

TechNova Sales Insights

Problem Statement

Conduct analysis to help determine how effective campaigns that have been run historically are. The following information should be collated:

Campaign Details:

- Each record should correspond to a particular Email or SMS campaign run
- Details of the campaign
- Date time when the campaign was run

Customer Details:

- Customer ID (who was sent Email or SMS)
- Age, Gender, Marital Status, Residential Pin code
- Number of Transactions (qty and value) in the Last 3/6/12 Months
- % of Premium/Mean/earValue Transactions

Campaign Outcome:

- Whether there was any transaction within 1 month of campaign - accordingly create an outcome variable = 0 for no transactions and 1 for transactions (the transaction should have connection to the campaign)
- Exclude records which were not sent any campaigns
- Transactions which were not related to the campaign should not be considered either 0 or 1 but should be excluded

```
In [1]: # Importing Essential Libraries to work on dataset
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
import warnings
warnings.filterwarnings("ignore")

In [2]: # Loading CSV file into DataFrame
df_cd = pd.read_csv('C:/Users/Dheem/OneDrive/Documents/Final Project/TechNova Sales Insights/Campaign_Data.csv')

In [3]: df_cd

Out[3]:
```

	CustID	status	Campaign_Exec_Date	Campaign_Channel	
	0	111111111RREFG1V6Gg1FejGEVVB88JM	viewed	2019-05-28	Email
	1	111111111RREFG1V6Gg1FejGEVVB88JM	delivered	2019-05-29	Email
	2	111111111RREFG1V6Gg1FejGEVVB88JM	delivered	2019-05-30	Email
	3	111111111RREFG1V6Gg1FejGEVVB88JM	delivered	2019-06-03	Email
	4	111111111RREFG1V6Gg1FejGEVVB88JM	delivered	2019-06-04	Email

	109125897	ttttttHGe	viewed	2020-11-24	Email
	10912588	ttttttHGe	viewed	2020-11-25	Email
	10912589	ttttttHGe	viewed	2020-11-26	Email
	10912590	ttttttHGe	delivered	2020-11-27	Email
	10912591	ttttttHGe	delivered	2020-11-28	Email

10912592 rows x 4 columns

```
In [4]: df_cd.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 10912592 entries, 0 to 10912591
Data columns (total 4 columns):
# Column      Dtype
---  --
0 CustID      object
1 status      object
2 Campaign_Exec_Date object
3 Campaign_Channel object
dtypes: object(4)
memory usage: 333.0+ MB

In [5]: # Loading CSV file into DataFrame
df_fd = pd.read_csv('C:/Users/Dheem/OneDrive/Documents/Final Project/TechNova Sales Insights/filtered_data.csv')

In [6]: df_fd

Out[6]:
```

	CustID	ItemDesc	MerchCategoryDescription	MerchClassDescription	MerchGroupDescription	SaleValue	OrderedQuantity	OrderDate	SalesChannelCode	Ecom_BnM_Indicator	St
	0	JFRGRFjeiVGBRfBRlFqBJu78d8h	APPLE PROTECTION PLAN MACBOOK PRO	APPLE PROTECTION PLAN	Mobile Computing	Computers Peripherals	25398.0001	1.0001	2019-01-01	RS	B&M
	1	JFRGRFjeiVGBRfBRlFqBJu78d8h	MacBook Pro 16 2.6GHz I7 9G 16GB S12GB SG	MacBook Pro 16 2.6GHz I7 9G 16GB S12GB SG	Mobile Computing	Computers Peripherals	193703.1001	1.0001	2019-01-01	RS	B&M
	2	VIAG7AFJ1E1HVU7V7H7APR1J1A8et	Samsung Galaxy A70 5 6GB White	Smart Phones (OS Based)	Communication	26518.9801	1.0001	2019-01-01	RS	B&M	
	3	IAEYH8EG1Ede1GE1JqEjVVMGFGfo	DELL New bag pack	DELL New bag pack	Mobile Computing	Computers Peripherals	1120.9801	1.0001	2019-01-01	RS	B&M
	4	RV1BRMMHMHVjoe78GdGAV1RAM7AoeE	HP Pav 15-0k046TX 9G5 8GBX4 256GB +CRF	HP Pav 15-0k046TX 9G5 8GBX4 256GB +CRF	Gaming Laptops	Computers Peripherals	66085.8001	1.0001	2019-01-01	RS	B&M

	453103	1ERAFMMGG7eABEEBJuAGFHVMEVBGR	Samsung Galaxy A53 (8GB+128GB) Black	Samsung Galaxy A53 (8GB+128GB) Black	Smart Phones (OS Based)	Communication	32460.5311	1.0001	2021-12-31	RS	B&M
	453104	1BRVAFJ1B1BVET7ABJAEQFHRH1J7H	APPLE PROTECT+ WITH ACS IPHONE 14 PRO MAX	APPLE PROTECT+ WITH ACS	Smart Phones (OS Based)	Communication	16472.0311	1.0001	2021-12-31	RS	B&M
	453105	1BRVAFJ1B1BVET7ABJAEQFHRH1J7H	Apple iPhone 14 Pro Max 256GB DeepPurple	Apple iPhone 14 Pro Max 256GB DeepPurple	Smart Phones (OS Based)	Communication	148399.8001	1.0001	2021-12-31	RS	B&M
	453106	VGBFJo7H7H718AF7HGV1FE1HJMUG	SAMSUNG LED 180CM 75AU7700 UHD4K	SAMSUNG LED 180CM 75AU7700 UHD4K	TV LCD	Entertainment	111424.8001	1.0001	2021-12-31	RS	B&M
	453107	oEVfHG1FRHHGEB077RHMR@F1e1Ji7	Apple iPhone 14 Plus (256GB, Blue)	Apple iPhone 14 Plus (256GB, Blue)	Smart Phones (OS Based)	Communication	94849.8001	1.0001	2021-12-31	TN	Ecom

453108 rows x 20 columns

```
In [7]: df_fd.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 453108 entries, 0 to 453107
Data columns (total 28 columns):
# Column      Dtype
---  --
0 CustID      object
1 status      object
2 MerchCategoryDescription  object
3 MerchClassDescription    object
4 MerchGroupDescription    object
5 SaleValue               float64
6 OrderedQuantity         float64
7 OrderDate              object
8 SalesChannelCode        object
9 Ecom_BnM_Indicator       object
10 StoreID               float64
11 StoreCode             object
12 StoreState            object
13 StorePincode           object
14 Log_OrderedQuantity    float64
15 Sort_SaleValue         float64
16 Month                 int64
17 Year                  int64
18 Segment               object
dtypes: float64(5), int64(2), object(13)
memory usage: 69.1+ MB

In [8]: # Making a List of columns to be dropped
columns_to_drop = ['ItemDesc', 'MerchCategoryDescription', 'MerchClassDescription', 'MerchGroupDescription', 'SalesChannelCode', 'Ecom_BnM_Indicator', 'StoreID', 'StoreCode', 'StoreCity', 'StoreState', 'StorePincode', 'Log_OrderedQuantity', 'Sort_SaleValue', 'Month', 'Year']

In [9]: df_fd.drop(columns=columns_to_drop, inplace=True)

In [10]: df_fd

Out[10]:
```

	CustID	SaleValue	OrderedQuantity	OrderDate	Segment	
	0	JFRGRFjeiVGBRfBRlFqBJu78d8h	25398.0001	1.0001	2019-01-01	Mainstream
	1	JFRGRFjeiVGBRfBRlFqBJu78d8h	193703.1001	1.0001	2019-01-01	Premium
	2	VIAG7AFJ1E1HVU7V7H7APR1J1A8et	26518.9801	1.0001	2019-01-01	Mainstream
	3	IAEYH8EG1Ede1GE1JqEjVVMGFGfo	1120.9801	1.0001	2019-01-01	Value
	4	RV1BRMMHMHVjoe78GdGAV1RAM7AoeE	66085.8001	1.0001	2019-01-01	Mainstream

	453103	1ERAFMMGG7eABEEBJuAGFHVMEVBGR	32460.5311	1.0001	2021-12-31	Premium
	453104	1BRVAFJ1B1BVET7ABJAEQFHRH1J7H	16472.0311	1.0001	2021-12-31	Mainstream
	453105	1BRVAFJ1B1BVET7ABJAEQFHRH1J7H	148399.8001	1.0001	2021-12-31	Premium
	453106	VGBFJo7H7H718AF7HGV1FE1HJMUG	111424.8001	1.0001	2021-12-31	Premium
	453107	oEVfHG1FRHHGEB077RHMR@F1e1Ji7	94849.8001	1.0001	2021-12-31	Premium

453108 rows x 5 columns

```
In [11]: # Loading CSV file into DataFrame
df_md = pd.read_csv('C:/Users/Dheem/OneDrive/Documents/Final Project/TechNova Sales Insights/Cleaned_Customer_Master_Data.csv')

In [12]: df_md

Out[12]:
```

	CustID	Gender	Pincode	State	
	0	76AATAh	DEFAULT	678506.0	Keral
	1	76BUJH7	DEFAULT	770001.0	Odisha
	2	787L7JH	DEFAULT	245101.0	Uttar Pradesh
	3	787J7HHH	DEFAULT	673008.0	Keral
	4	787BH8Ho	DEFAULT	501510.0	Telangana

	492176	EEEMV1EAEH7M7eMAGAGG7H1MGHG1JH	Male	141010.0	Punjab
	492177	EEEMeAMFMHMFj0R8HAAMRV1oeRe1F	Female	560059.0	Karnataka
	492178	EEEEeKHEMHQj0eGheGR1Fe771ERG	Male	560085.0	Karnataka
	492179	EEEE1e88H7HAFGBM81MG8MB@VU	Male	411007.0	Maharashtra
	492180	EEEEERH1RJ1E1G7F7MGBFEFP1MGOM	Male	400614.0	Maharashtra

492181 rows x 4 columns

```
In [13]: df_md.drop(columns=['State', 'Pincode'], inplace=True)

In [14]: # Merging Transaction Filtered Data with Customer Master Data on CustID
merged_2 = pd.merge(df_md, df_md, on='CustID', how='left')

In [15]: # Merging campaign data with above merged data on CustID
merged_data = pd.merge(df_md, merged_2, on='CustID', how='left')

In [16]: merged_data.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 1448008 entries, 0 to 1448007
Data columns (total 9 columns):
# Column      Dtype
---  --
0 CustID      object
1 status      object
2 Campaign_Exec_Date object
3 Campaign_Channel object
4 SaleValue   float64
5 OrderedQuantity float64
6 OrderDate   object
7 Segment     object
8 Gender      object
dtypes: float64(2), object(7)
memory usage: 94.1+ MB

In [17]: # Converting dates to datetime format
merged_data['Campaign_Exec_Date'] = pd.to_datetime(merged_data['Campaign_Exec_Date'])
merged_data['OrderDate'] = pd.to_datetime(merged_data['OrderDate'])

In [18]: # Converting other columns with dtype as object into category type
merged_data['CustID'] = merged_data['CustID'].astype('category')
merged_data['status'] = merged_data['status'].astype('category')
merged_data['Campaign_Channel'] = merged_data['Campaign_Channel'].astype('category')
merged_data['Segment'] = merged_data['Segment'].astype('category')
merged_data['Gender'] = merged_data['Gender'].astype('category')

In [19]: # Defining a function to calculate the number and value of transactions in the last N months
def filtered_transaction_metrics(data, months):
    end_date = data['Campaign_Exec_Date']
    start_date = end_date - pd.DateOffset(months=months)

    modified_data = data[(data['OrderDate'] >= start_date) & (data['OrderDate'] <= end_date)]
    transaction_count = modified_data.groupby('CustID')['OrderedQuantity'].sum().reset_index(name='Transactions_Last_months_Months')
    transaction_value = modified_data.groupby('CustID')['SaleValue'].sum().reset_index(name='Transaction_Value_Last_months_Months')

    return transaction_count, transaction_value

In [20]: # Calculating transaction metrics for the last 3, 6 and 12 months
metrics_3m = filtered_transaction_metrics(merged_data, 3)
metrics_6m = filtered_transaction_metrics(merged_data, 6)
metrics_12m = filtered_transaction_metrics(merged_data, 12)

In [21]: # Merging these metrics back into the merged data DataFrame
merged_data = merged_data.merge(metrics_3m[0], on='CustID', how='left')
merged_data = merged_data.merge(metrics_3m[1], on='CustID', how='left')
merged_data = merged_data.merge(metrics_6m[0], on='CustID', how='left')
merged_data = merged_data.merge(metrics_6m[1], on='CustID', how='left')
merged_data = merged_data.merge(metrics_12m[0], on='CustID', how='left')
merged_data = merged_data.merge(metrics_12m[1], on='CustID', how='left')

In [22]: merged_data

Out[22]:
```

	CustID	status	Campaign_Exec_Date	Campaign_Channel	SaleValue	OrderedQuantity	OrderDate	Segment	Gender	Transactions_Last_3_Months	Transaction_Value_Last_3_Months	Tn	
	0	111111111RREFG1V6Gg1FejGEVVB88JM	viewed	2019-05-28	Email	NaN	NaN	NaN	NaN	NaN	NaN	0.0	0.0
	1	111111111RREFG1V6Gg1FejGEVVB88JM	delivered	2019-05-29	Email	NaN	NaN	NaN	NaN	NaN	NaN	0.0	0.0
	2	111111111RREFG1V6Gg1FejGEVVB88JM	delivered	2019-05-30	Email	NaN	NaN	NaN	NaN	NaN	NaN	0.0	0.0
	3	111111111RREFG1V6Gg1FejGEVVB88JM	delivered	2019-06-03	Email	NaN	NaN	NaN	NaN	NaN	NaN	0.0	0.0
	4	111111111RREFG1V6Gg1FejGEVVB88JM	delivered	2019-06-04	Email	NaN	NaN	NaN	NaN	NaN	NaN	0.0	0.0

	14480596	ttttttHGe	viewed	2020-11-24	Email	NaN	NaN	NaN	NaN	NaN	NaN	0.0	0.0
	14480597	ttttttHGe	viewed	2020-11-25	Email	NaN	NaN	NaN	NaN	NaN	NaN	0.0	0.0
	14480598	ttttttHGe	viewed	2020-11-26	Email	NaN	NaN	NaN	NaN	NaN	NaN	0.0	0.0
	14480599	ttttttHGe	delivered	2020-11-27	Email	NaN	NaN	NaN	NaN	NaN	NaN	0.0	0.0
	14480600	ttttttHGe	delivered	2020-11-28	Email	NaN	NaN	NaN	NaN	NaN	NaN	0.0	0.0

14480601 rows x 15 columns

```
In [23]: merged_data.isnull().sum()

In [23]: status      0
         CustID     0
         Campaign_Exec_Date 0
         Campaign_Channel   0
         SaleValue         7459108
         OrderedQuantity    7459108
         OrderDate         7459108
         Segment          7459108
         Gender           7587808
         Transactions_Last_3_Months 0
         Transaction_Value_Last_3_Months 0
         Transaction_Last_6_Months 0
         Transaction_Value_Last_6_Months 0
         Transaction_Last_12_Months 0
         Transaction_Value_Last_12_Months 0
         dtype: int64

In [24]: # Making a List of columns with null values and eliminating them
fd_columns = ['SaleValue', 'OrderedQuantity', 'OrderDate', 'Segment']
merged_data = merged_data.dropna(subset=fd_columns)

In [25]: merged_data = merged_data.dropna(subset=['Gender'])

In [26]: # Defining a function to calculate segment percentage
def to_calculate_segment_percentage(data, segment):
    segment_data = data[data['Segment'] == segment]
    segment_percentage = segment_data.groupby('CustID')['OrderedQuantity'].sum() / data.groupby('CustID')['OrderedQuantity'].sum()
    segment_percentage = segment_percentage.reset_index(name='segment_Percentage')
    return segment_percentage

In [27]: # Calculating percentage for each segment
premium_percentage = to_calculate_segment_percentage(merged_data, 'Premium')
mainstream_percentage = to_calculate_segment_percentage(merged_data, 'Mainstream')
value_percentage = to_calculate_segment_percentage(merged_data, 'Value')

In [28]: # Merging these percentages back into the merged data DataFrame
merged_data = merged_data.merge(premium_percentage, on='CustID', how='left')
merged_data = merged_data.merge(mainstream_percentage, on='CustID', how='left')
merged_data = merged_data.merge(value_percentage, on='CustID', how='left')

In [29]: merged_data

Out[29]:
```

	CustID	status	Campaign_Exec_Date	Campaign_Channel	SaleValue	OrderedQuantity	OrderDate	Segment	Gender	Transactions_Last_3_Months	Transaction_Value_Last_3_Months	Transactions_Last_6_Months	Transaction_Value_Last_6_Months	
	0	111111111RREFG1V6Gg1FejGEVVB88JM	viewed	2019-05-28	Email	15298.9801	1.0001	2021-04-13	Mainstream	Male	0.0	0.0	0.0	0.0
	1	111111111RREFG1V6Gg1FejGEVVB88JM	viewed	2019-05-28	Email	45898.9801	1.0001	2021-04-13	Premium	Male	0.0	0.0	0.0	0.0
	2	111111111RREFG1V6Gg1FejGEVVB88JM	viewed	2019-05-28	Email	15298.9801	1.0001	2021-04-13	Mainstream	Male	0.0	0.0	0.0	0.0
	3	111111111RREFG1V6Gg1FejGEVVB88JM	delivered	2019-05-28	Email	15298.9801	1.0001	2021-04-13	Mainstream	Male	0.0	0.0	0.0	0.0
	4	111111111RREFG1V6Gg1FejGEVVB88JM	delivered	2019-05-29	Email	45898.9801	1.0001	2021-09-28	Premium	Male	0.0	0.0	0.0	0.0
	
	6973556	ttttt7o1H	viewed	2020-11-24	Email	9169.8001	1.0001	2021-04-26	Mainstream	DEFAULT	0.0	0.0	0.0	0.0
	6973557	ttttt7o1H	viewed	2020-11-25	Email	9169.8001	1.0001	2021-04-26	Mainstream	DEFAULT	0.0	0.0	0.0	0.0
	6973558	ttttt7o1H	viewed	2020-11-26	Email	9169.8001	1.0001	2021-04-26	Mainstream	DEFAULT	0.0	0.0	0.0	0.0
	6973559	ttttt7o1H	delivered	2020-11-27	Email	9169.8001	1.0001	2021-04-26	Mainstream	DEFAULT	0.0	0.0	0.0	0.0
	6973560	ttttt7o1H	delivered	2020-11-28	SMS	9169.8001	1.0001	2021-04-26	Mainstream	DEFAULT	0.0	0.0	0.0	0.0

6973561 rows x 18 columns

```
In [30]: # Creating the end date for the 1-month window over each campaign
merged_data['Campaign_End_Date'] = merged_data['Campaign_Exec_Date'] + pd.DateOffset(months=1)

In [31]: # Creating a flag for transactions within 1 month of the campaign
merged_data['Transaction_within_1_Month'] = merged_data.apply(
    lambda row: 1 if (row['OrderDate'] >= row['Campaign_Exec_Date']) and (row['OrderDate'] <= row['Campaign_End_Date'] + pd.DateOffset(months=1)) else 0,
    axis=1
)

In [32]: # Ensuring that there are both entries 1 and 0 in our required outcome column
merged_data['Transaction_within_1_Month'].unique()

Out[32]: array([0, 1], dtype=int64)

In [33]: # Dropping the temporary end date column as it's no longer needed
merged_data.drop(columns=['Campaign_End_Date'], inplace=True)

In [34]: # Filtering out rows where there was no campaign
merged_data = merged_data[merged_data['Campaign_Channel'].notnull()]

In [35]: # Dropping CustID column since to train the ML models only numeric data is allowed
merged_data = merged_data.drop(columns=['CustID'])

In [36]: merged_data

Out[36]:
```

	status	Campaign_Exec_Date	Campaign_Channel	SaleValue	OrderedQuantity	OrderDate	Segment	Gender	Transactions_Last_3_Months	Transaction_Value_Last_3_Months	Transactions_Last_6_Months	Transaction_Value_Last_6_Months	
	0	2019-05-28	Email	15298.9801	1.0001	2021-04-13	Mainstream	Male	0.0	0.0	0.0	0.0	
	1	2019-05-28	Email	45898.9801	1.0001	2021-09-28	Premium	Male	0.0	0.0	0.0	0.0	
	2	viewed	2019-05-28	Email	15298.9801	1.0001	2021-04-13	Mainstream	Male	0.0	0.0	0.0	0.0
	3	delivered	2019-05-29	Email	15298.9801	1.0001	2021-04-13	Mainstream	Male	0.0	0.0	0.0	0.0
	4	delivered	2019-05-29	Email	45898.9801	1.0001	2021-09-28	Premium	Male	0.0	0.0	0.0	0.0
	
	6973556	viewed	2020-11-24	Email	9169.8001	1.0001	2021-04-26	Mainstream	DEFAULT	0.0	0.0	0.0	0.0
	6973557	viewed	2020-11-25	Email	9169.8001	1.0001	2021-04-26	Mainstream	DEFAULT	0.0	0.0	0.0	0.0
	6973558	viewed	2020-11-26	Email	9169.8001	1.0001	2021-04-26	Mainstream	DEFAULT	0.0	0.0	0.0	0.0
	6973559	delivered	2020-11-27	Email	9169.8001	1.0001	2021-04-26	Mainstream	DEFAULT	0.0	0.0	0.0	0.0
	6973560	delivered	2020-11-28	SMS	9169.8001	1.0001	2021-04-26	Mainstream	DEFAULT	0.0	0.0	0.0	0.0

6973561 rows x 18 columns

```
In [37]: merged_data.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 6973561 entries, 0 to 6973560
Data columns (total 18 columns):
# Column      Dtype
---  --
0 status      object
1 Campaign_Exec_Date object
2 Campaign_Channel object
3 SaleValue   float64
4 OrderedQuantity float64
5 OrderDate   object
6 Segment     object
7 Gender      object
8 Transactions_Last_3_Months float64
9 Transaction_Value_Last_3_Months float64
10 Transactions_Last_6_Months float64
11 Transaction_Value_Last_6_Months float64
12 Transactions_Last_12_Months float64
13 Transaction_Value_Last_12_Months float64
14 Premium_Percentage float64
15 Mainstream_Percentage float64
16 Value_Percentage float64
17 Transaction_within_1_Month int64
dtypes: float64(14), datetime64[ns](2), int64(1), int64(1)
memory usage: 771.5+ MB

In [38]:
```

Problem Statement

- Using the above data, develop a campaign effectiveness model using Random Forest and XGBoost - use a 80:20 train/test split
- Create a confusion matrix for both the models and report the Precision, Recall and Accuracy
- Report Model Results appropriately using charts generated in Python - leveraging matplotlib/seaborn packages

```
In [39]: # Importing Essential Libraries to work on dataset
from sklearn.model_selection import train_test_split
from sklearn.ensemble import RandomForestClassifier
from sklearn.preprocessing import OneHotEncoder
import numpy as np
from sklearn.preprocessing import LabelEncoder

In [40]: le = LabelEncoder()

In [41]: merged_data['status'] = le.fit_transform(merged_data['status'])


```


