

CONTENTS

- 1 SOURCE
- 2 ABOUT PROJECT
- 3 IMPORTING LIBRARIES
- 4 IMPORTING DATA
 - 5.1 DATA DESCRIPTION
 - 5.2 INFORMATION
 - 5.2.1 IMPORTING UNIQUE COLUMNS
- 6 EXPLORATORY ANALYSIS
 - 6.1 CORRELATION ANALYSIS
 - 6.2 OUTLIERS
 - 6.2.1 DETECTING OUTLIERS
 - 6.2.2 TREATING OUTLIERS
- 7 PREDICTIVE ANALYSING
 - 7.1 RANDOM SAMPLING
 - 7.2 FEATURE SCALING
 - 7.3 LOGISTIC REGRESSION
 - 7.4 DECISION TREE
 - 7.5 RANDOM FOREST
 - 7.6 GRADIENT BOOSTING
 - 7.7 EXTREME GRADIENT BOOSTING

SOURCE

The following data set was extracted from the website : <https://www.brenda-enzymes.org/index.php>

ABOUT PROJECT

The objective of this classification project is to build predictive models using Logistic Regression, Decision Tree, and Ensemble Techniques to map the relationships between enzymes and their possible substrates. Enzymes, as catalysts, facilitate one or more specific reactions, and understanding the associations between enzymes and their potential substrates holds significant utility in various domains, including bioinformatics, drug discovery, and biotechnology.

IMPORTING LIBRARIES

```
In [1]: import warnings

import numpy as np

import pandas as pd

import seaborn as sns

import xgboost as XGB

import notebook_as_pdf

from scipy import stats

import ydata_profiling as p

import matplotlib.pyplot as plt

warnings.filterwarnings('ignore')

pd.set_option('display.max_columns', None)

from sklearn.tree import DecisionTreeClassifier

from sklearn.linear_model import LogisticRegression

from imblearn.under_sampling import RandomUnderSampler

from sklearn.model_selection import GridSearchCV, RFoid

from sklearn.preprocessing import LabelBinarizer, MinMaxScaler, Normalizer, StandardScaler

from sklearn.metrics import accuracy_score, auc, precision_score, confusion_matrix, f1_score, recall_score

from sklearn.ensemble import RandomForestClassifier, AdaBoostClassifier, GradientBoostingClassifier, ExtraTreeClassifier
```

Out[3]:

	id	BertzCT	Ch1	Ch1n	Ch1v	Ch2n	Ch2v	Ch3v	Ch4n	EState_VSA1	EState_VSA2	ExactMolWt	FpD
0	0	323.390782	9.877918	5.875576	5.875576	4.304757	4.304757	2.754513	1.749203	0.000000	0.000000	11.938294	222.068080
1	1	273.723798	1.259037	4.441467	1.874215	5.297097	5.297097	3.924155	2.569694	0.000000	0.000000	29.783175	315.210331
2	2	521.643822	1.0911303	8.527859	11.050864	6.665291	9.519706	5.824822	1.770579	15.645394	6.606882	382.121027	
3	3	567.431166	12.453343	7.089119	12.833709	6.478023	10.978151	7.914542	3.067181	95.639554	0.000000	530.070277	
4	4	112.770735	4.414719	2.866236	2.866236	1.875634	1.875634	1.036450	0.727664	17.980451	12.841643	118.062994	
...
14833	14833	632.207041	1.0911303	6.579933	9.179964	6.653983	6.030052	3.670528	1.770579	32.971529	6.606882	347.063084	
14834	14834	62.568425	2.6437384	1.446898	1.446898	0.879497	0.879497	0.174620	0.000000	0.000000	0.000000	74.024203	
14835	14835	981.327476	1.0911303	6.146219	6.146219	4.700576	4.700576	3.064846	2.13897	17.248535	0.000000	297.089560	
14836	14836	299.171248	9.949161	6.589761	7.848913	5.276568	5.476436	3.978973	2.299833	45.623794	0.000000	265.959270	
14837	14837	785.394062	15.671142	9.896164	10.234264	7.860296	8.522605	5.645502	3.312893	82.448246	5.687386	437.234828	

14838 rows x 38 columns

```
In [4]: DF_Test
```

Out[4]:

	id	BertzCT	Ch1	Ch1n	Ch1v	Ch2n	Ch2v	Ch3v	Ch4n	EState_VSA1	EState_VSA2	ExactMolWt	FpD
0	14838	344.632371	7.283603	4.473966	5.834958	5.412257	4.651530	2.096558	1.116433	49.458581	0.000000	204.077907	
1	14839	1432.430201	10.663869	7.079026	8.065215	5.297097	5.297097	3.924155	2.569694	0.000000	29.783175	315.210331	
2	14840	83.352608	3.931852	1.774415	1.073446	1.073446	4.050440	4.932842	4.392859	5.969305	6.420822	101.047679	
3	14841	150.255172	5.917290	3.548812	3.548812	2.595128	2.595128	1.642813	0.694113	0.000000	0.000000	187.115661	
4	14842	1817.276351	24.910940	15.540529	20.047314	12.535886	17.730988	11.979518	4.431173	84.554972	47.362026	891.167638	
...
9888	24726	246.422865	4.038581	2.816709	2.816709	1.875634	1.875634	1.225986	0.362743	24.146405	6.420822	146.105528	
9889	24727	591.069706	6.770857	5.682461	5.682461	4.050440	4.050440	2.167855	1.770579	5.969305	6.420822	304.240230	
9890	24728	378.113435	6.310349	3.402334	4.317724	2.817428	4.071978	1.970236	1.165747	36.705949	0.000000	260.029179	
9891	24729	737.653158	9.949161	7.337949	7.337949	4.428511	5.948361	3.974259	2.160881	36.992653	0.000000	314.068426	
9892	24730	785.394062	12.170505	7.565385	9.651755	5.842572	8.223500	5.664096	2.587586	49.458581	11.163878	348.045763	

9893 rows x 32 columns

STATISTICAL OVERVIEW

DATA DESCRIPTION

```
In [5]: DF_Train.describe().transpose()
```

Out[5]:

	count	mean	std	min	25%	50%	75%	max
id	14838	741.850000	4283.505982	0.000000	3709.250000	741.500000	11127.750000	14837.000000
BertzCT	14838	515.153604	542.456700	0.000000	149.103601	652.652685	343.090331	2237.316499
Ch1n	14838	9.135189	6.819989	0.000000	4.680739	6.485270	11.170477	69.551167
Ch1v	14838	5.854307	4.647064	0.000000	2.844556	4.052701	7.486791	50.174588
Ch2n	14838	6.738497	5.866444	0.000000	2.932442	4.392859	5.778859	53.431954
Ch2v	14838	4.432570	3.760516	0.000000	1.949719	2.970427	5.788793	32.195368
Ch3v	14838	5.253221	4.925065	0.000000	2.034668	3.242775	6.609350	34.579313
Ch4n	14838	3.418749	3.436208	0.000000	1.160763	1.948613	4.502070	22.880636
EState_VSA1	14838	29.202823	31.728679	0.000000	5.969305	17.353601	44.876559	363.705954
EState_VSA2	14838	1.0435316	1.365184	0.000000	0.000000	0.000000	12.841643	99.936429
ExactMolWt	14838	292.623087	225.384140	1.007276	148.037173	256.042653	343.090331	2237.316499
FpDensityMorgan1	14838	1.262674	5.491284	0.66600000	1.045455	1.250000	1.500000	3.000000
FpDensityMorgan2	14838	1.812070	4.959565	0.66600000	1.690009	1.865152	2.062153	3.200000
FpDensityMorgan3	14838	2.255470	5.501200	0.66600000	2.100000	2.358491	2.500000	3.400000
HallKierAlpha	14838	-1.207776	0.935314	-7.730000	-1.660000	-1.100000	-0.570000	0.820000
HeavyAtomMolWt	14838	274.955021	212.678755	0.000000	136.190000	194.275000	320.121000	1758.846000
Kappa3	14838	5.874372	45.730226	-10.040000	1.784008	3.261011	5.772840	1512.242231
MaxAbsEStateIndex	14838	10.556443	1.559331	0.000000	9.926190	10.421334	11.539743	15.630251
MinEStateIndex	14838	-2.119772	2.066415	-6.327514	-4.659604	-1.265370	-0.787037	6.000000
NumHeteroatoms	14838	8.584108	7.643769	0.000000	4.000000	6.000000	10.000000	42.000000
PEOE_VSA10	14838	11.021644	13.958962	0.000000	0.000000	6.041841	18.311899	97.663462
PEOE_VSA14	14838	1.7276011	3.4561655	0.000000	0.000000	5.969305	1.572860	15.000000
PEOE_VSA6	14838	8.962400	19.756727	0.000000	0.000000	0.000000	12.132734	482.434223
PEOE_VSA7	14838	11.318811	20.169745	0.000000	0.000000	0.000000	13.847474	211.501279
PEOE_VSA8	14838	6.704867	10.865415	0.000000	0.000000	0.000000	6.923737	100.348416
SMR_VSA10	14838	15.667765	18.080008	0.000000	5.969305	11.752550	17.721856	80.742293
SMR_VSA5	14838	31.066423	33.896638	0.000000	6.420822	20.075376	42.727765	492.729739
SlogP_VSA3	14838	1.936941	14.396554	-0.000000	4.794537	9.589074	14.912684	115.406157
VSA_EState9	14838	49.309959	29.174824	5.450056	30.000000	41.666667	56.096650	384.450519
fr_COO	14838	0.458215	0.667948	0.000000	0.000000	0.000000	1.000000	8.000000
fr_COO2	14838	0.582926	0.668111	0.000000	0.000000	0.000000	1.000000	8.000000
EC1	14838	0.667745	0.477038	0.000000	0.000000	1.000000	1.000000	1.000000
EC2	14838	0.798962	0.400790	0.000000	0.000000	1.000000	1.000000	1.000000
EC3	14838	0.313789	0.464047	0.000000	0.000000	0.000000	1.000000	1.000000
EC4	14838	0.279081	0.448562	0.000000	0.000000	0.000000	1.000000	1.000000
EC5	14838	0.144831	0.351942	0.000000	0.000000	0.000000	0.000000	1.000000
EC6	14838	0.151570	0.358616	0.000000	0.000000	0.000000	0.000000	1.000000

Out[6]:

	count	mean	std	min	25%	50%	75%	max
id	9893	19784.000000	2856.007440	14838.000000	17311.000000	19784.000000	22257.000000	24730.000000
BertzCT	9893	516.411916	544.327795	0.000000	150.285577	289.901774	652.758463	5175.541449
Ch1n	9893	9.106998	6.775361	0.000000	4.696377	6.447265	10.966690	69.551167
Ch1v	9893	5.848047	4.541662	0.000000	2.846050	4.009996	7.480880	42.282925
Ch2v	9893	6.732639	5.865886	0.000000	2.934030	4.337841	8.528316	53.990574
Ch2n	9893	4.428979	3.770031	0.000000	1.949719	2.930013	5.788793	36.368883
Ch3v	9893	5.247994	4.939702	0.000000	2.049137	3.168052	6.609350	34.579313
Ch4n	9893	3.401083	3.431766	0.000000	1.171060	1.923982	4.502610	26.736931
EState_VSA1	9893	17.50078	1.837159	0.000000	5.058512	1.058931	2.509394	15.620667
EState_VSA2	9893	28.956235	31.470865	0.000000	5.969305	17.282269	44.876559	363.705954
EState_VSA2	9893	10.534500	13.768117	0.000000	0.000000	6.420822	12.841643	99.936429
ExactMolWt	9893	292.006497	224.667454	15.007276	148.073559	206.021523	342.116212	1888.793611
FpDensityMorgan1	9893	1.280917	0.361229	-1.133333	1.043478	1.250000	1.500000	3.000000
FpDensityMorgan2	9893	1.861792	0.348650	0.416667	1.692308	1.866667	2.071429	3.200000
FpDensityMorgan3	9893	2.304095	5.692551	0.000000	2.000000	2.358491	2.500000	3.400000
HallKierAlpha	9893	-1.206781	0.932173	-7.620000	-1.680000	-1.100000	-0.570000	0.830000
HeavyAtomMolWt	9893	274.583106	212.321052	14.007000	140.050000	194.125000	320.121000	1758.846000
Kappa3	9893	5.280880	37.349006	-10.040000	1.788507	3.261011	5.772840	1512.242231
MaxAbsEStateIndex	9893	10.555415	1.572745	0.000000	9.946009	10.418624	11.526266	14.630251
MinEStateIndex	9893	-2.098765	2.059680	-6.117075	-4.638889	-1.252751	-0.787037	4.750000
NumHeteroatoms	9893	8.590215	7.660447	0.000000	4.000000	6.000000	10.000000	42.000000
PEOE_VSA10	9893	11.032300	13.905575	0.000000	0.000000	6.041841	18.311899	122.017202
PEOE_VSA14	9893	18.493154	35.832279	0.000000	0.000000	5.969305	17.907916	482.434223
PEOE_VSA6	9893	8.917260	19.358857	0.000000	0.000000	0.000000	12.132734	258.844527
PEOE_VSA7	9893	11.222822	20.015732	0.000000	0.000000	0.000000	13.847474	271.424271
PEOE_VSA8	9893	6.789234	10.778071					


```
[30]: from xgboost import XGBClassifier

XGB = XGBClassifier(random_state = 42 )

Parameters = { 'n_estimators': [ 100, 200 ],

               'learning_rate': [ 0.5 , 0.001],

               'max_depth': [3, 5 ] }

GS = GridSearchCV ( estimator = XGB , param_grid = Parameters , cv = 5 )

GS.fit ( X , Y_EC1 )

Best_Parameters = GS.best_params_

print ( "The best parameters are as follows : \n\n" , Best_Parameters )

XGB = GradientBoostingClassifier ( **Best_Parameters )

XGB.fit ( X , Y_EC1 )

XGB_Result = pd.DataFrame ( )

XGB_Result [ 'EC1_Predicted' ] = XGB.predict ( DF_Test )

display ( XGB_Result.EC1_Predicted.value_counts ( ) )

XGB_Result

The best parameters are as follows :

{'learning_rate': 0.001, 'max_depth': 5, 'n_estimators': 200}
1    9893
Name: EC1_Predicted, dtype: int64

Out[30]:
```

	EC1_Predicted
0	1
1	1
2	1
3	1
4	1
...	...
9888	1
9889	1
9890	1
9891	1
9892	1

9893 rows x 1 columns

```
In [32]: XGB = XGBClassifier( random_state = 42 )

Parameters = { 'n_estimators': [ 100, 200 ],

               'learning_rate': [ 0.5 , 0.001],

               'max_depth': [3, 5 ] }

GS = GridSearchCV ( estimator = XGB , param_grid = Parameters , cv = 5 )

GS.fit ( X , Y_EC2 )

Best_Parameters = GS.best_params_

print ( "The best parameters are as follows : \n\n" , Best_Parameters )

XGB = XGBClassifier ( **Best_Parameters )

XGB.fit ( X , Y_EC2 )

XGB_Result [ 'EC2_Predicted' ] = XGB.predict ( DF_Test )

display ( XGB_Result.EC2_Predicted.value_counts ( ) )

XGB_Result

The best parameters are as follows :

{'learning_rate': 0.001, 'max_depth': 3, 'n_estimators': 200}
1    9893
Name: EC2_Predicted, dtype: int64

Out[32]:
```

	EC1_Predicted	EC2_Predicted
0	1	1
1	1	1
2	1	1
3	1	1
4	1	1
...
9888	1	1
9889	1	1
9890	1	1
9891	1	1
9892	1	1

9893 rows x 2 columns

```
In [ ]:
```