

In [304]:

```

1  #Best Model - Simple SVM
2  #Best cross-validation score: 0.906
3  #Best Parameters {'C': 0.001, 'multi_class': 'crammer_singer', 'penalty': 'l1'}
4  #train score: 0.904
5  #test score: 1.0
6
7  #Three things i've tried that were not covered in class are-
8  #-->Randomized search CV with MLPClassifier model
9  #-->Saving and Uploading models using pickle
10 #-->GaussianProcessClassifier
11
12 #Comments
13 #--worked a lot on hyperparameter tuning and selected the most appropriate so that
14 #Referred https://scikit-learn.org/ for hyperparameters
15 #--used randomsampling of the train dataset to run models faster
16 #--I've split the train dataset into train and test so that i can verify how my
17 # before predicting the test dataset provided and uploading on kaggle.
18 #--Made sure to resolve all the warnings
19

```

In [399]:

```

1  #loading test dataset into "test"
2  import pandas as pd
3
4  test= pd.read_csv(r'/Users/Desktop/MLPROJECT2/test.csv')

```

In [400]:

```

1  #Random Sampling of train data and loading into "data"
2  import random
3
4  f = "/Users/Desktop/MLPROJECT2/train.csv"
5  num_lines = sum(1 for l in open(f))
6  size = int(num_lines /10)
7  skip_idx = random.sample(range(1, num_lines), num_lines - size)
8  data = pd.read_csv(f, skiprows=skip_idx)

```

In [401]:

```

1  #to see the data to get an understanding
2  data.head()

```

Out[401]:

	Id	V1	V2	V3	V4	V5	V6	V7	
0	245376	1.812653	-0.476162	-0.338988	1.386750	-0.745965	-0.449870	-0.492226	0.0299
1	202483	-1.014219	0.522775	-0.337978	-1.957797	3.578395	3.266965	0.602857	0.6446
2	221275	-0.600317	0.915481	0.835746	-0.621856	1.284060	-0.070605	1.042698	-0.0725
3	155702	-0.524867	0.751315	2.278360	-0.295041	0.303778	0.164136	0.454016	-0.1106
4	76252	-16.772703	-14.426415	-5.606166	2.828980	-4.856624	1.538447	3.147259	-0.8755

5 rows × 31 columns

In [402]:

```
1 # Find out number or rows, columns and datatypes of all variables
2 data.info()
```

<class 'pandas.core.frame.DataFrame'>

RangeIndex: 2483 entries, 0 to 2482

Data columns (total 31 columns):

#	Column	Non-Null Count	Dtype
0	Id	2483 non-null	int64
1	V1	2250 non-null	float64
2	V2	2483 non-null	float64
3	V3	2483 non-null	float64
4	V4	2483 non-null	float64
5	V5	2483 non-null	float64
6	V6	2483 non-null	float64
7	V7	2483 non-null	float64
8	V8	2483 non-null	float64
9	V9	2483 non-null	float64
10	V10	2483 non-null	float64
11	V11	2483 non-null	float64
12	V12	2483 non-null	float64
13	V13	2483 non-null	float64
14	V14	2483 non-null	float64
15	V15	2483 non-null	float64
16	V16	2483 non-null	float64
17	V17	2483 non-null	float64
18	V18	2483 non-null	float64
19	V19	2483 non-null	float64
20	V20	2232 non-null	float64
21	V21	2483 non-null	float64
22	V22	2483 non-null	float64
23	V23	2483 non-null	float64
24	V24	2483 non-null	float64
25	V25	2483 non-null	float64
26	V26	2483 non-null	float64
27	V27	2483 non-null	float64
28	V28	2483 non-null	float64
29	V29	2483 non-null	float64
30	Target	2483 non-null	int64

dtypes: float64(29), int64(2)

memory usage: 601.5 KB

In [403]:

```
1 #To find corr between variables and drop highly correlated values if any
2 data.corr()
```

Out[403]:

	Id	V1	V2	V3	V4	V5	V6	V7	
Id	1.000000	0.106002	-0.049306	-0.308952	-0.150191	0.184762	-0.062611	0.091785	-0
V1	0.106002	1.000000	0.045613	0.238650	-0.128216	0.149235	0.014607	0.058938	-0
V2	-0.049306	0.045613	1.000000	-0.048024	0.066016	-0.001432	-0.085869	-0.232418	0
V3	-0.308952	0.238650	-0.048024	1.000000	-0.146819	0.180421	0.008765	0.136767	-0
V4	-0.150191	-0.128216	0.066016	-0.146819	1.000000	-0.082313	-0.058040	-0.055540	0
V5	0.184762	0.149235	-0.001432	0.180421	-0.082313	1.000000	0.025338	0.052438	-0
V6	-0.062611	0.014607	-0.085869	0.008765	-0.058040	0.025338	1.000000	0.070925	-0
V7	0.091785	0.058938	-0.232418	0.136767	-0.055540	0.052438	0.070925	1.000000	-0
V8	-0.026564	-0.026575	0.053217	-0.083000	0.037203	-0.048890	-0.069330	-0.053344	1
V9	0.014134	0.014293	-0.115891	0.111096	-0.108221	0.028536	0.034843	0.149379	-0
V10	0.041763	0.132935	-0.162679	0.241818	-0.148472	0.119695	0.078622	0.224967	-0
V11	-0.219083	-0.120230	0.058261	-0.188455	0.146493	-0.111462	-0.033933	-0.130512	0
V12	0.119889	0.121350	-0.134130	0.262103	-0.219359	0.101383	0.093871	0.211309	-0
V13	-0.045369	-0.027993	0.011601	-0.048120	0.000784	-0.027940	0.000500	-0.004749	0
V14	-0.071924	0.118215	-0.111523	0.290048	-0.242955	0.142852	0.077512	0.207104	-0
V15	-0.170863	-0.032225	-0.044510	-0.042316	-0.003529	-0.027638	0.021803	0.017806	-0
V16	0.015783	0.075772	-0.115365	0.181669	-0.106318	0.123865	0.089122	0.209749	-0
V17	-0.004735	0.211444	-0.121221	0.265299	-0.161811	0.198831	0.074465	0.242020	-0
V18	0.108321	0.112264	-0.011928	0.146607	-0.093861	0.110579	-0.015746	0.107366	-0
V19	0.009758	-0.020820	0.028932	-0.051656	0.004942	-0.084806	0.003313	-0.026424	0
V20	-0.035656	0.158838	0.299548	0.044290	-0.053365	0.101502	-0.087152	-0.232085	0
V21	0.056391	0.015946	0.112762	0.015648	-0.037808	0.052674	-0.088302	-0.028338	0
V22	0.157150	0.013514	-0.011440	0.017954	-0.011248	0.008160	0.030719	0.014045	-0
V23	0.016082	0.161596	0.138701	0.066591	-0.074991	0.008635	0.002211	-0.003186	-0
V24	-0.014286	0.025835	-0.027349	-0.009417	-0.064187	0.004263	0.044203	0.026326	0
V25	-0.203261	0.038595	0.027026	0.001164	0.023302	-0.030131	-0.011469	-0.017525	0
V26	-0.036459	-0.051888	-0.023069	-0.010423	0.028182	0.030107	0.002508	-0.019956	0
V27	-0.018267	-0.127714	-0.091841	-0.105243	0.081911	0.007896	-0.018641	-0.064976	0
V28	-0.035357	0.161877	0.142665	0.037479	-0.022143	-0.024816	0.006421	-0.036404	0
V29	-0.011759	-0.278976	-0.544283	-0.227559	0.102688	-0.363458	0.228540	0.414344	-0
Target	-0.013973	-0.185015	0.195992	-0.374903	0.312903	-0.167710	-0.120355	-0.270939	0

31 rows × 31 columns

In [404]:

```
1 # Column with Null Values and it's count in Train dataset
2 null_values= data.isnull().sum()
3 for key,value in null_values.items():
4     if value >0:
5         print(key,":",value)
```

V1 : 233
V20 : 251

In [405]:

```
1 # Column with Null Values and it's count in test dataset
2 null_values= test.isnull().sum()
3 for key,value in null_values.items():
4     if value >0:
5         print(key,":",value)
```

V1 : 2499
V20 : 2504

In [406]:

```
1 #No of unique values in each column in train set
2 total_unique_values= data.nunique()
3 for key,value in total_unique_values.items():
4     if value >0:
5         print(key, ":", value)
```

```
Id : 2483
V1 : 2250
V2 : 2482
V3 : 2482
V4 : 2482
V5 : 2482
V6 : 2482
V7 : 2482
V8 : 2482
V9 : 2482
V10 : 2482
V11 : 2482
V12 : 2482
V13 : 2482
V14 : 2482
V15 : 2482
V16 : 2482
V17 : 2482
V18 : 2482
V19 : 2482
V20 : 2231
V21 : 2482
V22 : 2482
V23 : 2482
V24 : 2482
V25 : 2482
V26 : 2482
V27 : 2482
V28 : 2482
V29 : 1452
Target : 2
```

In [407]:

```
1 #Split Train dataset into train and test 70-30
2 from sklearn.model_selection import train_test_split
3 X_train, X_test, y_train, y_test = train_test_split(data.drop(['Target'], axis=1),
4                                                     data['Target'],
5                                                     test_size=0.3,
6                                                     random_state=0)
7
8 X_train.shape, X_test.shape
```

Out[407]:

```
((1738, 30), (745, 30))
```

In [408]:

```
1 #number of columns in test to check if they match with train dataset
2 print(test.shape)
```

```
(24846, 30)
```

In [409]:

```
1 #feature engineering
2 from sklearn.pipeline import Pipeline
3 from feature_engine import missing_data_imputers as mdi
4 from sklearn.model_selection import train_test_split
5 from sklearn.preprocessing import StandardScaler
6
7
8     #Replacing Null values with the column median values
9 preprocess = Pipeline([
10     ('imputer_num', mdi.MeanMedianImputer(imputation_method='median',
11     variables=['V1', 'V20'])),
12     # feature Scaling
13     ('scaler', StandardScaler())
14
15
16 ])
```

In [410]:

```
1 preprocess.fit(X_train,y_train)
```

Out[410]:

```
Pipeline(memory=None,
          steps=[('imputer_num',
                  MeanMedianImputer(imputation_method='median',
                                     variables=['V1', 'V20'])),
                 ('scaler',
                  StandardScaler(copy=True, with_mean=True, with_std=True)
                 )],
          verbose=False)
```

In [411]:

```
1 X_train=preprocess.transform(X_train)
2 X_test=preprocess.transform(X_test)
3 X_testt=preprocess.transform(test)
```

In [218]:

```
1 # Basic Algorithm
```

In [219]:

```
1 #f2 score
2 from sklearn.metrics import fbeta_score, make_scorer
3 ftwo_scorer = make_scorer(fbeta_score, beta=2)
```

In [278]:

```
1 #logistic regression
2 from sklearn.linear_model import LogisticRegression
3 from sklearn.model_selection import GridSearchCV
4
5 logreg=LogisticRegression(max_iter=10000)
6
7 logreg_param= {
8     'C': [0.01,0.1,1],
9     'penalty':['l2','l1'],
10    'solver': ['newton-cg', 'lbfgs', 'sag'],
11    'multi_class':['auto']
12 }
13 logreg_grid = GridSearchCV(logreg,logreg_param,cv=5,n_jobs=-1,scoring=f'two_score
14 logreg_grid.fit(X_train,y_train)
15 print(f'Best Mean Cross Validation Score is {logreg_grid.best_score_}')
16 print(f'Best Mean Cross Validation Score is {logreg_grid.best_params_}')
17 print(f'Train score is {logreg_grid.score(X_train,y_train)}')
18 print(f'Test score is {logreg_grid.score(X_test,y_test)}')
19
```

Best Mean Cross Validation Score is 0.7911340852130325

Best Mean Cross Validation Score is {'C': 1, 'multi_class': 'auto', 'penalty': 'l2', 'solver': 'newton-cg'}

Train score is 0.913978494623656

Test score is 0.9615384615384615

In [265]:

```

1  #Decision Tree
2
3  from sklearn.metrics import mean_squared_error
4  from sklearn.metrics import r2_score
5  from math import sqrt
6  from sklearn.tree import DecisionTreeClassifier
7
8  dtree = DecisionTreeClassifier(random_state=0)
9
10 param_dtree = {'max_depth': [1,2,3,5],
11                'splitter':['random','best'],
12                'max_features':['auto','log2','sqrt'],
13                'criterion':['gini','entropy'],
14
15                }
16
17
18 grid_dtree = GridSearchCV(dtree, param_dtree, cv=5,n_jobs=-1,scoring=ftwo_score)
19 grid_dtree.fit(X_train, y_train)
20
21 print("Best Mean Cross-validation score: {:.3f}".format(grid_dtree.best_score_))
22 print('Decision Tree parameters: ', grid_dtree.best_params_)
23 print('Train score: ', grid_dtree.score(X_train, y_train))
24 print("Test Score:", grid_dtree.score(X_test,y_test))
25
26 X_test_preds = grid_dtree.predict(X_testt)
27
28
29
30 pd.DataFrame({'Id': test.Id, 'Target':X_test_preds }).to_csv('solution_base15.csv')
31 print("Done :D")
32
33

```

Best Mean Cross-validation score: 0.745

Decision Tree parameters: {'criterion': 'gini', 'max_depth': 2, 'max_features': 'auto', 'splitter': 'best'}

Train score: 0.9042553191489363

Test Score: 1.0

Done :D

In [261]:

```
1 #knn
2 from sklearn.model_selection import GridSearchCV
3 from sklearn.neighbors import KNeighborsClassifier
4
5 knn = KNeighborsClassifier()
6
7 param_knn = {
8     'n_neighbors': [3],
9     'weights': ['uniform', 'distance'],
10    'algorithm' :['auto'],
11    'p':[2,3]
12
13    }
14
15 #apply grid search
16
17 grid_knn = GridSearchCV(knn, param_knn, cv=5,n_jobs=-1,scoring=ftwo_scorer)
18 grid_knn.fit(X_train, y_train)
19
20 # Mean Cross Validation Score
21 print("Best Mean Cross-validation score: {:.3f}".format(grid_knn.best_score_))
22 print('KNN parameters: ', grid_knn.best_params_)
23 print('Train score: ', grid_knn.score(X_train, y_train))
24 print("KNN Test Performance: ", grid_knn.score(X_test,y_test))
25
26
27
```

Best Mean Cross-validation score: 0.817

KNN parameters: {'algorithm': 'auto', 'n_neighbors': 3, 'p': 2, 'weights': 'uniform'}

Train score: 0.8152173913043478

KNN Test Performance: 1.0

In [263]:

```
1  #Simple SVM
2
3  import matplotlib.gridspec as gridspec
4  import itertools
5  from sklearn.model_selection import train_test_split
6  from sklearn.svm import SVC, LinearSVC
7  from sklearn.model_selection import GridSearchCV
8
9  #svc = SVC()
10
11 svc = LinearSVC(max_iter=100000)
12 svc.fit(X_train, y_train)
13
14 #define a list of parameters
15 param_svc_kernel = {'C': [0.001,0.01,0.1,1],
16                     'penalty':['l1','l2'],
17                     'multi_class':['ovr','crammer_singer']}
18
19 }
20
21 #apply grid search
22 grid_svc = GridSearchCV(svc,param_svc_kernel, cv=5,n_jobs=-1,scoring=ftwo_score)
23 grid_svc.fit(X_train, y_train)
24
25 print('Best cross-validation score:', grid_svc.best_score_)
26 print(grid_svc.best_params_)
27 print('train score: ', grid_svc.score(X_train, y_train))
28 print('test score: ', grid_svc.score(X_test, y_test))
29
30
31
32 X_test_preds = grid_svc.predict(X_testt)
33 pd.DataFrame({'Id': test.Id, 'Target':X_test_preds }).to_csv('solution_base3.csv')
34 print("Done :D")
35
36
```

Best cross-validation score: 0.906265664160401

{'C': 0.001, 'multi_class': 'crammer_singer', 'penalty': 'l1'}

train score: 0.9042553191489363

test score: 1.0

Done :D

In [312]:

```

1  #SVM with kernel='rbf'
2  import matplotlib.gridspec as gridspec
3  import itertools
4  from sklearn.model_selection import train_test_split
5  from sklearn.svm import SVC, LinearSVC
6  from sklearn.model_selection import GridSearchCV
7  import numpy as np
8
9
10 svc = SVC()
11 svc.fit(X_train, y_train)
12 svc_kernel = SVC(kernel = 'rbf')
13
14 #define a list of parameters
15 C_range = 10. ** np.arange(-3, 8)
16 gamma_range = 10. ** np.arange(-5, 4)
17
18 param_svc_kernel = {'C': C_range,
19                     'gamma':gamma_range}
20
21 #apply grid search
22 grid_svc_rbf = GridSearchCV(svc_kernel, param_svc_kernel, cv=5, n_jobs=-1,scoring='r2')
23 grid_svc_rbf.fit(X_train, y_train)
24 print('Best cross-validation score:', grid_svc_rbf.best_score_)
25 print(grid_svc_rbf.best_params_)
26 print('train score: ', grid_svc_rbf.score(X_train, y_train))
27 print('test score: ', grid_svc_rbf.score(X_test, y_test))
28
29
30 X_test_preds = grid_svc_rbf.predict(X_testtt)
31 pd.DataFrame({'Id': test.Id, 'Target':X_test_preds }).to_csv('solution_base.csv')
32 print("Done :D")
33
34
35

```

Best cross-validation score: 0.906265664160401

{'C': 10.0, 'gamma': 0.001}

train score: 0.913978494623656

test score: 1.0

Done :D

In [266]:

```

1  #Random Forest
2
3  from sklearn.ensemble import RandomForestClassifier
4  rfc = RandomForestClassifier(random_state=0)
5  rfc_param = {
6      'n_estimators': [100],
7      'max_features': ['auto', 'sqrt', 'log2'],
8      'max_depth' : [1,2,3,5],
9      'criterion' :['gini', 'entropy'],
10     'min_samples_split' :[2],
11     'min_samples_leaf':[3],
12
13 }
14 rfc_grid = GridSearchCV(rfc, rfc_param,cv=5,n_jobs=-1,scoring=ftwo_scorer)
15 rfc_grid.fit(X_train,y_train)
16 print(f'Best Mean Cross Validation Score is {rfc_grid.best_score_}')
17 print(f'Best param {rfc_grid.best_params_}')
18 print(f'Train score is {rfc_grid.score(X_train,y_train)}')
19 print(f'Test score is {rfc_grid.score(X_test,y_test)}')
20

```

Best Mean Cross Validation Score is 0.6660179861418251

Best param {'criterion': 'gini', 'max_depth': 2, 'max_features': 'auto', 'min_samples_leaf': 3, 'min_samples_split': 2, 'n_estimators': 100}

Train score is 0.8152173913043478

Test score is 0.8333333333333334

Done :D

In [226]:

```

1  #Extra trees
2  from sklearn.ensemble import ExtraTreesClassifier
3  etc = ExtraTreesClassifier(random_state=0)
4  etc_param = {
5      'n_estimators': [100],
6      'max_features': ['auto', 'sqrt', 'log2'],
7      'max_depth' : [1,2,3,5],
8      'criterion' :['gini', 'entropy'],
9      'min_samples_split':[2,3],
10
11
12 }
13 etc_grid = GridSearchCV(etc, etc_param,cv=5,n_jobs=-1,scoring=ftwo_scorer)
14 etc_grid.fit(X_train,y_train)
15
16 print(f'Best Mean Cross Validation Score is {etc_grid.best_score_}')
17 print(f'Best param {etc_grid.best_params_}')
18 print(f'Train score is {etc_grid.score(X_train,y_train)}')
19 print(f'Test score is {etc_grid.score(X_test,y_test)}')
20

```

Best Mean Cross Validation Score is 0.6255417956656346

Best param {'criterion': 'gini', 'max_depth': 5, 'max_features': 'auto', 'min_samples_split': 2, 'n_estimators': 100}

Train score is 0.8695652173913043

Test score is 0.6521739130434783

Done :D

In [227]:

```
1 #Gradient boost
2 from sklearn.ensemble import GradientBoostingClassifier
3 gbc= GradientBoostingClassifier(random_state=0)
4 gbc_param = {
5     'max_depth' : [1],
6     'n_estimators' : [200],
7     'learning_rate' : [0.1],
8     'subsample':[1],
9     'min_samples_split':[2],
10    'min_samples_leaf':[4],
11    'max_features': ['auto','sqrt','log2'],
12    'loss': ['deviance','exponential'],
13
14
15    }
16 gbc_grid = GridSearchCV(gbc, gbc_param,cv=5,n_jobs=-1,scoring=ftwo_scorer)
17 gbc_grid.fit(X_train,y_train)
18 print(f'Best Mean Cross Validation Score is {gbc_grid.best_score_}')
19 print(f'Best param {gbc_grid.best_params_}')
20 print(f'Train score is {gbc_grid.score(X_train,y_train)}')
21 print(f'Test score is {gbc_grid.score(X_test,y_test)}')
22
23
```

Best Mean Cross Validation Score is 0.8236842105263158

Best param {'learning_rate': 0.1, 'loss': 'exponential', 'max_depth': 1, 'max_features': 'auto', 'min_samples_leaf': 4, 'min_samples_split': 2, 'n_estimators': 200, 'subsample': 1}

Train score is 0.913978494623656

Test score is 0.9615384615384615

Done :D

In [268]:

```

1  #XGBoost
2
3  from xgboost import XGBClassifier
4  xgbc= XGBClassifier(random_state=0,early_stopping_rounds=4,objective= 'binary:logit')
5  xgbc_param = {
6      'max_depth' : [1,2,3,5],
7      'n_estimators' : [150],
8      'learning_rate' : [0.01,0.1,0.5],
9      'min_child_weight' : [3],
10     # 'subsample':[0.8]
11 }
12 xgbc_grid = GridSearchCV(xgbc, xgbc_param,cv=5,n_jobs=-1,scoring=ftwo_scorer)
13 xgbc_grid.fit(X_train,y_train)
14 print(f'Best Mean Cross Validation Score is {xgbc_grid.best_score_}')
15 print(f'Best param {xgbc_grid.best_params_}')
16 print(f'Train score is {xgbc_grid.score(X_train,y_train)}')
17 print(f'Test score is {xgbc_grid.score(X_test,y_test)}')
18

```

Best Mean Cross Validation Score is 0.8094820384294069

Best param {'learning_rate': 0.01, 'max_depth': 1, 'min_child_weight': 3, 'n_estimators': 150}

Train score is 0.8241758241758242

Test score is 0.8333333333333334

Done :D

In [393]:

```

1  #Dummy Classifier
2
3  from sklearn.dummy import DummyClassifier
4
5  dummy= DummyClassifier()
6
7  dummy_param= {
8      "strategy":["prior','uniform','stratified','most_frequent','constant']
9  }
10
11
12 dummy_grid = GridSearchCV(dummy, dummy_param,cv=5,n_jobs=-1,scoring=ftwo_scorer)
13 dummy_grid.fit(X_train,y_train)
14 print(f'Best Mean Cross Validation Score is {dummy_grid.best_score_}')
15 print(f'Best param {dummy_grid.best_params_}')
16 print(f'Train score is {dummy_grid.score(X_train,y_train)}')
17 print(f'Test score is {dummy_grid.score(X_test,y_test)}')
18

```

Best Mean Cross Validation Score is 0.06134344811097836

Best param {'strategy': 'uniform'}

Train score is 0.05241090146750525

Test score is 0.04705882352941176

In [229]:

```

1  #Cost Sensitive Algorithms

```

In [412]:

```

1  #logistic regression
2  from sklearn.linear_model import LogisticRegression
3  from sklearn.model_selection import GridSearchCV
4  from sklearn.model_selection import RepeatedStratifiedKFold
5
6  logreg=LogisticRegression(max_iter=1000)
7
8  logreg_param= {
9      'C': [0.01],
10     'penalty':['l1', 'l2'],
11     'class_weight': ['balanced', 'none'],
12     'solver': ['lbfgs'],
13     'max_iter': [200],
14     'multi_class': ['auto'],
15     'class_weight': [{0:1, 1:10}]
16 }
17
18 cv = RepeatedStratifiedKFold(n_splits=7, n_repeats=3, random_state=1)
19
20 clogreg_grid = GridSearchCV(logreg, logreg_param, cv=cv, n_jobs=-1, scoring='f2_score')
21 clogreg_grid.fit(X_train, y_train)
22 print(f'Best Mean Cross Validation Score is {clogreg_grid.best_score_}')
23 print(f'Best param {clogreg_grid.best_params_}')
24 print(f'Train score is {clogreg_grid.score(X_train, y_train)}')
25 print(f'Test score is {clogreg_grid.score(X_test, y_test)}')
26
27
28
29 X_train_preds = clogreg_grid.predict(X_train)
30 X_test_preds = clogreg_grid.predict(X_test)
31
32 print('train rmse: {}'.format(sqrt(mean_squared_error(y_train, X_train_preds))))
33 print('test rmse: {}'.format(sqrt(mean_squared_error(y_test, X_test_preds))))
34
35

```

Best Mean Cross Validation Score is 0.6744916804901324

Best param {'C': 0.01, 'class_weight': {0: 1, 1: 10}, 'max_iter': 200, 'multi_class': 'auto', 'penalty': 'l2', 'solver': 'lbfgs'}

Train score is 0.7575757575757575

Test score is 0.9090909090909091

train rmse: 0.07955572841757301

test rmse: 0.08192319205190404

In [413]:

```

1  #Decision Tree
2
3  from sklearn.metrics import mean_squared_error
4  from sklearn.metrics import r2_score
5  from math import sqrt
6  from sklearn.tree import DecisionTreeClassifier
7  dtree = DecisionTreeClassifier(random_state=0)
8
9
10 #define a list of parameters
11 param_dtree = {'max_depth': [3],
12                'class_weight': [{0:100,1:1}, {0:10,1:1}, {0:1,1:1}, {0:1,1:10}],
13                'splitter':['random','best'],
14                'max_features':['auto','sqrt','log2']}
15
16
17 cgrid_dtree = GridSearchCV(dtree, param_dtree, cv=5,n_jobs=-1,scoring=ftwo_score)
18 cgrid_dtree.fit(X_train, y_train)
19
20 print("Best Mean Cross-validation score: {:.3f}".format(cgrid_dtree.best_score_))
21 print('Decision Tree parameters: ', cgrid_dtree.best_params_)
22 print('Train score: ', cgrid_dtree.score(X_train, y_train))
23 print("Test Score:", cgrid_dtree.score(X_test,y_test))
24
25 X_train_preds = cgrid_dtree.predict(X_train)
26 X_test_preds = cgrid_dtree.predict(X_test)
27
28 print('train rmse: {}'.format(sqrt(mean_squared_error(y_train, X_train_preds))))
29 print('test rmse: {}'.format(sqrt(mean_squared_error(y_test, X_test_preds))))
30

```

Best Mean Cross-validation score: 0.673

Decision Tree parameters: {'class_weight': {0: 1, 1: 1}, 'max_depth': 3, 'max_features': 'log2', 'splitter': 'best'}

Train score: 0.7142857142857142

Test Score: 0.6382978723404255

train rmse: 0.0719608639321375

test rmse: 0.08192319205190404

In [414]:

```

1  #Simple SVM
2  import matplotlib.gridspec as gridspec
3  import itertools
4  from sklearn.model_selection import train_test_split
5  from sklearn.svm import SVC, LinearSVC
6  from sklearn.model_selection import GridSearchCV
7
8  #svc = SVC()
9
10 svc = LinearSVC(max_iter=100000)
11 svc.fit(X_train, y_train)
12
13
14
15 #define a list of parameters
16 param_svc_kernel = {'C': [0.01],
17                     'class_weight': [{0:100,1:1}, {0:10,1:1}, {0:1,1:1}, {0:1,1:1}],
18                     'penalty':['l1','l2'],
19                     'multi_class':['ovr','crammer_singer']}
20
21
22 #apply grid search
23
24 cgrid_svc = GridSearchCV(svc,param_svc_kernel, cv=5,n_jobs=-1,scoring=ftwo_score)
25 cgrid_svc.fit(X_train, y_train)
26 print('Best cross-validation score:', cgrid_svc.best_score_)
27 print(cgrid_svc.best_params_)
28 print('train score: ', cgrid_svc.score(X_train, y_train))
29 print('test score: ', cgrid_svc.score(X_test, y_test))
30
31 X_train_preds = cgrid_svc.predict(X_train)
32 X_test_preds = cgrid_svc.predict(X_test)
33
34 print('train rmse: {}'.format(sqrt(mean_squared_error(y_train, X_train_preds))))
35 print('test rmse: {}'.format(sqrt(mean_squared_error(y_test, X_test_preds))))
36
37

```

```

Best cross-validation score: 0.6836496836496837
{'C': 0.01, 'class_weight': {0: 1, 1: 1}, 'multi_class': 'crammer_singer', 'penalty': 'l1'}
train score:  0.703125
test score:   0.980392156862745
train rmse: 0.07955572841757301
test rmse: 0.03663716527236559

```

In [415]:

```
1  #Random Forest
2  from sklearn.model_selection import cross_val_score
3  from sklearn.model_selection import RepeatedStratifiedKFold
4
5  crf = RandomForestClassifier(n_estimators=150)
6  crf_param={
7
8      'max_depth' : [2,3,5],
9      'criterion' : ["gini","entropy"],
10     'max_features': ["auto","sqrt","log2"],
11     'class_weight':["balanced"]
12
13 }
14 crf_grid = GridSearchCV(crf, crf_param,cv=5,n_jobs=-1, scoring=ftwo_scorer)
15
16 crf_grid.fit(X_train,y_train)
17
18 print(f'Best Mean Cross Validation Score is {crf_grid.best_score_}')
19 print(f'Best param {crf_grid.best_params_}')
20 print(f'Train score is {crf_grid.score(X_train,y_train)}')
21 print(f'Test score is {crf_grid.score(X_test,y_test)}')
22
23 X_train_preds = crf_grid.predict(X_train)
24 X_test_preds = crf_grid.predict(X_test)
25
26 print('train rmse: {}'.format(sqrt(mean_squared_error(y_train, X_train_preds))))
27 print('test rmse: {}'.format(sqrt(mean_squared_error(y_test, X_test_preds))))
28
29
30
```

Best Mean Cross Validation Score is 0.6991216512955644

Best param {'class_weight': 'balanced', 'criterion': 'entropy', 'max_depth': 3, 'max_features': 'log2'}

Train score is 0.8076923076923078

Test score is 0.8823529411764706

train rmse: 0.06346351669793114

test rmse: 0.06345743169703524

In [416]:

```

1  #XGBoost
2  from xgboost import XGBClassifier
3  xgbc= XGBClassifier(random_state=0,early_stopping_rounds=2,objective= 'binary:logistic')
4
5  param_grid = {
6      # 'scale_pos_weight': [1, 10],
7          'max_depth' : [1,2,3,5],
8          'n_estimators' : [150],
9          'learning_rate' : [0.1],
10         'min_child_weight' : [1,2],
11         'subsample': [0.5],
12         'class_weight': [{0:100,1:1}, {0:10,1:1}, {0:1,1:1}, {0:1,1:10},
13     ]
14 }
15
16 cxgbc_grid = GridSearchCV(xgbc, param_grid,cv=10,n_jobs=-1, scoring=ftwo_scorer)
17 cxgbc_grid.fit(X_train,y_train)
18 print(f'Best Mean Cross Validation Score is {cxgbc_grid.best_score_}')
19 print(f'Best param {cxgbc_grid.best_params_}')
20 print(f'Train score is {cxgbc_grid.score(X_train,y_train)}')
21 print(f'Test score is {cxgbc_grid.score(X_test,y_test)}')
22
23
24 X_train_preds = cxgbc_grid.predict(X_train)
25 X_test_preds = cxgbc_grid.predict(X_test)
26
27 print('train rmse: {}'.format(sqrt(mean_squared_error(y_train, X_train_preds))))
28 print('test rmse: {}'.format(sqrt(mean_squared_error(y_test, X_test_preds))))
29
30

```

Best Mean Cross Validation Score is 0.6864468864468865

Best param {'class_weight': {0: 100, 1: 1}, 'learning_rate': 0.1, 'max_depth': 1, 'min_child_weight': 2, 'n_estimators': 150, 'subsample': 0.5}

Train score is 0.7421875000000001

Test score is 0.9615384615384615

train rmse: 0.0719608639321375

test rmse: 0.05181277601508398

In [417]:

```

1  #Extratrees
2  from sklearn.ensemble import ExtraTreesClassifier
3
4
5  etc= ExtraTreesClassifier(random_state=0,class_weight='balanced')
6  etc_param = {
7      'n_estimators': [100],
8      'max_features': ['auto','sqrt','log2'],
9      'max_depth' : [1,2,3,5],
10     'criterion' :['gini','entropy'],
11     'min_samples_split':[2,3],
12
13 }
14 cetc_grid = GridSearchCV(etc, etc_param,cv=5,n_jobs=-1,scoring=ftwo_scorer)
15 cetc_grid.fit(X_train,y_train)
16
17 print(f'Best Mean Cross Validation Score is {cetc_grid.best_score_}')
18 print(f'Best param {cetc_grid.best_params_}')
19 print(f'Train score is {cetc_grid.score(X_train,y_train)}')
20 print(f'Test score is {cetc_grid.score(X_test,y_test)}')
21
22 X_train_preds = cetc_grid.predict(X_train)
23 X_test_preds = cetc_grid.predict(X_test)
24
25 print('train rmse: {}'.format(sqrt(mean_squared_error(y_train, X_train_preds))))
26 print('test rmse: {}'.format(sqrt(mean_squared_error(y_test, X_test_preds))))
27
28

```

Best Mean Cross Validation Score is 0.7160528682267813
 Best param {'criterion': 'gini', 'max_depth': 3, 'max_features': 'auto', 'min_samples_split': 2, 'n_estimators': 100}
 Train score is 0.7835820895522388
 Test score is 0.9433962264150945
 train rmse: 0.07955572841757301
 test rmse: 0.06345743169703524

In [307]:

```

1  #Bagging Decision Tree with undersampling
2  from imblearn.ensemble import BalancedBaggingClassifier
3  from sklearn.model_selection import RepeatedStratifiedKFold
4  from sklearn.model_selection import cross_val_score
5
6  b_dtree = BalancedBaggingClassifier()
7  cv = RepeatedStratifiedKFold(n_splits=10, n_repeats=3, random_state=1)
8  scores = cross_val_score(b_dtree, X_train, y_train, scoring=ftwo_scorer, cv=cv,
9  print('Mean f2: %.3f' % scores.mean())
10

```

Mean f2: 0.691

In [423]:

```
1 scores1= cross_val_score(etc, X_train, y_train, scoring=ftwo_scorer, cv=5, n_jo
2 scores2= cross_val_score(b_dtree, X_train, y_train, scoring=ftwo_scorer, cv=5,
3 scores3= cross_val_score(xgbc, X_train, y_train, scoring=ftwo_scorer, cv=5, n_jo
4 scores4= cross_val_score(crf, X_train, y_train, scoring=ftwo_scorer, cv=5, n_jo
5 scores5= cross_val_score(svc, X_train, y_train, scoring=ftwo_scorer, cv=5, n_jo
6 scores6= cross_val_score(dtrees, X_train, y_train, scoring=ftwo_scorer, cv=5, n_
7 scores7= cross_val_score(logreg, X_train, y_train, scoring=ftwo_scorer, cv=5, n
8 print('mean scores of all cost sensitive algorithms')
9 print(f'etc {scores1.mean()}')
10 print(f'b_tree {scores2.mean()}')
11 print(f'xgbc {scores3.mean()}')
12 print(f'randomforest {scores4.mean()}')
13 print(f'svc {scores5.mean()}')
14 print(f'dtree {scores6.mean()}')
15 print(f'logreg {scores7.mean()}')
16
```

mean scores of all cost sensitive algorithms

etc 0.563338959890684

b_tree 0.5525844184139883

xgbc 0.637052685778323

randomforest 0.6914293436032567

svc 0.6556433904259992

dtree 0.6619252432155658

logreg 0.6223100570926658

In [237]:

```
1 #3. Data Sampling Algorithms
```

In [428]:

```

1  #Logistic Regression
2  from sklearn.linear_model import LogisticRegression
3  from imblearn.over_sampling import SMOTE
4  from imblearn.pipeline import Pipeline as p
5
6  pipe_roc_lg = p([('smote', SMOTE()), ('lg', LogisticRegression())])
7  param_roc_lg = {'smote__k_neighbors': [2],
8                  'lg__C': [0.0001],
9                  'lg__penalty': ['l1', 'l2'],
10                 'lg__max_iter': [150],
11                 'lg__solver': ['newton-cg'],
12                 'lg__class_weight': ["balanced"],
13                 'lg__multi_class': ['auto', 'ovr', 'multinomial']}
14
15
16  ogrid_lg= GridSearchCV(pipe_roc_lg,param_roc_lg, cv=5, n_jobs=-1, scoring=ftwo_s
17  ogrid_lg.fit(X_train, y_train)
18
19  print("Best Mean cross-validation score: {:.3f}".format(ogrid_lg.best_score_))
20  print("Best parameters: {}".format(ogrid_lg.best_params_))
21  print(f'Train score is {ogrid_lg.score(X_train,y_train)}')
22  print(f'Test score is {ogrid_lg.score(X_test,y_test)}')
23
24
25  # let's get the predictions
26  X_test_preds = ogrid_lg.predict(X_testt)
27
28
29
30  pd.DataFrame({'Id': test.Id, 'Target':X_test_preds }).to_csv('solution_base1.csv')
31  print("Done :D")
32
33

```

Best Mean cross-validation score: 0.699

Best parameters: {'lg__C': 0.0001, 'lg__class_weight': 'balanced', 'lg__max_iter': 150, 'lg__multi_class': 'multinomial', 'lg__penalty': 'l2', 'lg__solver': 'newton-cg', 'smote__k_neighbors': 2}

Train score is 0.7835820895522388

Test score is 0.9259259259259259

Done :D

In [343]:

```

1  #Decision Tree
2  pipe_roc_dtree = p([('smote', SMOTE()), ('dtree', DecisionTreeClassifier(random_state=42)),
3  param_roc_dtree = {'smote__k_neighbors': [2,3,5],
4                      'dtree__max_depth': [2,3,5],
5                      'dtree__splitter': ['random', 'best'],
6                      'dtree__max_features': ['auto', 'log2', 'sqrt'],
7                      'dtree__criterion': ['gini', 'entropy'],
8                      }
9
10 ogrid_dtree= GridSearchCV(pipe_roc_dtree,param_roc_dtree, cv=5, n_jobs=-1, scoring='roc_auc')
11 ogrid_dtree.fit(X_train, y_train)
12
13 print("Best Mean cross-validation score: {:.3f}".format(ogrid_dtree.best_score_))
14 print("Best parameters: {}".format(ogrid_dtree.best_params_))
15 print(f'Train score is {ogrid_dtree.score(X_train,y_train)}')
16 print(f'Test score is {ogrid_dtree.score(X_test,y_test)}')
17
18 # let's get the predictions
19 X_test_preds = ogrid_dtree.predict(X_test)
20
21
22
23 pd.DataFrame({'Id': test.Id, 'Target':X_test_preds }).to_csv('solution_base111.csv')
24 print("Done :D")
25

```

Best Mean cross-validation score: 0.868

Best parameters: {'dtree__criterion': 'entropy', 'dtree__max_depth': 3, 'dtree__max_features': 'log2', 'dtree__splitter': 'best', 'smote__k_neighbors': 5}

Train score is 0.9042553191489363

Test score is 0.8928571428571429

Done :D

In [398]:

```

1  #KNN
2
3  pipe_roc_knn = p([('smote', SMOTE()), ('knn', KNeighborsClassifier())])
4  param_roc_knn = {'smote__k_neighbors': [2, 3, 5],
5                   'knn__weights': ['uniform'],
6                   'knn__algorithm': ['kd_tree', 'brute'],
7                   'knn__p': [3]
8                   }
9
10 ogrid_knn= GridSearchCV(pipe_roc_knn, param_roc_knn, cv=5, n_jobs=-1, scoring='f1')
11 ogrid_knn.fit(X_train, y_train)
12
13 print("Best Mean cross-validation score: {:.3f}".format(ogrid_knn.best_score_))
14 print("Best parameters: {}".format(ogrid_knn.best_params_))
15 print(f'Train score is {ogrid_knn.score(X_train, y_train)}')
16 print(f'Test score is {ogrid_knn.score(X_test, y_test)}')
17
18

```

Best Mean cross-validation score: 0.873

Best parameters: {'knn__algorithm': 'kd_tree', 'knn__p': 3, 'knn__weights': 'uniform', 'smote__k_neighbors': 3}

Train score is 0.9693877551020408

Test score is 0.8928571428571429

In [241]:

```

1  #SVM rbf
2  from sklearn import svm
3
4  from sklearn.svm import SVC, LinearSVC
5
6  pipe_roc_svm = p([('smote', SMOTE()), ('svm', svm.SVC(kernel='rbf'))])
7  param_roc_svm = {'smote__k_neighbors': [2, 3, 5],
8                   'svm__C': [0.001, 0.01, 0.1],
9                   'svm__gamma': ['scale', 'auto'],
10                  }
11
12 ogrid_svm_rbf= GridSearchCV(pipe_roc_svm, param_roc_svm, cv=5, n_jobs=-1, scoring='f1')
13 ogrid_svm_rbf.fit(X_train, y_train)
14
15 print("Best Mean cross-validation score: {:.3f}".format(ogrid_svm_rbf.best_score_))
16 print("Best parameters: {}".format(ogrid_svm_rbf.best_params_))
17 print(f'Train score is {ogrid_svm_rbf.score(X_train, y_train)}')
18 print(f'Test score is {ogrid_svm_rbf.score(X_test, y_test)}')

```

Best Mean cross-validation score: 0.874

Best parameters: {'smote__k_neighbors': 2, 'svm__C': 0.1, 'svm__gamma': 'scale'}

Train score is 0.913978494623656

Test score is 1.0

In [347]:

```

1  #Random Forest
2  from sklearn.ensemble import RandomForestClassifier
3  from imblearn.over_sampling import SMOTE
4
5  pipe_roc_rf = p([('smote',SMOTE()),('rf',RandomForestClassifier(random_state=0))
6  param_roc_rf = {
7
8      'smote__k_neighbors': [2],
9      'rf__n_estimators' : [150],
10     'rf__max_depth' : [1],
11     'rf__criterion' : ["gini","entropy"],
12     'rf__min_samples_split' :[1,2],
13     'rf__max_features': ["auto","sqrt","log2"]
14
15     }
16
17  ogrid_rf= GridSearchCV(pipe_roc_rf,param_roc_rf, cv=7, n_jobs=-1, scoring=ftwo_s
18  ogrid_rf.fit(X_train, y_train)
19
20  print("Best Mean cross-validation score: {:.3f}".format(ogrid_rf.best_score_))
21  print("Best parameters: {}".format(ogrid_rf.best_params_))
22  print(f'Train score is {ogrid_rf.score(X_train,y_train)}')
23  print(f'Test score is {ogrid_rf.score(X_test,y_test)}')
24
25
26  # let's get the predictions
27  X_test_preds = ogrid_rf.predict(X_testt)
28
29
30
31  pd.DataFrame({'Id': test.Id, 'Target':X_test_preds }).to_csv('solution_base222.c
32  print("Done :D")
33
34

```

Best Mean cross-validation score: 0.862

Best parameters: {'rf__criterion': 'gini', 'rf__max_depth': 1, 'rf__max_features': 'log2', 'rf__min_samples_split': 2, 'rf__n_estimators': 150, 'smote__k_neighbors': 2}

Train score is 0.8602150537634409

Test score is 1.0

Done :D

In [352]:

```

1 #easy ensemble classifier
2 from imblearn.ensemble import EasyEnsembleClassifier
3 pipe_roc_ee = p([('smote', SMOTE()), ('ee', EasyEnsembleClassifier(random_state=0))
4 param_roc_ee = {'smote__k_neighbors': [1],
5                 'ee__n_estimators': [50]
6                 }
7
8 ogrid_ee= GridSearchCV(pipe_roc_ee,param_roc_ee, cv=5, n_jobs=-1, scoring=ftwo_s
9 ogrid_ee.fit(X_train, y_train)
10
11 print("Best Mean cross-validation score: {:.3f}".format(ogrid_ee.best_score_))
12 print("Best parameters: {}".format(ogrid_ee.best_params_))
13 print(f'Train score is {ogrid_ee.score(X_train,y_train)}')
14 print(f'Test score is {ogrid_ee.score(X_test,y_test)}')
15

```

Best Mean cross-validation score: 0.841

Best parameters: {'ee__n_estimators': 50, 'smote__k_neighbors': 1}

Train score is 1.0

Test score is 1.0

In [309]:

```

1 #XGBoost
2 from xgboost import XGBClassifier
3
4 pipe_roc_xgb = p([('smote', SMOTE()), ('xgb', XGBClassifier(random_state=0,
5
6
7
8
9 param_roc_xgb = {
10     'smote__k_neighbors': [1,2,3,5],
11     'xgb__learning_rate' : [0.01],
12     'xgb__max_depth' : [2,3,5]
13 }
14
15 ogrid_xgb= GridSearchCV(pipe_roc_xgb,param_roc_xgb, cv=5, n_jobs=-1, scoring=ftv
16 ogrid_xgb.fit(X_train, y_train)
17
18 print("Best Mean cross-validation score: {:.3f}".format(ogrid_xgb.best_score_))
19 print("Best parameters: {}".format(ogrid_xgb.best_params_))
20
21 print(f'Train score is {ogrid_xgb.score(X_train,y_train)}')
22 print(f'Test score is {ogrid_xgb.score(X_test,y_test)}')
23

```

Best Mean cross-validation score: 0.821

Best parameters: {'smote__k_neighbors': 3, 'xgb__learning_rate': 0.01, 'xgb__max_depth': 3}

Train score is 0.9693877551020408

Test score is 0.9259259259259259

Done :D

In [251]:

```

1  #Neural Network(scikit learn MLPClassifier) with RandomizedSearchCV
2  from sklearn.neural_network import MLPClassifier
3  from sklearn.model_selection import GridSearchCV
4  from sklearn.model_selection import RandomizedSearchCV
5
6
7  nn = MLPClassifier(random_state=0)
8
9  param_nn = {
10      'activation': ['identity', 'logistic', 'tanh', 'relu'],
11      'solver': ['adam', 'lbfgs', 'sgd'],
12      'alpha': [0.001, 0.01],
13      'max_iter': [10000],
14      # 'learning_rate': ['constant', 'invscaling', 'adaptive']
15  }
16
17
18  #apply grid search
19  grid_nn = RandomizedSearchCV(nn, param_nn, cv=5, n_jobs=-1, scoring=ftwo_scorer)
20  grid_nn.fit(X_train, y_train)
21
22  # Mean Cross Validation Score
23  print("Best Mean Cross-validation score: {:.3f}".format(grid_nn.best_score_))
24  print('parameters: ', grid_nn.best_params_)
25  print('Train score: ', grid_nn.score(X_train, y_train))
26  print("Test score: ", grid_nn.score(X_test, y_test))
27
28
29  # let's get the predictions
30  X_test_preds = grid_nn.predict(X_test)
31
32
33
34  pd.DataFrame({'Id': test.Id, 'Target': X_test_preds }).to_csv('solution_base5.csv')
35  print("Done :D")

```

Best Mean Cross-validation score: 0.864

parameters: {'solver': 'adam', 'max_iter': 10000, 'alpha': 0.001, 'activation': 'logistic'}

Train score: 0.913978494623656

Test score: 1.0

Done :D

In [382]:

```

1 #Stacking 1
2 #Stacking of Data Sampling Algorithms(randomforest, xgboost,easyensemble) with
3
4 from sklearn.ensemble import StackingClassifier
5 sclf1 = StackingClassifier(estimators=
6                             [ ('randomforest', ogrid_rf.best_estimator_),
7                               ('xgboost', ogrid_xgb.best_estimator_),
8                               ('easyensemble', ogrid_ee.best_estimator_),
9                               ], final_estimator=LogisticRegression())
10 sclf1_param = {
11                 'stack_method':['auto', 'predict_proba']
12             }
13
14 sclf1_grid = GridSearchCV(sclf1, sclf1_param,cv=5, n_jobs=-1, scoring=ftwo_score)
15 sclf1_grid.fit(X_train,y_train)
16 print(f'Best Mean Cross Validation Score is {sclf1_grid.best_score_}')
17 print(f'Best param {sclf1_grid.best_params_}')
18 print(f'Train score is {sclf1_grid.score(X_train,y_train)}')
19 print(f'Test score is {sclf1_grid.score(X_test,y_test)}')
```

Best Mean Cross Validation Score is 0.6023391812865496

Best param {'stack_method': 'auto'}

Train score is 0.8695652173913043

Test score is 1.0

In [378]:

```

1 #Stacking 2
2 #Stacking of cost sensitive algorithms(logistic reg, decision tree, simple svm)
3
4 from sklearn.ensemble import StackingClassifier
5 sclf2 = StackingClassifier(estimators=
6                             [ ('logreg', clogreg_grid.best_estimator_),
7                               ('dtree', cgrid_dtree.best_estimator_),
8                               ('simplesvm', cgrid_svc.best_estimator_),
9                               ], final_estimator=XGBClassifier())
10 sclf2_param = {
11                 'logreg__C':[0.01],
12                 'dtree__max_depth': [1,2,3],
13                 'simplesvm__penalty':['l1','l2'],
14             }
15
16 sclf2_grid = GridSearchCV(sclf2, sclf2_param,cv=5, n_jobs=-1, scoring=ftwo_score)
17 sclf2_grid.fit(X_train,y_train)
18 print(f'Best Mean Cross Validation Score is {sclf2_grid.best_score_}')
19 print(f'Best param {sclf2_grid.best_params_}')
20 print(f'Train score is {sclf2_grid.score(X_train,y_train)}')
21 print(f'Test score is {sclf2_grid.score(X_test,y_test)}')
```

Best Mean Cross Validation Score is 0.906265664160401

Best param {'dtree__max_depth': 1, 'logreg__C': 0.01, 'simplesvm__penalty': 'l1'}

Train score is 0.9042553191489363

Test score is 1.0

In [389]:

```

Stacking 3
Stacking of basic Algorithms(knn, gradient boost and svm) with XGB as my final estimator
3
4
5
from sklearn.ensemble import StackingClassifier
clf3 = StackingClassifier(estimators=
8     [ ('knn', grid_knn.best_estimator_),
9       ('gbc', gbc_grid.best_estimator_),
10      ('dtree', grid_dtree.best_estimator_),
11      ], final_estimator=XGBClassifier(random_state=42,early_s
clf3_param = {
13     # 'final_estimator__C' : [0.01,0.1],
14     # 'knn__n_neighbors': [2,3,5],
15     # 'gbc__learning_rate' : [0.01,0.1]
16     'stack_method':['predict_proba']
17
18
clf3_grid = GridSearchCV(clf3, clf3_param,cv=5, n_jobs=-1, scoring=ftwo_scorer)
clf3_grid.fit(X_train,y_train)
print(f'Best Mean Cross Validation Score is {clf3_grid.best_score_}')
print(f'Best param {clf3_grid.best_params_}')
print(f'Train score is {clf3_grid.score(X_train,y_train)}')
print(f'Test score is {clf3_grid.score(X_test,y_test)}')

```

Best Mean Cross Validation Score is 0.8236842105263158
 Best param {'stack_method': 'predict_proba'}
 Train score is 0.9895833333333334
 Test score is 1.0

In [356]:

```

1  #Saving and uploading model using pickle for logistic regression model
2
3  import pickle
4
5  #logistic regression
6  from sklearn.linear_model import LogisticRegression
7  from sklearn.model_selection import GridSearchCV
8
9  logreg=LogisticRegression(max_iter=1000)
10 logreg_param= {
11     'C': [0.001,0.01,0.1,10,100],
12     'penalty':['l1', 'l2']
13 }
14 logreg_grid = GridSearchCV(logreg,logreg_param,cv=5,n_jobs=-1,scoring=ftwo_scorer)
15 logreg_grid.fit(X_train,y_train)
16
17 model = pickle.dumps(logreg_grid)
18 model_from_pickle = pickle.loads(model)
19 trainscore=model_from_pickle.score(X_train,y_train)
20
21 testscore=model_from_pickle.score(X_test,y_test)
22 print (f'train score {trainscore}')
23

```

train score 0.8695652173913043

In [255]:

```
1  #GaussianProcessClassifier
2  from sklearn.gaussian_process import GaussianProcessClassifier
3
4  gpc=GaussianProcessClassifier(random_state=0)
5
6  gpc_param= { 'n_restarts_optimizer':[2,3,4],
7              'max_iter_predict':[150],
8              'warm_start':[True]
9  }
10 #isoforest_grid = RandomizedSearchCV(isoforest, isoforest_param,cv=5,n_jobs=-1,s
11
12 gpc_grid = GridSearchCV(gpc, gpc_param,cv=5,n_jobs=-1,scoring=ftwo_scorer)
13 gpc_grid.fit(X_train,y_train)
14 print(f'Best Mean Cross Validation Score is {gpc_grid.best_score_}')
15 print(f'Best param {gpc_grid.best_params_}')
16 print(f'Train score is {gpc_grid.score(X_train,y_train)}')
17 print(f'Test score is {gpc_grid.score(X_test,y_test)}')
```

Best Mean Cross Validation Score is 0.8048767752715122

Best param {'max_iter_predict': 150, 'n_restarts_optimizer': 2, 'warm_start': True}

Train score is 1.0

Test score is 0.9615384615384615

In []:

1