

# Bank Customer Behaviour Prediction

We have been provided the data from a credit card processor company which has approximately 20000 records of the same number of users which are represented by their ID to keep their privacy. The data is of customer credit/debit card spend of three months(April, May & June) and we have to predict their expected average spend in the coming 3 months (July, August & September). We have been provided a separate dataset that has the average of three months (July, August & September) for some entries and not all which we have to predict based on the previous month data and other factors in the dataset. The entries we have been provided can be our training dataset and we have to match our predicted average with the actual values. The model with the lowest Root Mean Square Percentage Error (RMSPE) will be selected.

## Step 1: Setting up the work environment

We will download the necessary packages and obtain the datasets. We will have a quick overview of the datatypes.

```
In [64]: #Getting packages
import pandas as pd # For data manipulation and analysis
import numpy as np # For Numerical operations
import matplotlib.pyplot as plt # For basic data visualisation
import seaborn as sns # For advanced data visualization
from scipy import stats # For statistics
from sklearn.linear_model import LinearRegression #for linear regression
from sklearn.metrics import mean_squared_error, make_scorer #for Error measurement
from sklearn.model_selection import train_test_split #for model training and testing
from sklearn.metrics import r2_score #for model efficiency
from sklearn.linear_model import Ridge #for ridge regularization
from sklearn.linear_model import Lasso #for lasso regularization
from sklearn.tree import DecisionTreeRegressor #for decision tree regression
from sklearn.ensemble import RandomForestRegressor #for random forest regression
from sklearn.model_selection import GridSearchCV, RandomizedSearchCV #for Hyperparameter tuning
```

```
In [2]: #getting the dataset
bod=pd.read_excel('C:\Users\ujaydutta\Desktop\Data analysis\Projects\Bank customer analysis\CreditCard copy head1')
bod.head()
```

```
Out[2]:
```

	ID	cc_cons_apr	dc_cons_apr	cc_cons_may	dc_cons_may	cc_cons_jun	dc_cons_jun	cc_count_apr	cc_count_may	cc_count_jun	...	credit
0	17051	3412.60	2909.34	1077.36	3084.50	4615.97	6693.0	19.0	10.0	10.0	2.0	...
1	11491	18133.76	3877.00	6123.78	5597.43	64620.00	6117.00	4.0	60	60	1	...
2	7433	6363.65	735.49	8799.00	13768.00	38266.00	2638.0	2.0	83	47	...	...
3	14606	12765.66	4429.16	16745.86	6360.00	29063.97	3711.0	12.0	2	16	...	...
4	8381	27819.70	1944.00	7006.50	2228.50	1096.25	1065.0	15.0	13	61	...	...

5 rows x 39 columns

```
In [3]: #Examining the dataset
bod.shape
```

```
Out[3]:
```

```
(20000, 39)
```

```
In [4]: #seeing datatypes
bod.dtypes
```

```
Out[4]:
```

	ID	cc_cons_apr	dc_cons_apr	cc_cons_may	dc_cons_may	cc_cons_jun	dc_cons_jun	cc_count_apr	cc_count_may	cc_count_jun	...	credit
cc_cons_apr	float64	float64	float64	float64	float64	float64	float64	float64	float64	float64	float64	float64
dc_cons_apr	float64	float64	float64	float64	float64	float64	float64	float64	float64	float64	float64	float64
cc_cons_may	float64	float64	float64	float64	float64	float64	float64	float64	float64	float64	float64	float64
dc_cons_may	float64	float64	float64	float64	float64	float64	float64	float64	float64	float64	float64	float64
cc_cons_jun	float64	float64	float64	float64	float64	float64	float64	float64	float64	float64	float64	float64
dc_cons_jun	float64	float64	float64	float64	float64	float64	float64	float64	float64	float64	float64	float64
cc_count_apr	float64	float64	float64	float64	float64	float64	float64	float64	float64	float64	float64	float64
cc_count_may	float64	float64	float64	float64	float64	float64	float64	float64	float64	float64	float64	float64
cc_count_jun	float64	float64	float64	float64	float64	float64	float64	float64	float64	float64	float64	float64
dc_count_apr	float64	float64	float64	float64	float64	float64	float64	float64	float64	float64	float64	float64
dc_count_may	float64	float64	float64	float64	float64	float64	float64	float64	float64	float64	float64	float64
dc_count_jun	float64	float64	float64	float64	float64	float64	float64	float64	float64	float64	float64	float64
card_lim	int64	int64	int64	int64	int64	int64	int64	int64	int64	int64	int64	int64
personal_loan_active	int64	int64	int64	int64	int64	int64	int64	int64	int64	int64	int64	int64
vehicle_loan_active	int64	int64	int64	int64	int64	int64	int64	int64	int64	int64	int64	int64
personal_loan_closed	float64	float64	float64	float64	float64	float64	float64	float64	float64	float64	float64	float64
vehicle_loan_closed	float64	float64	float64	float64	float64	float64	float64	float64	float64	float64	float64	float64
investment_1	float64	float64	float64	float64	float64	float64	float64	float64	float64	float64	float64	float64
investment_2	float64	float64	float64	float64	float64	float64	float64	float64	float64	float64	float64	float64
investment_3	float64	float64	float64	float64	float64	float64	float64	float64	float64	float64	float64	float64
investment_4	float64	float64	float64	float64	float64	float64	float64	float64	float64	float64	float64	float64
debit_amount_apr	float64	float64	float64	float64	float64	float64	float64	float64	float64	float64	float64	float64
credit_amount_apr	float64	float64	float64	float64	float64	float64	float64	float64	float64	float64	float64	float64
debit_count_apr	float64	float64	float64	float64	float64	float64	float64	float64	float64	float64	float64	float64
credit_count_apr	float64	float64	float64	float64	float64	float64	float64	float64	float64	float64	float64	float64
max_credit_amount_apr	float64	float64	float64	float64	float64	float64	float64	float64	float64	float64	float64	float64
debit_amount_may	float64	float64	float64	float64	float64	float64	float64	float64	float64	float64	float64	float64
credit_amount_may	float64	float64	float64	float64	float64	float64	float64	float64	float64	float64	float64	float64
debit_count_may	float64	float64	float64	float64	float64	float64	float64	float64	float64	float64	float64	float64
credit_count_may	float64	float64	float64	float64	float64	float64	float64	float64	float64	float64	float64	float64
max_credit_amount_may	float64	float64	float64	float64	float64	float64	float64	float64	float64	float64	float64	float64
debit_amount_jun	float64	float64	float64	float64	float64	float64	float64	float64	float64	float64	float64	float64
credit_amount_jun	float64	float64	float64	float64	float64	float64	float64	float64	float64	float64	float64	float64
debit_count_jun	float64	float64	float64	float64	float64	float64	float64	float64	float64	float64	float64	float64
credit_count_jun	float64	float64	float64	float64	float64	float64	float64	float64	float64	float64	float64	float64
max_credit_amount_jun	float64	float64	float64	float64	float64	float64	float64	float64	float64	float64	float64	float64
loan_enq	object	object	object	object	object	object	object	object	object	object	object	object
emi_active	float64	float64	float64	float64	float64	float64	float64	float64	float64	float64	float64	float64
dtype:	object	object	object	object	object	object	object	object	object	object	object	object

## Step 2: Data cleaning/formatting

We are going to remove abnormally high/low values from the dataset in order to make sure we can take out better insights from the data. The columns in the dataset are going to be converted into proper format and all irregularities would be removed.

```
In [5]: #dropping missing values
bod.isnull().sum()
```

```
Out[5]:
```

	ID	cc_cons_apr	dc_cons_apr	cc_cons_may	dc_cons_may	cc_cons_jun	dc_cons_jun	cc_count_apr	cc_count_may	cc_count_jun	...	credit
cc_cons_apr	0	0	0	0	0	0	0	0	0	0	0	0
dc_cons_apr	0	0	0	0	0	0	0	0	0	0	0	0
cc_cons_may	0	0	0	0	0	0	0	0	0	0	0	0
dc_cons_may	0	0	0	0	0	0	0	0	0	0	0	0
cc_cons_jun	0	0	0	0	0	0	0	0	0	0	0	0
dc_cons_jun	0	0	0	0	0	0	0	0	0	0	0	0
cc_count_apr	1	1	1	1	1	1	1	1	1	1	1	1
cc_count_may	0	0	0	0	0	0	0	0	0	0	0	0
cc_count_jun	0	0	0	0	0	0	0	0	0	0	0	0
dc_count_apr	0	0	0	0	0	0	0	0	0	0	0	0
dc_count_may	0	0	0	0	0	0	0	0	0	0	0	0
dc_count_jun	0	0	0	0	0	0	0	0	0	0	0	0
card_lim	0	0	0	0	0	0	0	0	0	0	0	0
personal_loan_active	0	0	0	0	0	0	0	0	0	0	0	0
vehicle_loan_active	0	0	0	0	0	0	0	0	0	0	0	0
personal_loan_closed	1	1	1	1	1	1	1	1	1	1	1	1
vehicle_loan_closed	0	0	0	0	0	0	0	0	0	0	0	0
investment_1	0	0	0	0	0	0	0	0	0	0	0	0
investment_2	0	0	0	0	0	0	0	0	0	0	0	0
investment_3	2	2	2	2	2	2	2	2	2	2	2	2
investment_4	0	0	0	0	0	0	0	0	0	0	0	0
debit_amount_apr	0	0	0	0	0	0	0	0	0	0	0	0
credit_amount_apr	0	0	0	0	0	0	0	0	0	0	0	0
debit_count_apr	0	0	0	0	0	0	0	0	0	0	0	0
credit_count_apr	0	0	0	0	0	0	0	0	0	0	0	0
max_credit_amount_apr	0	0	0	0	0	0	0	0	0	0	0	0
debit_amount_may	0	0	0	0	0	0	0	0	0	0	0	0
credit_amount_may	0	0	0	0	0	0	0	0	0	0	0	0
debit_count_may	0	0	0	0	0	0	0	0	0	0	0	0
credit_count_may	0	0	0	0	0	0	0	0	0	0	0	0
max_credit_amount_may	0	0	0	0	0	0	0	0	0	0	0	0
debit_amount_jun	0	0	0	0	0	0	0	0	0	0	0	0
credit_amount_jun	0	0	0	0	0	0	0	0	0	0	0	0
debit_count_jun	0	0	0	0	0	0	0	0	0	0	0	0
credit_count_jun	0	0	0	0	0	0	0	0	0	0	0	0
max_credit_amount_jun	0	0	0	0	0	0	0	0	0	0	0	0
loan_enq	2	2	2	2	2	2	2	2	2	2	2	2
emi_active	1	1	1	1	1	1	1	1	1	1	1	1
dtype:	int64	int64	int64	int64	int64	int64	int64	int64	int64	int64	int64	int64

```
In [6]: #dropping null values
bod=bod.dropna()
bod
```

```
Out[6]:
```

	ID	cc_cons_apr	dc_cons_apr	cc_cons_may	dc_cons_may	cc_cons_jun	dc_cons_jun	cc_count_apr	cc_count_may	cc_count_jun	...	credit
0	17051	3412.60	2909.34	1077.36	3084.50	4615.97	6693.0	19.0	10.0	10.0	2.0	...
1	11491	18133.760	3877.00	6123.78	5597.43	64620.00	6117.00	4.0	60	60	1	...
2	7433	6363.65	735.49	8799.00	13768.00	38266.00	2638.0	2.0	83	47	...	...
3	14606	12765.660	4429.16	16745.86	6360.00	29063.97	3711.0	12.0	2	16	...	...
4	8381	27819.70	1944.00	7006.50	2228.50	1096.25	1065.0	15.0	13	61	...	...

5 rows x 39 columns

```
In [7]: #dropping irrelevant columns
columns_to_drop = ['loan_enq', 'personal_loan_active', 'vehicle_loan_active', 'personal_loan_closed', 'vehicle_loan_closed']
bod = bod.drop(columns=columns_to_drop)
# Display the modified DataFrame
bod.head()
```

```
Out[7]:
```

	ID	cc_cons_apr	dc_cons_apr	cc_cons_may	dc_cons_may	cc_cons_jun	dc_cons_jun	cc_count_apr	cc_count_may	cc_count_jun	...	credit
0	17051	3412.60	2909.34	1077.36	3084.50	4615.97	6693.0	19.0	10.0	10.0	2.0	...
1	11491	18133.760	3877.00	6123.78	5597.43	64620.00	6117.00	4.0	60	60	1	...
2	7433	6363.65	735.49	8799.00	13768.00	38266.00	2638.0	2.0	83	47	...	...
3	14606	12765.66	4429.16	16745.86	6360.00	29063.97	3711.0	12.0	2	16	...	...
4	8381	27819.70	1944.00	7006.50	2228.50	1096.25	1065.0	15.0	13	61	...	...

5 rows x 34 columns

```
In [8]: #making the data type appropriate
bod['ID'] = bod['ID'].astype(str)
```

```
In [9]: #seeing new datatypes
bod.dtypes
```

```
Out[9]:
```

	ID	cc_cons_apr	dc_cons_apr	cc_cons_may	dc_cons_may	cc_cons_jun	dc_cons_jun	cc_count_apr	cc_count_may	cc_count_jun	...	credit
cc_cons_apr	object	float64	float64	float64	float64	float64	float64	float64	float64	float64	float64	float64
dc_cons_apr	float64	float64	float64	float64	float64	float64	float64	float64	float64	float64	float64	float64
cc_cons_may	float64	float64	float64	float64	float64	float64	float64	float64	float64	float64	float64	float64
dc_cons_may	float64	float64	float64	float64	float64	float64	float64	float64	float64	float64	float64	float64
cc_cons_jun	float64	float64	float64	float64	float64	float64	float64	float64	float64	float64	float64	float64
dc_cons_jun	float64	float64	float64	float64	float64	float64	float64	float64	float64	float64	float64	float64
cc_count_apr	float64	float64	float64	float64	float64	float64	float64	float64	float64	float64	float64	float64
cc_count_may	float64	float64	float64	float64	float64	float64	float64	float64	float64	float64	float64	float64
cc_count_jun	float64	float64	float64	float64	float64	float64	float64	float64	float64	float64	float64	float64
dc_count_apr	float64	float64	float64	float64	float64	float64	float64	float64	float64	float64	float64	float64
dc_count_may	float64	float64	float6									



```
[83]: # Calculate RMSE on the test data
rmse_value = rmse(y_test, y_pred_best_tree)
print("RMSE on Test Set:", rmse_value)

RMSE on Test Set: 5549.876172232738
```

We can see that after Hyperparameter tuning we are able to bring down the error to almost half the original value and hence we must select this model to predict the missing values.

### Step 7:Model validation

This is the final step in this exercise and in this step we are going to predict the missing values of the dataset using the decision tree regressor model we perfected through hyperparameter tuning.

```
In [89]: #Getting the dependent variables
X_missing = missing_data[['cc_cons_apr', 'dc_cons_apr', 'cc_cons_may', 'dc_cons_may', 'cc_cons_jun', 'dc_cons_jun', 'debit_amount_apr', 'debit_count_apr', 'credit_count_apr', 'debit_amount_may', 'debit_count_may', 'credit_count_may', 'debit_amount_jun', 'debit_count_jun', 'credit_count_jun', 'emi_active', 'total_customer_investment', 'trans_cnt_apr', 'trans_cnt_may', 'trans_cnt_jun', 'card_utilization_apr', 'card_utilization_may', 'card_utilization_jun', 'Credit_Card_Utilization_Ratio']]

In [90]: # Predict the missing values
X_missing = missing_data[['cc_cons']] = best_tree.predict(X_missing)
missing_data['cc_cons']

C:\Users\anjoy\Documents\LocalTemp\ipykernel_12120\1594784132.py:2: SettingWithCopyWarning:
A value is trying to be set on a copy of a slice from a DataFrame.
Try using .loc[row_indexer,col_indexer] = value instead
See the caveats in the documentation: https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#returning-a-view-versus-a-copy
missing_data['cc_cons'] = best_tree.predict(X_missing)
14995    12708.0
14996    11711.0
14997    13396.5
14998    12708.0
14999    13396.5
...
19988    13396.5
19989    12708.0
19990    12711.5
19991    13396.5
19992    12708.0
Name: cc_cons, Length: 4998, dtype: float64
```

```
In [ ]: # Output data frame with predicted values
output_df = missing_data[['ID', 'cc_cons']]
output_df.head()
```

```
In [95]: # Making the final dataset
final_data = pd.concat([known_data, output_df], axis=0)

# Reset the index of the final DataFrame
final_data.reset_index(drop=True, inplace=True)
```

```
In [96]: # Saving the final DataFrame to an Excel file
final_data.to_excel('final_data_with_missing_predictions.xlsx', index=False)
```

```
In [ ]:
```