0	0	0	0	1
156254	156255	39.0	551.209007	39.0	3.0	1.0	10.0	1.0	0.0	3.0	...	1	0	0	0	0	0	1	1	1	1
156255	156256	23.0	310.179361	23.0	2.0	0.0	3.0	1.0	0.0	1.0	...	0	0	1	0	0	0	1	0	0	0
156256	156257	22.0	295.095691	22.0	3.0	1.0	4.0	0.0	0.0	1.0	...	0	0	0	0	0	0	0	0	0	1
156257	156258	17.0	231.137162	17.0	2.0	1.0	4.0	0.0	0.0	1.0	...	0	0	0	0	0	0	0	0	1	1

156258 rows x 136 columns

## Data preparation

The open-source machine learning library Scikit-learn [2] was exploited to prepare the data. Firstly, the percentage of active substances in the test set was evaluated.

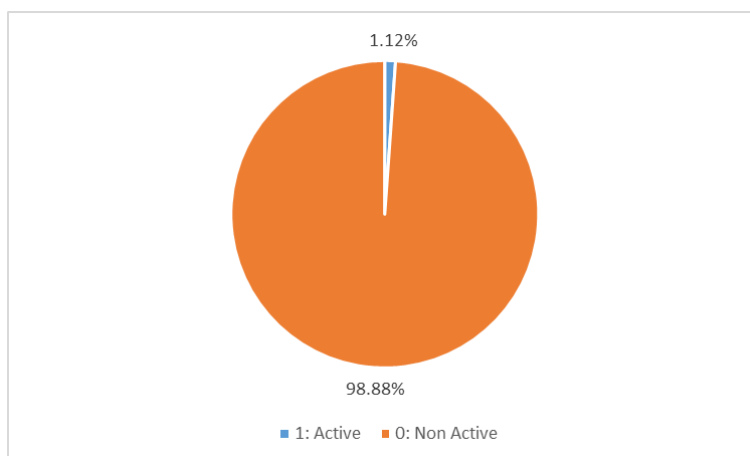


Figure 3 – Percentage of active substances in the test set

Afterwards, the columns "INDEX", "SMILES", "HeavyAtom", "ACTIVE" (stored first) were dropped. "HeavyAtom" was dropped since it was containing the same values as the column "NumAtoms".

The remaining data frame was normalized thanks to the functions shown in Figure 4.

```
imp_mean = SimpleImputer(missing_values=np.nan, strategy='mean')
scaler = MinMaxScaler()
kbd = KBinsDiscretizer(n_bins=10, encode="ordinal")
```

Figure 4

NumAtoms	MolWt	HeavyAtom	AroRing	AmideBond	RotatableBond	SaturatedRing	AL_COO	Benzene	ACTIVE
19.0	256.132411	19.0	2.0	0.0	2.0	0.0	0.0	1.0	0.0
30.0	420.114378	30.0	4.0	0.0	6.0	0.0	0.0	2.0	0.0
27.0	364.142307	27.0	3.0	0.0	7.0	0.0	0.0	2.0	0.0
26.0	381.172227	26.0	1.0	1.0	5.0	2.0	0.0	1.0	0.0
21.0	283.157229	21.0	2.0	0.0	7.0	0.0	0.0	2.0	0.0
...	...	...	...	...	...	...	...	...	...

Figure 5 - Before data preparation

	NumAtoms	MolWt	AroRing	AmideBond	RotatableBond	SaturatedRing	AL_COO	Benzene
0	1.0	1.0	2.0	0.0	1.0	0.0	0.0	1.0
1	8.0	8.0	3.0	0.0	5.0	0.0	0.0	2.0
2	7.0	5.0	3.0	0.0	6.0	0.0	0.0	2.0
3	6.0	6.0	1.0	1.0	4.0	1.0	0.0	1.0
4	3.0	2.0	2.0	0.0	6.0	0.0	0.0	2.0
...	...	...	...	...	...	...	...	...

Figure 6 - After data preparation

## Modelling

---

### Introduction

All the models were compared based on the AUC score. The data separation was done using stratified k-fold technique, on the training set that was given, and the optimization of the hyperparameters for each model was done thanks to GridSearchCV. This function pairs the input hyper-parameters for classifier, generates and learns the input data separated by given cross validation method. Then based on the auc score for each paired parameters returns the model and the best parameters with the highest performance. Then, according to the comparison between the predicted output and the ACTIVE column for the validation set, roc\_auc\_score is used to calculate the auc score.

Function in code: `grid = GridSearchCV(clf, param_grid=parameters, cv=skf, scoring="roc_auc")`

`auc = metrics.roc_auc_score(y_val, prediction[:, 1])`

`grid.best_params_`

`grid.best_estimator`

Stratified k-fold ( $k = 5$ ) cross-validation has been implemented to reduce the risk of overfitting. The model is trained using  $k-1$  of the folds and validated on the remaining fold. Stratified is there to ensure that each fold of dataset has the same proportion of observations for a given label, this was done since the class distribution was extremely uneven, as depicted in Figure 3. A graphical explanation is presented below [3].

Function in code: `skf = StratifiedKFold(n_splits=5)`

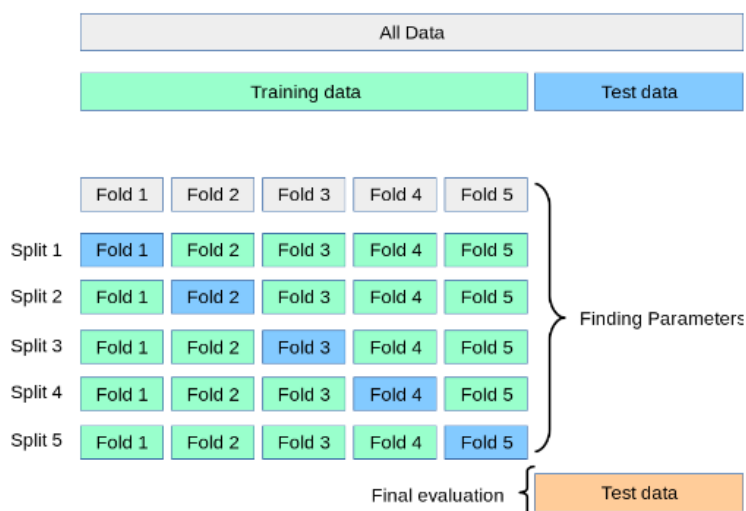


Figure 7 - k-fold cross-validation

Supervised Learning: For this project of SMILES activity prediction model training, supervised learning methods are suitable for it. Because the desired output labels for activities of compounds are already given on the training sets. Supervised learning algorithms can be trained by mapping the pairs of input compound features vectors with the supervisory active signals. Therefore, we use random forest classifier, multilayer perceptron classifier, multinomial Naïve Bayes and Bernoulli Naïve Bayes for model training.

### Random Forest Classifier (RFC)

RFC offers very good predicting capabilities together with a lower chance of overfitting, robustness to outliers and runs efficiently on large datasets, as the one that we are considering [4], [5].

Function in code: `model = RandomForestClassifier (**best_params_grid)`

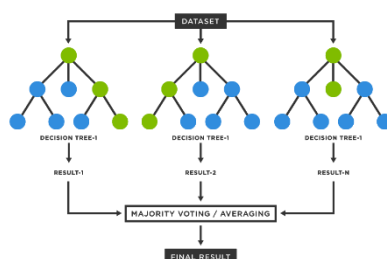


Figure 8 – RFC

Below it is the model selection with basic and fingerprint features for different parameter pairs based on GridSearchCV method.



m_state	params	split0_test_score	split1_test_score	split2_test_score	split3_test_score	split4_test_score	mean_test_score	std_test_score	rank_test_score
0	{'max_depth': 10, 'n_estimators': 400, 'random...'	0.772856	0.759384	0.746942	0.754230	0.759567	0.758596	0.008483	1
0	{'max_depth': 10, 'n_estimators': 200, 'random...'	0.773608	0.758398	0.746875	0.753095	0.760625	0.758520	0.008912	2
0	{'max_depth': 10, 'n_estimators': 300, 'random...'	0.772830	0.760018	0.746288	0.753745	0.759678	0.758512	0.008727	3
0	{'max_depth': 5, 'n_estimators': 400, 'random...'	0.764216	0.767899	0.742286	0.738539	0.746499	0.751888	0.011897	4
0	{'max_depth': 5, 'n_estimators': 300, 'random...'	0.763862	0.767853	0.742509	0.738255	0.746135	0.751723	0.011875	5
0	{'max_depth': 5, 'n_estimators': 200, 'random...'	0.763775	0.767773	0.740932	0.737815	0.745577	0.751174	0.012239	6
0	{'max_depth': 15, 'n_estimators': 400, 'random...'	0.688763	0.674773	0.685518	0.704456	0.689147	0.688531	0.009512	7
0	{'max_depth': 15, 'n_estimators': 300, 'random...'	0.688856	0.671688	0.684814	0.706097	0.686368	0.687564	0.011011	8
0	{'max_depth': 15, 'n_estimators': 200, 'random...'	0.682612	0.670549	0.681859	0.706535	0.685447	0.685401	0.011727	9

## Multilayer Perceptron Classifier (MLP)

It is a feedforward supervised artificial neural network (ANN) model with backpropagation techniques. The layers can be separated into three main parts: input layer, hidden layers and output layer.

Function in code: `mlp = MLPClassifier(**best_params_grid)`

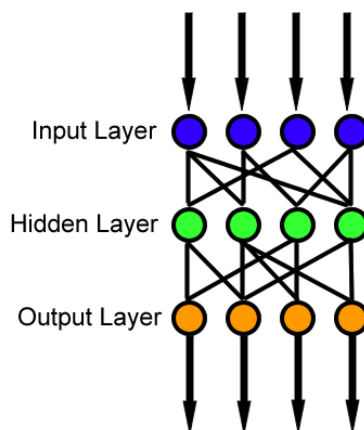


Figure 9

Below it is the model selection with basic and fingerprint features for different parameter pairs based on GridSearchCV method.

params	split0_test_score	split1_test_score	split2_test_score	split3_test_score	split4_test_score	mean_test_score	std_test_score	rank_test_score
{'activation': 'tanh', 'hidden_layer_sizes': (30,), 'learning_rate': 'adaptive', 'max_iter': 50, 'solver': 'adam'}	0.771860712082769	0.7771516483508260	0.7578901846322	0.7547069005987960	0.7683707202661900	0.7659960331861560	0.008457312130285240	1
{'activation': 'tanh', 'hidden_layer_sizes': (30,), 'learning_rate': 'constant', 'max_iter': 20, 'solver': 'adam'}	0.7717129584430730	0.7757098346692880	0.7554697368750820	0.7529639123437200	0.7698106167230740	0.7651334118108470	0.009148887699624510	2
{'activation': 'tanh', 'hidden_layer_sizes': (30,), 'learning_rate': 'constant', 'max_iter': 60, 'solver': 'adam'}	0.7695914780189340	0.772209307357120	0.7567343307039470	0.7538355935072330	0.767140716917310	0.7649006227315110	0.008572549421381560	3
{'activation': 'tanh', 'hidden_layer_sizes': (30,), 'learning_rate': 'adaptive', 'max_iter': 60, 'solver': 'adam'}	0.7714153003195090	0.775294939999760	0.7565375842242400	0.7549723981659250	0.7655774314910600	0.764854423580154	0.00811823720498160	4
{'activation': 'tanh', 'hidden_layer_sizes': (25,), 'learning_rate': 'adaptive', 'max_iter': 30, 'solver': 'adam'}	0.7730880959566250	0.7750239959930260	0.7557063297285190	0.7529034354767070	0.7675138149447360	0.7648271344092160	0.00905436633296300	5
{'activation': 'tanh', 'hidden_layer_sizes': (30,), 'learning_rate': 'adaptive', 'max_iter': 40, 'solver': 'adam'}	0.7695147957502310	0.7771494399917420	0.7575813982282280	0.7531406613405060	0.7665343945434880	0.7647821307890790	0.008552368393840040	6
{'activation': 'tanh', 'hidden_layer_sizes': (30,), 'learning_rate': 'constant', 'max_iter': 50, 'solver': 'adam'}	0.76912240696702810	0.7771886217009030	0.7562971572573780	0.7516017293346290	0.7693921237280390	0.7647202559962660	0.00937762636487070	7
{'activation': 'relu', 'hidden_layer_sizes': (30,), 'learning_rate': 'constant', 'max_iter': 60, 'solver': 'adam'}	0.7715038589884150	0.772458064702850	0.7534171769778130	0.7568162923297480	0.7690215176473120	0.76470347048209150	0.007968181566206880	8
{'activation': 'tanh', 'hidden_layer_sizes': (30,), 'learning_rate': 'adaptive', 'max_iter': 30, 'solver': 'adam'}	0.7727137555646030	0.77597831812471	0.7548349288108000	0.753079968166477	0.7668185087050570	0.7648450964745300	0.009297574923785010	9
{'activation': 'tanh', 'hidden_layer_sizes': (30,), 'learning_rate': 'constant', 'max_iter': 30, 'solver': 'adam'}	0.7705502624216860	0.77452908276149	0.757853908270060	0.7532338182202970	0.7671170181228510	0.7645568015739980	0.007865431659873910	10
{'activation': 'tanh', 'hidden_layer_sizes': (30,), 'learning_rate': 'adaptive', 'max_iter': 20, 'solver': 'adam'}	0.7708625333801660	0.7753116479618040	0.7541055406118650	0.754444395807273	0.7679628883217540	0.7645374071765720	0.00870078306392276	11
{'activation': 'relu', 'hidden_layer_sizes': (30,), 'learning_rate': 'constant', 'max_iter': 30, 'solver': 'adam'}	0.7729704541348390	0.7767391538541820	0.752786512075213	0.753078658914649	0.7668742454267370	0.7644888048811240	0.0095948214246918200	12
{'activation': 'tanh', 'hidden_layer_sizes': (25,), 'learning_rate': 'adaptive', 'max_iter': 60, 'solver': 'adam'}	0.7701722994468780	0.7754925994129510	0.7562631178112710	0.7531747018880350	0.7672528062410160	0.7644711049600300	0.008446236984224820	13
{'activation': 'tanh', 'hidden_layer_sizes': (25,), 'learning_rate': 'constant', 'max_iter': 50, 'solver': 'adam'}	0.7704860421438270	0.7757058421905090	0.7546682201372130	0.7528674017803950	0.7685272693776290	0.7644511551259140	0.009051123892324890	14
{'activation': 'tanh', 'hidden_layer_sizes': (25,), 'learning_rate': 'constant', 'max_iter': 20, 'solver': 'adam'}	0.7707137510758900	0.7756349391848800	0.7549612025126350	0.7521363716704090	0.768432726920460	0.7643757954271680	0.009185145477946170	15
{'activation': 'relu', 'hidden_layer_sizes': (30,), 'learning_rate': 'adaptive', 'max_iter': 30, 'solver': 'adam'}	0.7727368537601800	0.77273223739408963	0.7534846947802570	0.7549062809486090	0.7675013770523690	0.7643705891900760	0.0085586250245315240	16
{'activation': 'tanh', 'hidden_layer_sizes': (25,), 'learning_rate': 'adaptive', 'max_iter': 40, 'solver': 'adam'}	0.7699057818056540	0.776527248900680	0.7570595847412010	0.751846839080436	0.7664814698783240	0.7643641850712570	0.008867620045537820	17
{'activation': 'tanh', 'hidden_layer_sizes': (30,), 'learning_rate': 'constant', 'max_iter': 40, 'solver': 'adam'}	0.7679158769328140	0.7761698347665440	0.7555520705593580	0.7533864615879540	0.768480228298070	0.7643615347352950	0.008496283712288840	18
{'activation': 'tanh', 'hidden_layer_sizes': (25,), 'learning_rate': 'adaptive', 'max_iter': 20, 'solver': 'adam'}	0.7716110271346750	0.7758255126770750	0.7533014064550640	0.7528402815639570	0.7681244874045300	0.7643405428470600	0.009520517571288840	19
{'activation': 'tanh', 'hidden_layer_sizes': (30,), 'learning_rate': 'adaptive', 'max_iter': 60, 'solver': 'adam'}	0.7717546681160510	0.7739223896888160	0.75484747946879010	0.732128422641453	0.7690045851869380	0.7643369716642320	0.009025192607993520	20

## Multinomial Naïve Bayes

It is a kind of Naive Bayes classifier which is working on discrete numeric features for classification. It would be better to feed integer type of features into this classifier. For other Naïve Bayes methods like BernoulliNB for binary features, multinomial is more suitable for this project.

Function in code: model = MultinomialNB(\*\*best\_params\_grid)

Below it is the model selection with basic and fingerprint features for different parameter pairs based on GridSearchCV method.

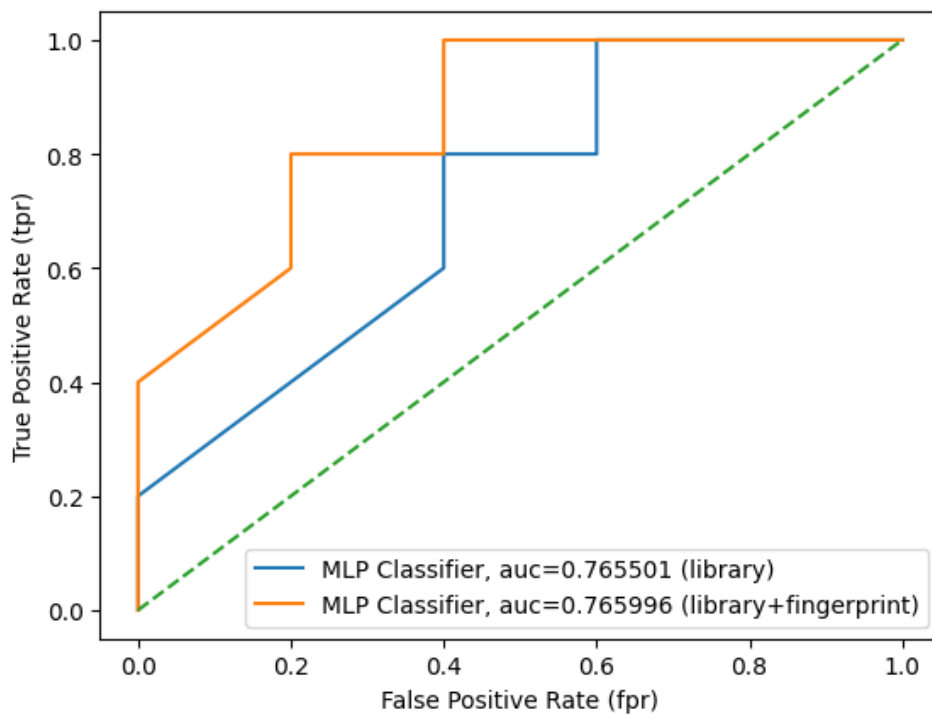
params	split0_test_score	split1_test_score	split2_test_score	split3_test_score	split4_test_score	mean_test_score	std_test_score	rank_test_score
{'alpha': 0, 'fit_prior': True}	0.752244	0.758346	0.737104	0.743091	0.756964	0.749550	0.008200	1
{'alpha': 0, 'fit_prior': False}	0.752244	0.758346	0.737104	0.743091	0.756964	0.749550	0.008200	1
{'alpha': 1, 'fit_prior': True}	0.747372	0.754115	0.733785	0.739245	0.753913	0.745686	0.008058	3
{'alpha': 1, 'fit_prior': False}	0.747372	0.754115	0.733785	0.739245	0.753913	0.745686	0.008058	3

## Comparison and results

Based on the auc scores for different models with different input features, we found out that the fingerprint features help to improve the accuracies of the trained models but not a lot. The Naïve Bayes approaches auc scores don't change based on the input features, which is probably due to the priorities ranking of fingerprint is very low. The highest auc score in this project is 0.766, which is produced from MLP with parameters: activation = 'tanh', hidden\_layer\_sizes = (30,), learning\_rate = 'adaptive', max\_iter = 50).

	Basic Information	Basic Information + Fingerprint
Random Forest Classifier	0.758384 { 'max_depth': 10, 'n_estimators': 400, 'random_state': 0 }	0.758596 { 'max_depth': 10, 'n_estimators': 400, 'random_state': 0 }
Multilayer Perceptron Classifier	0.765501 { activation='tanh', hidden_layer_sizes=(25,), learning_rate='adaptive', max_iter=40 }	<b>0.765996</b> <b>(activation='tanh', hidden_layer_sizes=(30,), learning_rate='adaptive', max_iter=50)</b>
Multinomial Naïve Bayes	0.749550 { alpha=0, fitprior=True }	0.749550 { alpha=0, fitprior=True }
BernoulliNB Naïve Bayes	0.660803 { alpha=2, fitprior=False, binarize = 2 }	0.660803 { alpha=2, fitprior=False, Binarize = 2 }

## Results



The plot shows the comparison between both Multilayer Perceptron classifiers when using just the original library and the combination of both library and fingerprints. As this is the best model classifier we tested and seen performs with the best results, we decided to plot just this one.

## Future work

---

Different paths can be taken toward future work. One would be to increase the number of features, extract and use more such that can provide even more useful information about chemical components and molecules.

When it comes to the results obtained from the current models used, we know that the use and combination of other models will give different results. A good approach would be to combine the outputs of the models as we did with combining the original features and the fingerprints features, such that we could see if there is an improvement on that.

Finally, a better representation of the data could be achieved, looking for new ways that the results could be represented more graphically and easier to understand.

## References

---

- [1]. “RDKit.” <https://www.rdkit.org/> (accessed Dec. 14, 2022).
- [2]. “scikit-learn: machine learning in Python — scikit-learn 1.2.0 documentation.” <https://scikit-learn.org/stable/> (accessed Dec. 14, 2022).
- [3]. “3.1. Cross-validation: evaluating estimator performance — scikit-learn 1.2.0 documentation.” [https://scikit-learn.org/stable/modules/cross\\_validation.html](https://scikit-learn.org/stable/modules/cross_validation.html) (accessed Dec. 14, 2022).
- [4]. “Why Choose Random Forest and Not Decision Trees – Towards AI.” <https://towardsai.net/p/machine-learning/why-choose-random-forest-and-not-decision-trees> (accessed Dec. 14, 2022).
- [5]. “sklearn.ensemble.RandomForestClassifier — scikit-learn 1.2.0 documentation.” <https://scikit-learn.org/stable/modules/generated/sklearn.ensemble.RandomForestClassifier.html> (accessed Dec. 14, 2022).
- [6]. “Scikit-learn: sklearn.naive\_bayes.MultinomialNB” [https://scikit-learn.org/stable/modules/generated/sklearn.naive\\_bayes.MultinomialNB.html#sklearn.naive\\_bayes.MultinomialNB](https://scikit-learn.org/stable/modules/generated/sklearn.naive_bayes.MultinomialNB.html#sklearn.naive_bayes.MultinomialNB) (accessed Dec 15, 2022).
- [7]. “Angelo.edu: Aromatic Rings” [https://www.angelo.edu/faculty/kboudrea/molecule\\_gallery/04\\_aromatics/00\\_aromatics.htm#:~:text=Aromatic%20rings%20\(also%20known%20as,picture%20has%20some%20complications%2C%20however](https://www.angelo.edu/faculty/kboudrea/molecule_gallery/04_aromatics/00_aromatics.htm#:~:text=Aromatic%20rings%20(also%20known%20as,picture%20has%20some%20complications%2C%20however) (accessed Dec 15, 2022).
- [8]. “cham.libretexts.com: calculating degrees of unsaturation (DoU)” [https://chem.libretexts.org/Bookshelves/Organic\\_Chemistry/Supplemental\\_Modules\\_\(Organic\\_Chemistry\)/Alkenes/Properties\\_of\\_Alkenes/Degree\\_of\\_Unsaturation#:~:text=As%20stated%20before%2C%20a%20saturated,to%20be%20an%20acyclic%20alkane](https://chem.libretexts.org/Bookshelves/Organic_Chemistry/Supplemental_Modules_(Organic_Chemistry)/Alkenes/Properties_of_Alkenes/Degree_of_Unsaturation#:~:text=As%20stated%20before%2C%20a%20saturated,to%20be%20an%20acyclic%20alkane) (accessed Dec. 15 2022).
- [9]. “Wiley Online Library: aliphatic carboxylic acids” <https://onlinelibrary.wiley.com/doi/abs/10.1002/0471435139.tox070#:~:text=Aliphatic%20carboxylic%20acids%20include%20a,intermediates%20in%20normal%20biochemical%20processes>. (accessed Dec. 15, 2022)
- [10]. “sp8rks: RDKit Tutorial” [https://github.com/sp8rks/MaterialsInformatics/blob/main/worked\\_examples/RDKit\\_tutorial/RDkit\\_tutorial.ipynb](https://github.com/sp8rks/MaterialsInformatics/blob/main/worked_examples/RDKit_tutorial/RDkit_tutorial.ipynb) (accessed Dec. 14, 2022)

## Appendix:

---

Used libraries:

# basic

import numpy as np

import pandas as pd

import time

import itertools

# sklearn

import sklearn

from sklearn.tree import DecisionTreeClassifier

from sklearn.model\_selection import train\_test\_split

from sklearn.impute import SimpleImputer

from sklearn.preprocessing import MinMaxScaler

from sklearn.preprocessing import KBinsDiscretizer

from sklearn.ensemble import RandomForestClassifier

from sklearn.model\_selection import KFold, cross\_val\_score

from sklearn.neural\_network import MLPClassifier

from sklearn import metrics

from sklearn.model\_selection import GridSearchCV

from sklearn.naive\_bayes import MultinomialNB

from sklearn.naive\_bayes import BernoulliNB

# rdkit

from rdkit import Chem

import rdkit.Chem.rdMolDescriptors as d

import rdkit.Chem.Fragments as f

import rdkit.Chem.Lipinski as l

from rdkit.Chem import AllChem