

▼ Introduction: Business Problem

When a customer orders a product which is not available in the store or out-of-stock temporarily or due to lack of supply, the customer decides to wait until the desired product is available and there is a guaranteed delivery, then this scenario is called Backorder of the specific product. If the backorders are not handled promptly, there is a high chance of losing a customer to its competitor, along with impact on revenue, share market price etc. On the other hand, taking actions to satisfy or reduce backorders puts tremendous amounts of pressure on the different stages of supply chain management. Taking steps to satisfy or reduce backorders will lead to increased labor/production/transport/Warehouse costs etc. Machine Learning can identify patterns related to backorders before the customer orders. With this, the production can adjust to minimise delays for customer service and provide accurate dates to keep the customers informed. This predictive analysis approach gives the company ample time to react and enables them to satisfy the demands of customers and smooth the supply chain process. The case study "Backorder Prediction" deals with predicting the backorders of products by applying Machine Learning techniques to overcome or reduce the cost of backorders. We will identify parts with the highest chances of shortage so that we could present a high opportunity to improve the company's overall performance. This case study deals with investigating Machine Learning classifiers for imbalanced dataset where the chance of an item going into backorder is very rare when compared to items that do not.

Why Machine Learning is needed?

Through Machine Learning classification techniques, supply chain parameters could be analysed for different product lines. This allows Manufacturer/Company to predict and identify the products that would be unavailable in near future and then take the right step to neutralise the shortage. The predictions based on the ML Models improve with time as they get more data to work on. The degree of performance metrics of these predictive models can be improved by using ensembling techniques which uses diverse models combines them to produce improved results.

Problem Statement

Classify the product whether they would go into backorder (Yes or No) based on data like forecast sales, sales quantity, threshold values, past performance etc.

ML Formulation

The task at hand is classifying whether a product will go to backorder or not for a given input data. This is a Binary Classification Problem hence consists of two target values :

- Yes: Represents that product will go to backorder.
- No: Represents that product will not go to backorder

Business Constraints

- No strict latency constraints.
- Misclassification may result in less effective supply chain management systems such as inaccurate demand forecasting and misclassification of backorder products.

Dataset Analysis

Dataset consists of 23 columns namely : SKU : stands for Stock Keeping Unit. It is a unique ID for each row

National_inv : Present Inventory level of the product

Lead time : Transit time of the product which means how long it takes for a shipment to be delivered at its final destination after it has been picked from the start point.

In_transit_qty : it is calculated based on the last picking slip or based on cumulative quantity. In simple words it is the amount of product in transit.

Forecast_3_months, Forecast_6_months, Forecast_9_months : Forecast of sales of the product for next 3, 6, 9 months respectively.

Sales_1_month, Sales_3_month, Sales_6_month : Sales of product in last 1, 3, 6, 9 months respectively.

Min_bank : Minimum amount of stock recommended.

potential_issue : Problem/issue identified in the product/part

Pieces_past_due : Amount of parts of the product overdue if any.

Perf_6_months_avg, perf_12_months_avg: Product performance over past 6 months and 12 months respectively.

local_bo_qty : Amount of stock orders overdue

deck_risk, oe_constraint, ppap_risk, stop_auto_buy, rev_stop : Yes or No flags set for the products

went_to_backorder: Target Variable

Performance Metrics

In the case of Backorders, it happens very rarely that the products go to backorder. So in this type of problems, Accuracy is not a useful metric since it works well with a Balanced dataset. Since the dataset is imbalanced, the ideal performance metric would be Precision, Recall and F1 score, especially not to miss failure cases. F1 Score might be a better measure to use if we need to seek balance between Precision and Recall and also given that there is uneven class distribution (large number of Actual Negatives).

$$F1_score = 2 * ((Precision * Recall) / (Precision + Recall))$$

Precision : This can be thought of as , Out of postively predicted values, how many are actually positive.It is calculated as

$$\text{Precision} = \text{True Positive} / \text{True Postive} + \text{False Positive}$$

Recall : Recall can be thought of as, out of actual positives how many were detected by our model.It is calculated as

$$\text{Recall} = \text{True Positive} / \text{True Positive} + \text{False Negative}$$

Terms used in Precision and Recall are :

- True Positive : Product detected as Backorder and it turns out to be True.Benefits by predicting correctly the Backorders as Profit generated from such items is benefit.
- False Positive : Product detected as Backorder and it Non-Backorder product.This would cost the company as the fact that we predicted few items as Backorder items but they were not actually in Backorder list. The warehousing cost for such items is the cost associated with the false positives.
- False Negative : A miss. The product is missed to be detected by the model and product stays as Non-Backorder.This will cost the company as cost associated with incorrectly missing items when actual demand was there for them.

F1 Score will be used when False Positives and False Negatives are more important as they are crucial for Bussiness cost.

Among F1 Score, we will choose Macro Averaged F1 Score since we would like to treat both class as equals. Micro F1 Score will be used when we want to maximise the classification of a particular class which is not the case in this problem. Due to imbalanced dataset, we will get high Micro Average F1 Score which is not called for.

▼ Import Libraries

```
1 import numpy as np
2 import pandas as pd
3 import matplotlib.pyplot as plt
4 import warnings
5 warnings.filterwarnings('ignore')
6 import seaborn as sns
7 from sklearn.model_selection import train_test_split
8 from sklearn.impute import KNNImputer
9 from sklearn.decomposition import TruncatedSVD
10 from sklearn.preprocessing import MinMaxScaler
11 import tensorflow as tf
12 from pandas import read_csv
```

```
13 from numpy.random import seed
14 from keras.layers import Input, Dense
15 from keras.models import Model
16 from tensorflow import keras
17 from tensorflow.keras import layers
18
19 from sklearn.experimental import enable_iterative_imputer
20 from sklearn.impute import IterativeImputer
21
22 from tensorflow.keras.layers import Dense, Input, Conv2D, MaxPool2D, Activation, Dropout, F
23 from tensorflow.keras.models import Model
24 import random as rn
25 from keras.models import Sequential
26 from keras.layers import LeakyReLU
27 from keras.layers import BatchNormalization
28 from tensorflow.keras.optimizers import Adam
29
30 from tensorflow.keras.callbacks import ReduceLROnPlateau
31 from tensorflow.keras.callbacks import EarlyStopping
32 from tensorflow.keras.callbacks import ModelCheckpoint
33
34 import imblearn
35 from imblearn.over_sampling import SMOTE
36 from imblearn.over_sampling import SMOTENC
37
38 from sklearn.metrics import confusion_matrix, ConfusionMatrixDisplay
39 from sklearn.metrics import f1_score
40
41 from scipy.stats import chi2_contingency
42 from scipy.stats import chi2
43 import scipy.stats as stats
44 from scipy.stats import f_oneway
45 from scipy import stats
46 from xgboost import XGBClassifier
47 from sklearn.model_selection import RandomizedSearchCV
48 from sklearn.tree import DecisionTreeClassifier
49 from sklearn.calibration import CalibratedClassifierCV
50 from sklearn.neighbors import KNeighborsClassifier
51 from sklearn.metrics import log_loss
52 from sklearn.metrics import confusion_matrix
53 from sklearn.model_selection import train_test_split
54 from sklearn.linear_model import LogisticRegression
55 from sklearn.ensemble import RandomForestClassifier
56
57 from imblearn.under_sampling import RandomUnderSampler
58 from xgboost import XGBClassifier
59
60 from sklearn.ensemble import AdaBoostClassifier
61 from sklearn.tree import DecisionTreeRegressor
62 from sklearn.model_selection import GridSearchCV
63 import random
64 import pickle
65 from joblib import dump, load
```

```
1 from google.colab import drive
2 drive.mount('/content/drive')
```

Mounted at /content/drive

▼ Load Dataset

```
1 df_train = pd.read_csv("/content/drive/MyDrive/Kaggle_Training_Dataset_v2.csv")
2 df_test = pd.read_csv("/content/drive/MyDrive/Kaggle_Test_Dataset_v2.csv")
```

▼ Count rows and columns for train and test data

```
1 print("Number of rows in train data : ",df_train.shape[0])
2 print("Number of rows in test data : ",df_test.shape[0])
3 print("Number of columns in train data : ",df_train.shape[1])
4 print("Number of columns in test data : ",df_test.shape[1])
```

```
Number of rows in train data : 1687861
Number of rows in test data : 242076
Number of columns in train data : 23
Number of columns in test data : 23
```

```
1 # Removing the last rows of Train as well as Test data which contains Nan For all col
2 df_train.drop(df_train.tail(1).index,inplace=True)
3 df_test.drop(df_test.tail(1).index,inplace=True)
```

▼ Merge Train and Test Data

```
1 df = df_train.append(df_test,ignore_index=True)
```

```
1 df.shape

(1929935, 23)
```

```
1 df.describe()
```

	national_inv	lead_time	in_transit_qty	forecast_3_month	forecast_6_mont
count	1.929935e+06	1.814318e+06	1.929935e+06	1.929935e+06	1.929935e+0
mean	4.965683e+02	7.878627e+00	4.306440e+01	1.785399e+02	3.454659e+0
std	2.957343e+04	7.054212e+00	1.295420e+03	5.108770e+03	9.831562e+0
min	-2.725600e+04	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+0

1

50% 1.500000e+04 0.000000e+00 0.000000e+00 0.000000e+00 0.000000e+0

► Check for Unique rows

[] ↳ 2 cells hidden

► Check for Nan Values

[] ↳ 1 cell hidden

Comment

- lead_time column has 115617 Nan values. Need to perform imputation because of this.

► Count of Target values of Imbalanced dataset

[] ↳ 5 cells hidden

► Exploratory Data Analysis

[] ↳ 114 cells hidden

► Finding Correlation between 2 Categorical Features

[] ↳ 5 cells hidden

► Finding Correlation between Categorical and Numerical Feature using Point Biserial Correlation

The Point Biserial Correlation is used to measure the correlation between a Categorical Variable(Binary Category) and Continous Variable

The Correlation coeffecients varies between -1 to +1 with 0 implying No Correlation.

Null Hypothesis : There is no correlation between the two features.

Alternate Hypothesis : There is a correlation between the two features.

We can verify the hypothesis using **p-value**.

[] ↳ 7 cells hidden

▼ Data Preprocessing

▼ Converting Target Variable to 0 and 1

```
1 #Converting Target Variable to 0 and 1
2 df['went_on_backorder'].replace({'Yes':1,'No':0},inplace=True)
3 df['went_on_backorder'].astype(int)
```

```
0      0
1      0
2      0
3      0
4      0
..
1929930 0
1929931 0
1929932 0
1929933 0
1929934 0
Name: went_on_backorder, Length: 1929935, dtype: int64
```

▼ Check for Number of 0s in Numerical Features and drop columns if 0s are more than 95%.

```
1 # Drop Columns if 95% values are 0
2 zero_columns = []
3 for col in df.columns:
4     if col!='went_on_backorder' and col!='sku':
5         value = df[col].value_counts()
6         if 0 in value:
7             if value[0] >= (0.95*df.shape[0]):
8                 zero_columns.append(col)
9                 print(col,"-> percentage of zeros : ",(value[0]*100/df.shape[0]),"%")

pieces_past_due -> percentage of zeros :  98.53347392528764 %
local_bo_qty -> percentage of zeros :  98.64223406487783 %

1 df = df.drop(columns=zero_columns,axis=1)
```

▼ Replacing -99 by Nan in performance column

```

1 # replacing -99 by Nan in performance column
2 df.perf_6_month_avg.replace({-99.0 : np.nan},inplace=True)
3 df.perf_12_month_avg.replace({-99.0 : np.nan},inplace=True)

```

▼ Converting Categorical Features

```

1 categorical_columns = ['rev_stop','stop_auto_buy','ppap_risk','oe_constraint','deck_r
2 for col in categorical_columns:
3     df[col].replace({'Yes':1,'No':0},inplace=True)
4     df[col]=df[col].astype(int)

```

▼ Removing Outlier Datapoints from DataFrame

```

1 old_shape = df.shape[0]

1 df = df[(df.national_inv >= 0.000) & (df.national_inv <= 5487.000) & (df.in_transit_q
2     (df.forecast_3_month <= 2280.000) & (df.forecast_6_month <= 4335.659999999916) &\
3     (df.forecast_9_month <= 6316.000) & (df.sales_1_month <= 693.000) & (df.sales_3_mon
4     (df.sales_6_month <= 4410.000) & (df.sales_9_month <= 6698.000) & (df.min_bank <= 6

1 new_shape = df.shape[0]

1 print("Number of Outliers removed : ",old_shape-new_shape)

Number of Outliers removed : 48564

1 df.describe()

```

	national_inv	lead_time	in_transit_qty	forecast_3_month	forecast_6_month
count	1.881371e+06	1.768920e+06	1.881371e+06	1.881371e+06	1.881371e+06
mean	1.376750e+02	7.879902e+00	1.145663e+01	3.950405e+01	7.661432e+00
std	4.669358e+02	7.043846e+00	7.915361e+01	1.639866e+02	3.039033e+01
min	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00
25%	4.000000e+00	4.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00
50%	1.400000e+01	8.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00
75%	7.200000e+01	9.000000e+00	0.000000e+00	3.000000e+00	1.000000e+00
max	5.487000e+03	5.200000e+01	5.300000e+03	2.280000e+03	4.320000e+03

▼ Train Test Split

```

1  # Assigning the Target Variable Column to y_true variable and dropping it from the DF
2  y_true = df['went_on_backorder']

1  df = df.drop(['sku', 'went_on_backorder'], axis=1)

1  X_train, X_test, y_train, y_test = train_test_split(df, y_true, stratify=y_true, test_size=
2  X_train, X_cv, y_train, y_cv = train_test_split(X_train, y_train, stratify=y_train, test_si
3  print(X_train.shape, X_cv.shape, X_test.shape)

(1204076, 19) (301020, 19) (376275, 19)

```

▼ Median Imputation

```

1  median_values = X_train.median()
2
3  X_train_median = X_train.fillna(median_values)
4  X_cv_median = X_cv.fillna(median_values)
5  X_test_median = X_test.fillna(median_values)
6
7  print(X_train_median.shape, X_cv_median.shape, X_test_median.shape)

(1204076, 19) (301020, 19) (376275, 19)

```

▼ Feature Engineering

▼ Adding SVD to Median Imputation

```

1  Trunc_SVD = TruncatedSVD(n_components=2, n_iter=20)
2  X_train_median_SVD = Trunc_SVD.fit_transform(X_train_median)
3  X_cv_median_SVD = Trunc_SVD.transform(X_cv_median)
4  X_test_median_SVD = Trunc_SVD.transform(X_test_median)
5  print(X_train_median_SVD.shape, X_cv_median_SVD.shape, X_test_median_SVD.shape)

(1204076, 2) (301020, 2) (376275, 2)

```

▼ AutoEncoder for Feature Engineering

```

1  tf.keras.backend.clear_session()
2  np.random.seed(0)
3  rn.seed(0)

```

```

1  earlystop = EarlyStopping(monitor='val_loss', min_delta=0.1, patience=3, verbose=1)
2  reduce_lr_on_val_loss = ReduceLROnPlateau(monitor='val_loss', factor=0.8, patience=1, vo

```

```

2 reduce_lr_on_val_loss = ReduceLROnPlateau(monitor= val_loss , factor=0.5, patience=1, ve
3
4 filepath="model_save_median/weights-{epoch:02d}-{val_loss:.4f}.hdf5"
5 checkpoint = ModelCheckpoint(filepath=filepath, monitor='val_loss', verbose=1, save_b

1 #https://machinelearningmastery.com/autoencoder-for-classification/
2 number_of_columns = X_train_median.shape[1]
3 encoder_input_dim = Input(shape=(number_of_columns,))
4
5 #encoder level 1
6 e = Dense(number_of_columns)(encoder_input_dim)
7 e = BatchNormalization()(e)
8 e = LeakyReLU()(e)
9
10 # encoder level 2
11 e = Dense(10)(e)
12 e = BatchNormalization()(e)
13 e = LeakyReLU()(e)
14
15 #Bottleneck
16 n_bottleneck = 2
17 bottleneck = Dense(n_bottleneck)(e)

1 d = Dense(10)(bottleneck)
2 d = BatchNormalization()(d)
3 d = LeakyReLU()(d)
4
5 # decoder level 2
6 d = Dense(number_of_columns)(d)
7 d = BatchNormalization()(d)
8 d = LeakyReLU()(d)
9 # output
10 encoder_output = Dense(number_of_columns,activation='linear')(d)
11 # Defining auto encoder model
12 model = Model(inputs=encoder_input_dim,outputs = encoder_output)
13
14 callback_list = [earlystop,reduce_lr_on_val_loss,checkpoint]
15 model.compile(optimizer='adam',loss='mse')
16 model.fit(X_train_median,y_train,epochs=25,batch_size=100,shuffle=True,validation_dat
17

Epoch 1/25
12041/12041 [=====] - 37s 3ms/step - loss: 0.6548 - val_loss:

Epoch 00001: val_loss improved from inf to 0.00688, saving model to model_save_mediar
Epoch 2/25
12041/12041 [=====] - 34s 3ms/step - loss: 0.0069 - val_loss:

Epoch 00002: val_loss improved from 0.00688 to 0.00643, saving model to model_save_me
Epoch 3/25
12041/12041 [=====] - 34s 3ms/step - loss: 0.0063 - val_loss:

Epoch 00003: ReduceLROnPlateau reducing learning rate to 0.0009000000427477062.

Epoch 00003: val_loss did not improve from 0.00643

```

Epoch 4/25

12041/12041 [=====] - 34s 3ms/step - loss: 0.0062 - val_loss:

Epoch 00004: val_loss improved from 0.00643 to 0.00612, saving model to model_save_m

Epoch 00004: early stopping

<tensorflow.python.keras.callbacks.History at 0x7f6338198790>



```
1 encoder = Model(inputs =encoder_input_dim,outputs = bottleneck)
2 encoder.save('encoder.h5')
```

WARNING:tensorflow:Compiled the loaded model, but the compiled metrics have yet to be



```
1 encoded_output_Train = encoder.predict(X_train_median)
2 encoded_output_cv = encoder.predict(X_cv_median)
3 encoded_output_test = encoder.predict(X_test_median)
```

```
1 print(encoded_output_Train.shape,encoded_output_cv.shape,encoded_output_test.shape)

(1204076, 2) (301020, 2) (376275, 2)
```

► Feature Binning using Decision Tree

[] ↳ 11 cells hidden

► Adding the Feature Engineered Data to the Main DataFrame

[] ↳ 5 cells hidden

► Over Sampling Technique : SMOTE-NC

[] ↳ 4 cells hidden

► Normalise All Columns after SMOTE

[] ↳ 1 cell hidden

► Under Sampling

[] ↳ 7 cells hidden

Key points and takeways from EDA and Feature Engineering

- We are solving binary classification problem with very high data imbalance. Positive class being the majority.
- Categorical features consists of Yes and No Category.
- SKU is unique and is not important for classification and hence will be dropped.
- All numerical features had extreme skewness on Right indicating them as Outliers.
- in Lead_time feature there are many Nan values. performed imputations on that feature.
- Performance average features had -99 as their value. implemented Median imputation.
- From the Correlation Matrix The Forecast, performance average and sales features are extremely correlated to each other.
- Performed Chi Square test on 2 categorical features. Made Null hypothesis for this correlation and performed statistical tests. Results state that there is no correlation between rev_stop feature and the target variable.
- Performed Point Biserial correlation test to check for correlation between categorical and numerical features. Forecast features , sales 1 month and national_inv feature are not correlated to the target variable.
- local_bo_qty and piecies_past_due Features have close to 98% of 0s in them. They have been removed.
- Removed close to 48k outlier datapoints as a part of data cleaning process.
- Performed SMOTE-NC as we have categorical data as well , an oversampling technique to increase minority class duplicate datapoints to increase percentage of minority class.
- Used RandomUnderSampler as an undersampling technique to reduce the datapoints from majority class.
- Apart from given features, also added 2-2 Truncated SVD features, AutoEncoder Features, Discretised bin features as a feature extraction method.

► Random Model

[] ↳ 1 cell hidden

Machine Learning Models

► Logistic Regression

[] ↳ 15 cells hidden

► Random Forest

[] ↳ 17 cells hidden

► XGBOOST

[] ↳ 31 cells hidden

▼ Ensemble Model

▼ Without Feature Engineering

```

1  # Assigning the Target Variable Column to y_true variable and dropping it from the DF
2  y_true = df['went_on_backorder']

1  df = df.drop(['sku', 'went_on_backorder'], axis=1)

1  df.shape

(1881371, 19)

1  X_train_en, X_test_en, y_train_en, y_test_en = train_test_split(df, y_true, stratify=y_true,
2  print(X_train_en.shape, X_test_en.shape, y_train_en.shape, y_test_en.shape)

(1505096, 19) (376275, 19) (1505096,) (376275,)

1  data_d1, data_d2, y_d1, y_d2 = train_test_split(X_train_en, y_train_en, stratify=y_train_e
2  print(data_d1.shape, data_d2.shape, y_d1.shape, y_d2.shape)

(752548, 19) (752548, 19) (752548,) (752548,)

1  median_values = data_d1.median()
2
3  columns = data_d1.columns
4
5  data_d1 = data_d1.fillna(median_values)
6  data_d2 = data_d2.fillna(median_values)
7  X_test_en = X_test_en.fillna(median_values)
8
9  print(data_d1.shape, data_d2.shape, X_test_en.shape)

(752548, 19) (752548, 19) (376275, 19)

1  data_d1 = pd.DataFrame(data_d1, columns=columns)
2  data_d2 = pd.DataFrame(data_d2, columns=columns)
3  X_test_en = pd.DataFrame(X_test_en, columns=columns)

1  # Normalise Median Imputed Data
2  MinMaxSc = MinMaxScaler()

```

```
3
4 columns = data_d1.columns
5 data_d1 = pd.DataFrame(MinMaxSc.fit_transform(data_d1),columns=columns)
6 data_d2 = pd.DataFrame(MinMaxSc.transform(data_d2),columns=columns)
7 X_test_en = pd.DataFrame(MinMaxSc.transform(X_test_en),columns=columns)
8
9 print(data_d1.shape,data_d2.shape,X_test_en.shape)

(752548, 19) (752548, 19) (376275, 19)

1 def generating_samples_with_replacement(input_data,output_data):
2     rows_selected = random.sample(range(0,len(input_data)-1),int(0.6*(len(input_data)
3     sampled_input = input_data.iloc[rows_selected]
4     sampled_output = output_data.iloc[rows_selected]
5
6     sampled_input = sampled_input.values.tolist()
7     sampled_output = sampled_output.values.tolist()
8
9     replicate = random.sample(range(0,len(sampled_input)-1),int(0.4*(len(input_data)))
10    for num in replicate:
11        sampled_input.append(sampled_input[num])
12        sampled_output.append(sampled_output[num])
13
14    return sampled_input,sampled_output

1 def baseModel_DecisionTree(number_of_models,data_d1,y_d1,data_d2,y_d2):
2     input_data_list=[]
3     output_data_list=[]
4     clf_list=[]
5     output_list=[]
6     final_data = []
7
8     print("Generating samples with 60% sample and 40% sample with replacement")
9     print("In Sample Generate")
10    for num in range(number_of_models):
11        print("Sample generation for Model : ",num)
12        inp,outp = generating_samples_with_replacement(data_d1,y_d1)
13        input_data_list.append(inp)
14        output_data_list.append(outp)
15    print("Out Sample Generate")
16
17    print("Fit Decision Tree Models on samples generated")
18    print("In Fitting Models")
19    for num in range(number_of_models):
20        print("Fitting for Model : ",num)
21        DT = DecisionTreeRegressor(max_depth=None)
22        clf = DT.fit(input_data_list[num],output_data_list[num])
23        clf_list.append(clf)
24    print("Out Fitting Models")
25
26    print("Predict values for second dataset")
27    print("In Predict")
28    for num in range(number_of_models):
29        print("Predict for Model : ",num)
```

```

30         clf = clf_list[num]
31         output = clf.predict(data_d2)
32         output_list.append(output)
33     print("Out Predict")
34
35     print("Zip each row of the column to form dataset")
36     print(" In form Meta Data")
37     for num in range(len(output_list[0])):
38         output=[]
39         for i in range(number_of_models):
40             output.append(output_list[i][num])
41         final_data.append(output)
42     print("Out form Meta Data")
43
44     final_data = pd.DataFrame(final_data)
45
46     return clf_list,final_data

1  number_of_models=15
2  modellist,data = BaseModel_DecisionTree(number_of_models,data_d1,y_d1,data_d2,y_d2)

```

Generating samples with 60% sample and 40% sample with replacement

In Sample Generate

Sample generation for Model : 0
Sample generation for Model : 1
Sample generation for Model : 2
Sample generation for Model : 3
Sample generation for Model : 4
Sample generation for Model : 5
Sample generation for Model : 6
Sample generation for Model : 7
Sample generation for Model : 8
Sample generation for Model : 9
Sample generation for Model : 10
Sample generation for Model : 11
Sample generation for Model : 12
Sample generation for Model : 13
Sample generation for Model : 14

Out Sample Generate

Fit Decision Tree Models on samples generated

In Fitting Models

Fitting for Model : 0
Fitting for Model : 1
Fitting for Model : 2
Fitting for Model : 3
Fitting for Model : 4
Fitting for Model : 5
Fitting for Model : 6
Fitting for Model : 7
Fitting for Model : 8
Fitting for Model : 9
Fitting for Model : 10
Fitting for Model : 11
Fitting for Model : 12
Fitting for Model : 13
Fitting for Model : 14

Out Fitting Models

Predict values for second dataset

```

In Predict
Predict for Model : 0
Predict for Model : 1
Predict for Model : 2
Predict for Model : 3
Predict for Model : 4
Predict for Model : 5
Predict for Model : 6
Predict for Model : 7
Predict for Model : 8
Predict for Model : 9
Predict for Model : 10
Predict for Model : 11
Predict for Model : 12
Predict for Model : 13
Predict for Model : 14
Out Predict
Zip each row of the column to form dataset
In form Meta Data
Out form Meta Data

```

1 # Meta Classifier - Logistic Regression

```

1 lg_clf = LogisticRegression()
2 params = {'C':[10**x for x in range(-5,4)]}
3 grid_log_clf = GridSearchCV(lg_clf,param_grid=params,n_jobs=-1,scoring='f1_macro',ver
4 grid_log_clf.fit(data,y_d2.values.ravel())

```

```

Fitting 5 folds for each of 9 candidates, totalling 45 fits
[Parallel(n_jobs=-1)]: Using backend LokyBackend with 2 concurrent workers.
[Parallel(n_jobs=-1)]: Done 28 tasks      | elapsed: 27.9s
[Parallel(n_jobs=-1)]: Done 45 out of 45 | elapsed: 45.9s finished
GridSearchCV(cv=None, error_score=nan,
              estimator=LogisticRegression(C=1.0, class_weight=None, dual=False,
              fit_intercept=True,
              intercept_scaling=1, l1_ratio=None,
              max_iter=100, multi_class='auto',
              n_jobs=None, penalty='l2',
              random_state=None, solver='lbfgs',
              tol=0.0001, verbose=0,
              warm_start=False),
              iid='deprecated', n_jobs=-1,
              param_grid={'C': [1e-05, 0.0001, 0.001, 0.01, 0.1, 1, 10, 100,
              1000]}},
              pre_dispatch='2*n_jobs', refit=True, return_train_score=True,
              scoring='f1_macro', verbose=3)

```

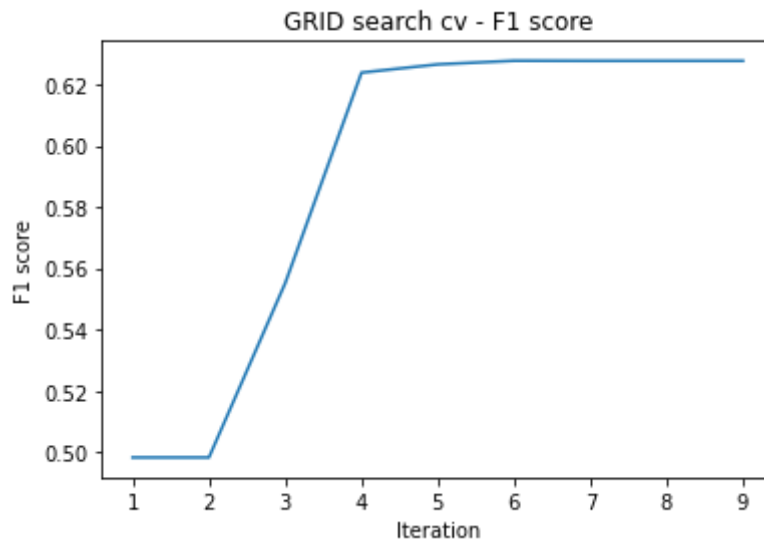
```

1 labels = grid_log_clf.cv_results_['params']
2 x_axis = range(1,10)
3 y_axis = grid_log_clf.cv_results_['mean_test_score']
4 for i,label in enumerate(labels):
5     print(label,":",y_axis[i])
6
7 plt.xlabel("Iteration")
8 plt.ylabel("F1 score")
9 plt.title("GRID search cv - F1 score")
10 plt.plot(x_axis,y_axis)

```



```
{'C': 1e-05} : 0.4982939749796286
{'C': 0.0001} : 0.4982939749796286
{'C': 0.001} : 0.5551692329818663
{'C': 0.01} : 0.6238061616131717
{'C': 0.1} : 0.6264931744337635
{'C': 1} : 0.6276808492590912
{'C': 10} : 0.6276425035563356
{'C': 100} : 0.6276425035563356
{'C': 1000} : 0.6276425035563356
[<matplotlib.lines.Line2D at 0x7f47a58af4d0>]
```



```
1 grid_log_clf.best_estimator_
```

```
LogisticRegression(C=1, class_weight=None, dual=False, fit_intercept=True,
                    intercept_scaling=1, l1_ratio=None, max_iter=100,
                    multi_class='auto', n_jobs=None, penalty='l2',
                    random_state=None, solver='lbfgs', tol=0.0001, verbose=0,
                    warm_start=False)
```

```
1 clf_best = LogisticRegression(C=1)
2 clf_best.fit(data,y_d2.values.ravel())
```

```
LogisticRegression(C=1, class_weight=None, dual=False, fit_intercept=True,
                    intercept_scaling=1, l1_ratio=None, max_iter=100,
                    multi_class='auto', n_jobs=None, penalty='l2',
                    random_state=None, solver='lbfgs', tol=0.0001, verbose=0,
                    warm_start=False)
```

```
1 def perform_test(number_of_models,clf_list,X_test_en):
2     test_output_list=[]
3     test_data = []
4
5     print("In Predicting Output for Test Data")
6     for num in range(number_of_models):
7         clf = clf_list[num]
8         output = clf.predict(X_test_en)
9         test_output_list.append(output)
10    print("Out Predicting Output for Test Data")
11
12    for i in range(len(test output list[0])):
```

```

13     outp=[]
14     for j in range(number_of_models):
15         outp.append(test_output_list[j][i])
16     test_data.append(outp)
17     test_data = pd.DataFrame(test_data)
18
19     return test_data

```

```

1 test_data = perform_test(number_of_models,modelList,X_test_en)
2 y_test_pred = clf_best.predict(test_data)

```

In Predicting Output for Test Data
Out Predicting Output for Test Data

```
1 print("Test F1 score : ",f1_score(y_test_en,y_test_pred,average='macro'))
```

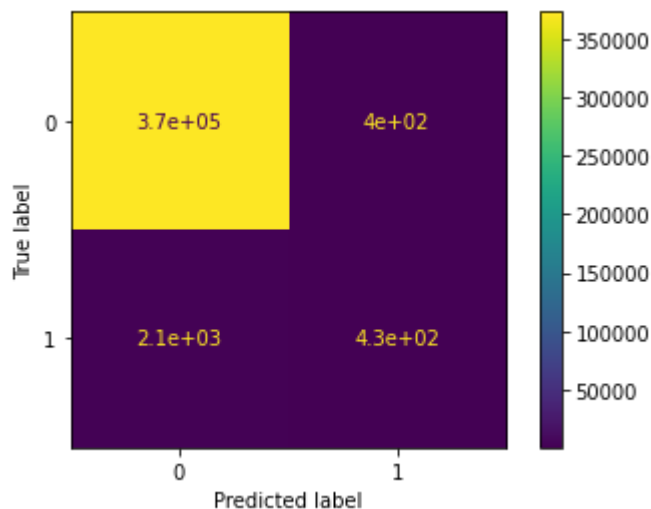
Test F1 score : 0.6361531436954092

```

1 confusionMatrix = confusion_matrix(y_test_en,y_test_pred,labels=[0,1])
2 cm_display = ConfusionMatrixDisplay(confusion_matrix = confusionMatrix,display_labels
3 cm_display.plot()

```

<sklearn.metrics._plot.confusion_matrix.ConfusionMatrixDisplay at 0x7f2c03eb9a90>



1

▼ With Feature Engineering

```

1 # Assigning the Target Variable Column to y_true variable and dropping it from the DF
2 y_true = df['went_on_backorder']

```

```
1 df = df.drop(['sku', 'went_on_backorder'],axis=1)
```

```
1 df.shape
```

```
(1881371, 19)
```

```
1 X_train_en,X_test_en,y_train_en,y_test_en=train_test_split(df,y_true,stratify=y_true,
2 print(X_train_en.shape,X_test_en.shape,y_train_en.shape,y_test_en.shape)
```

```
(1505096, 19) (376275, 19) (1505096,) (376275,)
```

```
1 data_d1,data_d2,y_d1,y_d2 = train_test_split(X_train_en,y_train_en,stratify=y_train_e
2 print(data_d1.shape,data_d2.shape,y_d1.shape,y_d2.shape)
```

```
(752548, 19) (752548, 19) (752548,) (752548,)
```

```
1 median_values = data_d1.median()
2
3 columns = data_d1.columns
4
5 data_d1 = data_d1.fillna(median_values)
6 data_d2 = data_d2.fillna(median_values)
7 X_test_en = X_test_en.fillna(median_values)
8
9 print(data_d1.shape,data_d2.shape,X_test_en.shape)
```

```
(752548, 19) (752548, 19) (376275, 19)
```

```
1 data_d1 = pd.DataFrame(data_d1,columns=columns)
2 data_d2 = pd.DataFrame(data_d2,columns=columns)
3 X_test_en = pd.DataFrame(X_test_en,columns=columns)
```

▼ SVD

```
1 Trunc_SVD = TruncatedSVD(n_components=2,n_iter=20)
2 data_d1_SVD = Trunc_SVD.fit_transform(data_d1)
3 data_d2_SVD = Trunc_SVD.transform(data_d2)
4 X_test_en_SVD = Trunc_SVD.transform(X_test_en)
5 print(data_d1_SVD.shape,data_d2_SVD.shape,X_test_en_SVD.shape)
```

```
(752548, 2) (752548, 2) (376275, 2)
```

▼ Auto Encoders

```
1 tf.keras.backend.clear_session()
2 np.random.seed(0)
3 rn.seed(0)

1 earlystop = EarlyStopping(monitor='val_loss', min_delta=0.1, patience=3, verbose=1)
2 reduce_lr_on_val_loss = ReduceLROnPlateau(monitor='val_loss',factor=0.9,patience=1,ve
3
4 filepath="model_save_median/weights-{epoch:02d}-{val_loss:.4f}.hdf5"
5 checkpoint = ModelCheckpoint(filepath=filepath, monitor='val_loss', verbose=1, save_b
```

```

1  #https://machinelearningmastery.com/autoencoder-for-classification/
2  number_of_columns = data_d1.shape[1]
3  encoder_input_dim = Input(shape=(number_of_columns,))
4
5  #encoder level 1
6  e = Dense(number_of_columns)(encoder_input_dim)
7  e = BatchNormalization()(e)
8  e = LeakyReLU()(e)
9
10 # encoder level 2
11 e = Dense(10)(e)
12 e = BatchNormalization()(e)
13 e = LeakyReLU()(e)
14
15 #Bottleneck
16 n_bottleneck = 2
17 bottleneck = Dense(n_bottleneck)(e)

1  d = Dense(10)(bottleneck)
2  d = BatchNormalization()(d)
3  d = LeakyReLU()(d)
4
5  # decoder level 2
6  d = Dense(number_of_columns)(d)
7  d = BatchNormalization()(d)
8  d = LeakyReLU()(d)
9  # output
10 encoder_output = Dense(number_of_columns,activation='linear')(d)
11 # Defining auto encoder model
12 model = Model(inputs=encoder_input_dim,outputs = encoder_output)
13
14 callback_list = [earlystop,reduce_lr_on_val_loss,checkpoint]
15 model.compile(optimizer='adam',loss='mse')
16 model.fit(data_d1,y_d1,epochs=25,batch_size=100,shuffle=True,validation_data=(data_d2
17

```

Epoch 1/25

7526/7526 [=====] - 28s 3ms/step - loss: 0.6895 - val_loss:

Epoch 00001: val_loss improved from inf to 0.00816, saving model to model_save_mediar

Epoch 2/25

7526/7526 [=====] - 27s 4ms/step - loss: 0.0079 - val_loss:

Epoch 00002: val_loss improved from 0.00816 to 0.00681, saving model to model_save_me

Epoch 3/25

7526/7526 [=====] - 24s 3ms/step - loss: 0.0075 - val_loss:

Epoch 00003: ReduceLROnPlateau reducing learning rate to 0.0009000000427477062.

Epoch 00003: val_loss improved from 0.00681 to 0.00676, saving model to model_save_me

Epoch 4/25

7526/7526 [=====] - 25s 3ms/step - loss: 0.0067 - val_loss:

Epoch 00004: val_loss improved from 0.00676 to 0.00669, saving model to model_save_me

Epoch 00004: early stopping
 <tensorflow.python.keras.callbacks.History at 0x7f3a3fde1090>

```
1 encoder = Model(inputs =encoder_input_dim,outputs = bottleneck)
2 encoder.save('encoder.h5')
```

WARNING:tensorflow:Compiled the loaded model, but the compiled metrics have yet to be

```
1 encoded_output_data_d1 = encoder.predict(data_d1)
2 encoded_output_data_d2 = encoder.predict(data_d2)
3 encoded_output_test = encoder.predict(X_test_en)

1 print(encoded_output_data_d1.shape,encoded_output_data_d2.shape,encoded_output_test.s

(752548, 2) (752548, 2) (376275, 2)
```

Feature Binning

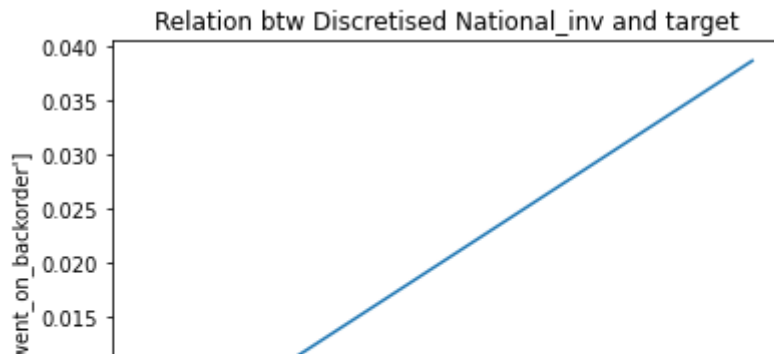
```
1 # Bins and Probabilities for National_inv Feature
2 clf_inv = DecisionTreeClassifier(max_depth=2)
3 clf_inv.fit(data_d1.national_inv.to_frame(),y_d1)
4 data_d1['national_inv_prob'] = clf_inv.predict_proba(data_d1.national_inv.to_frame())
5 print(data_d1.national_inv_prob.head(10))
```

```
844820    0.007317
197031    0.002152
897553    0.002152
1399467   0.002152
1451080   0.007317
178399    0.002152
488875    0.007317
798025    0.007317
839330    0.002152
694483    0.002152
Name: national_inv_prob, dtype: float64
```

```
1 # Adding target variable to check relation
2 data_d1['went_on_backorder'] = y_d1
```

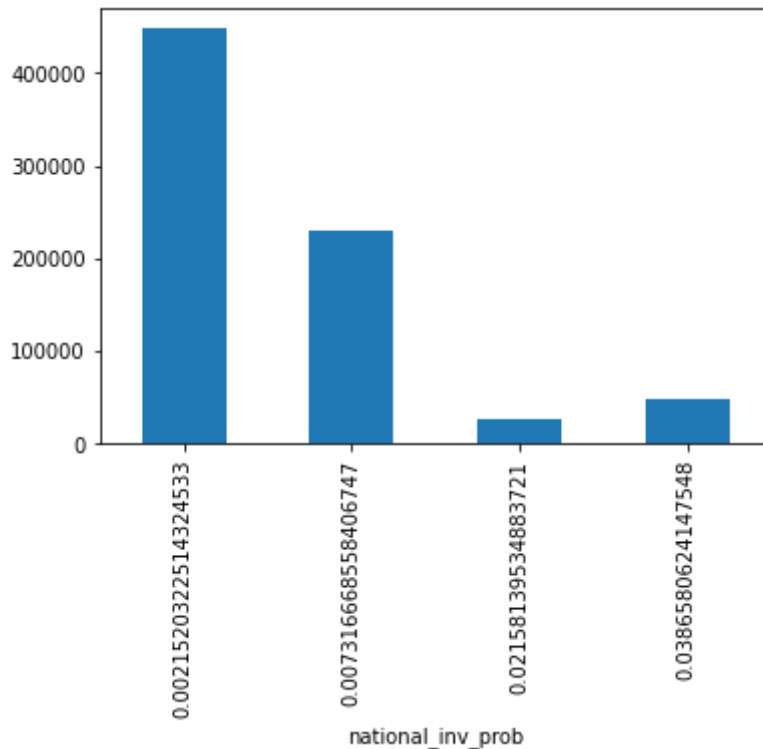
```
1 fig = plt.figure()
2 fig = data_d1.groupby(['national_inv_prob'])['went_on_backorder'].mean().plot()
3 fig.set_title("Relation btw Discretised National_inv and target")
4 fig.set_ylabel(['went_on_backorder'])
```

```
Text(0, 0.5, "['went_on_backorder']")
```



- 1 #check the number of products per probabilistic bin to understand distribution of dis
- 2 data_d1.groupby(['national_inv_prob'])['went_on_backorder'].count().plot.bar()

```
<matplotlib.axes._subplots.AxesSubplot at 0x7f39ec101410>
```



- 1 print("Unique values of National_inv prob : ",data_d1.national_inv_prob.unique())
- Unique values of National_inv prob : [0.00731667 0.00215203 0.03865806 0.0215814]

- 1 pd.concat([data_d1.groupby(['national_inv_prob'])['national_inv'].min(),
- 2 data_d1.groupby(['national_inv_prob'])['national_inv'].max()],axis=1)

	national_inv	national_inv
national_inv_prob		
0.002152	10.0	5486.0
0.007317	2.0	9.0
0.021581	1.0	1.0
0.038658	0.0	0.0

```

1 data_d2['national_inv_prob'] = clf_inv.predict_proba(data_d2.national_inv.to_frame())
2 X_test_en['national_inv_prob'] = clf_inv.predict_proba(X_test_en.national_inv.to_fram

1 data_d1['national_inv_bins'] = 0
2 data_d1.loc[(data_d1['national_inv'] == 0.0), 'national_inv_bins'] = 1
3 data_d1.loc[(data_d1['national_inv'] == 1.0), 'national_inv_bins'] = 2
4 data_d1.loc[(data_d1['national_inv'] >= 2.0) & (data_d1['national_inv'] <= 9.0), 'nati
5 data_d1.loc[(data_d1['national_inv'] >= 10.0), 'national_inv_bins'] = 4
6
7 data_d2['national_inv_bins'] = 0
8 data_d2.loc[(data_d2['national_inv'] == 0.0), 'national_inv_bins'] = 1
9 data_d2.loc[(data_d2['national_inv'] == 1.0), 'national_inv_bins'] = 2
10 data_d2.loc[(data_d2['national_inv'] >= 2.0) & (data_d2['national_inv'] <= 9.0), 'nati
11 data_d2.loc[(data_d2['national_inv'] >= 10.0), 'national_inv_bins'] = 4
12
13 X_test_en['national_inv_bins'] = 0
14 X_test_en.loc[(X_test_en['national_inv'] == 0.0), 'national_inv_bins'] = 1
15 X_test_en.loc[(X_test_en['national_inv'] == 1.0), 'national_inv_bins'] = 2
16 X_test_en.loc[(X_test_en['national_inv'] >= 2.0) & (X_test_en['national_inv'] <= 9.0)
17 X_test_en.loc[(X_test_en['national_inv'] >= 10.0), 'national_inv_bins'] = 4

1 data_d1 = data_d1.drop(['went_on_backorder'], axis=1)

```

▼ Adding SVD and AutoEncoder Features in the main data

```

1 for i in range(2):
2     data_d1['T_SVD_'+str(i)] = data_d1_SVD[:,i]
3     data_d2['T_SVD_'+str(i)] = data_d2_SVD[:,i]
4     X_test_en['T_SVD_'+str(i)] = X_test_en_SVD[:,i]
5     print(data_d1.shape, data_d2.shape, X_test_en.shape)

(752548, 23) (752548, 23) (376275, 23)

1 for i in range(2):
2     data_d1['AutoEncoder_'+str(i)] = encoded_output_data_d1[:,i]
3     data_d2['AutoEncoder_'+str(i)] = encoded_output_data_d2[:,i]
4     X_test_en['AutoEncoder_'+str(i)] = encoded_output_test[:,i]
5     print(data_d1.shape, data_d2.shape, X_test_en.shape)

(752548, 25) (752548, 25) (376275, 25)

```

▼ Normalise Data

```

1 # Normalise Median Imputed Data
2 MinMaxSc = MinMaxScaler()
3
4 columns = data_d1.columns
5 data_d1 = pd.DataFrame(MinMaxSc.fit_transform(data_d1), columns=columns)
6 data_d2 = pd.DataFrame(MinMaxSc.transform(data_d2), columns=columns)

```

```
6 data_d2 = pd.DataFrame(MinMaxSc.transform(data_d2),columns=columns)
7 X_test_en = pd.DataFrame(MinMaxSc.transform(X_test_en),columns=columns)
8
9 print(data_d1.shape,data_d2.shape,X_test_en.shape)

(752548, 25) (752548, 25) (376275, 25)

1 number_of_models=15
2 modelList,data = baseModel_DecisionTree(number_of_models,data_d1,y_d1,data_d2,y_d2)
```

Generating samples with 60% sample and 40% sample with replacement

In Sample Generate

Sample generation for Model : 0
Sample generation for Model : 1
Sample generation for Model : 2
Sample generation for Model : 3
Sample generation for Model : 4
Sample generation for Model : 5
Sample generation for Model : 6
Sample generation for Model : 7
Sample generation for Model : 8
Sample generation for Model : 9
Sample generation for Model : 10
Sample generation for Model : 11
Sample generation for Model : 12
Sample generation for Model : 13
Sample generation for Model : 14

Out Sample Generate

Fit Decision Tree Models on samples generated

In Fitting Models

Fitting for Model : 0
Fitting for Model : 1
Fitting for Model : 2
Fitting for Model : 3
Fitting for Model : 4
Fitting for Model : 5
Fitting for Model : 6
Fitting for Model : 7
Fitting for Model : 8
Fitting for Model : 9
Fitting for Model : 10
Fitting for Model : 11
Fitting for Model : 12
Fitting for Model : 13
Fitting for Model : 14

Out Fitting Models

Predict values for second dataset

In Predict

Predict for Model : 0
Predict for Model : 1
Predict for Model : 2
Predict for Model : 3
Predict for Model : 4
Predict for Model : 5
Predict for Model : 6
Predict for Model : 7
Predict for Model : 8
Predict for Model : 9
Predict for Model : 10
Predict for Model : 11


```

Predict for Model : 12
Predict for Model : 13
Predict for Model : 14
Out Predict
Zip each row of the column to form dataset
In form Meta Data
Out form Meta Data

```

```
1 # Meta Classifier - Logistic Regression
```

```

1 lg_clf = LogisticRegression()
2 params = {'C':[10**x for x in range(-5,4)]}
3 grid_log_clf = GridSearchCV(lg_clf,param_grid=params,n_jobs=-1,scoring='f1_macro',ver
4 grid_log_clf.fit(data,y_d2.values.ravel())

```

```

Fitting 5 folds for each of 9 candidates, totalling 45 fits
[Parallel(n_jobs=-1)]: Using backend LokyBackend with 2 concurrent workers.
[Parallel(n_jobs=-1)]: Done 28 tasks      | elapsed: 24.3s
[Parallel(n_jobs=-1)]: Done 45 out of 45 | elapsed: 39.4s finished
GridSearchCV(cv=None, error_score=nan,
              estimator=LogisticRegression(C=1.0, class_weight=None, dual=False,
              fit_intercept=True,
              intercept_scaling=1, l1_ratio=None,
              max_iter=100, multi_class='auto',
              n_jobs=None, penalty='l2',
              random_state=None, solver='lbfgs',
              tol=0.0001, verbose=0,
              warm_start=False),
              iid='deprecated', n_jobs=-1,
              param_grid={'C': [1e-05, 0.0001, 0.001, 0.01, 0.1, 1, 10, 100,
              1000]}},
              pre_dispatch='2*n_jobs', refit=True, return_train_score=False,
              scoring='f1_macro', verbose=3)

```

```

1 labels = grid_log_clf.cv_results_['params']
2 x_axis = range(1,10)
3 y_axis = grid_log_clf.cv_results_['mean_test_score']
4 for i,label in enumerate(labels):
5     print(label,":",y_axis[i])
6
7 plt.xlabel("Iteration")
8 plt.ylabel("F1 score")
9 plt.title("GRID search cv - F1 score")
10 plt.plot(x_axis,y_axis)

```

```
{'C': 1e-05} : 0.4982939749796286
{'C': 0.0001} : 0.4982939749796286
{'C': 0.001} : 0.5534706671930081
{'C': 0.01} : 0.621628415859751
{'C': 0.1} : 0.6283068586396076
{'C': 1} : 0.6287836210301228
{'C': 10} : 0.6290223855418355
{'C': 100} : 0.6292815696831074
{'C': 1000} : 0.6292815696831074
[<matplotlib.lines.Line2D at 0x7f3a3017db50>]
```



```
1 grid_log_clf.best_estimator_
```

```
LogisticRegression(C=100, class_weight=None, dual=False, fit_intercept=True,
                    intercept_scaling=1, l1_ratio=None, max_iter=100,
                    multi_class='auto', n_jobs=None, penalty='l2',
                    random_state=None, solver='lbfgs', tol=0.0001, verbose=0,
                    warm_start=False)
```



```
1 clf_best = LogisticRegression(C=10)
2 clf_best.fit(data,y_d2.values.ravel())
```

```
LogisticRegression(C=10, class_weight=None, dual=False, fit_intercept=True,
                    intercept_scaling=1, l1_ratio=None, max_iter=100,
                    multi_class='auto', n_jobs=None, penalty='l2',
                    random_state=None, solver='lbfgs', tol=0.0001, verbose=0,
                    warm_start=False)
```

```
1 test_data = perform_test(number_of_models,modelList,X_test_en)
2 y_test_pred = clf_best.predict(test_data)
```

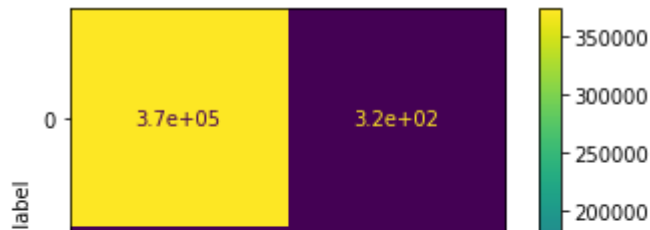
In Predicting Output for Test Data
Out Predicting Output for Test Data

```
1 print("Test F1 score : ",f1_score(y_test_en,y_test_pred,average='macro'))
```

Test F1 score : 0.6207567400765969

```
1 confusionMatrix = confusion_matrix(y_test_en,y_test_pred,labels=[0,1])
2 cm_display = ConfusionMatrixDisplay(confusion_matrix = confusionMatrix,display_labels
3 cm_display.plot())
```

```
<sklearn.metrics._plot.confusion_matrix.ConfusionMatrixDisplay at 0x7f3a30049a10>
```



- With use of Custom Ensemble models we see that with increasing number of base models ie Decision Tree the F1 scores initially increase and then tapers off.
- With Custom Ensembles applied on with and without feature engineering data , the f1 score almost in the same range.
- Amongst all the models applied Random Forest with Oversampling data has given the best F1 score of 0.7 on cv and test data.

1

▼ Summary

- We observe that with every model the improvement in f1 scores.
- We saw that the Model on oversampled data perform well as compared to Undersampled data, this could be because of loss of data while undersampling.
- Among all the models , Random Forest Model provided a f1 score of 0.69 and models including XGBoost and Custom ensembles gave a f1 score of 0.62.

```
1 from prettytable import PrettyTable
2 table = PrettyTable()
3 table.field_names = ['Model', 'Parameter', 'Technique', 'Train F1 Score', 'CV F1 Score', '
4 table.add_row(['Logistic Regression', 'Best Alpha = ' + str(10), 'SMOTE', 0.79, 0.50, 0.50]
5 table.add_row(['Logistic Regression', 'Best Alpha = ' + str(1e-05), 'RUS', 0.52, 0.49, 0.49]
6 table.add_row(['Random Forest', "n_estimators = " + str(500) + " , max_depth = " + str(50)
7 table.add_row(['Random Forest', "n_estimators = " + str(100) + " , max_depth = " + str(20)
8 table.add_row(['XGBOOST', "n_estimators = " + str(500) + " , learning rate = " + str(0.1)
9 table.add_row(['XGBOOST', "n_estimators = " + str(500) + " , learning rate = " + str(0.2)
10 table.add_row(['AdaBoost', "n_estimators = " + str(500) + " , learning rate = " + str(1), '
11 table.add_row(['AdaBoost', "n_estimators = " + str(500) + " , learning rate = " + str(0.1)
12 table.add_row(['Ensemble Models', "Decision Tree + Logistic Regression ", 'With FE', 0.6
13 table.add_row(['Ensemble Models', "Decision Tree + Logistic Regression ", 'Without FE',
14 print(table)
```

Model	Parameter
Logistic Regression	Best Alpha = 10
Logistic Regression	Best Alpha = 1e-05
Random Forest	n_estimators = 500 , max_depth = 50 , min_samples_split = 5
Random Forest	n_estimators = 100 , max_depth = 20 , min_samples_split = 2
XGBOOST	n_estimators = 500 , learning rate = 0.1 , max_depth = 3

```
|      XGBOOST      | n_estimators = 500 , learning rate = 0.2 , max_depth = 10
|      AdaBoost     | n_estimators = 500 , learning rate = 1
|      AdaBoost     | n_estimators = 500 , learning rate = 0.1
| Ensemble Models   | Decision Tree + Logistic Regression
| Ensemble Models   | Decision Tree + Logistic Regression
+-----+-----+
|<----->|
```

1

1

1

1

1

1

1

1

1