



Micro Credit Loan

Submitted by:

Mahaboob Basha Shaik

ACKNOWLEDGMENT

I would like to express my special thanks of gratitude to my SME Mr. Shubham Yadav as well Fliprobo who gave me the golden opportunity to do this wonderful project on the Micro Credit Loan Project, which also helped me in doing a lot of Research and i came to know about so many new things I am thankful to them.

Secondly, I would also like to thank my parents and friends who helped me a lot in finalizing this project within the limited time frame.

INTRODUCTION

- **Business Problem Framing**

Identify the target customers who are credible to payback the loans which can be identified with their usage of their mobile phone recharges and monthly balance maintenance.

- **Conceptual Background of the Domain Problem**

Communication has become important factor in everyone day to day life. MFI wants to collaborate with telecom industry to reach out all poor customers and low-income families who can use this opportunity and payback the sanctioned amount in stipulated time frame.

- **Review of Literature**

we study and discuss the MFI and Telecom Credit strategies in a telecommunication industry. The goal of this study is to apply various analysis techniques along with data mining techniques to predict customer credibility in repay of loans taken from a telecommunications company. The study will help telecommunications companies to understand customer repayable capacity.

- **Motivation for the Problem Undertaken**

MFI working to Telecom Industry to provide loans to low level income which help in not happening their connectivity with their near and dear once.

Analytical Problem Framing

- Mathematical/ Analytical Modeling of the Problem

Linear Regression is applied for the project completion. In which I used Lasso Regressor, Ridge Regressor, Decision Tree, Random Forest Regressor and GridsearchCV for finding the best model.

- Data Sources and their formats

	int64	object	float64
Count	13	3	21

```
: Unnamed: 0      int64
label            int64
msisdn          object
aon             float64
daily_decr30    float64
daily_decr90    float64
rental30        float64
rental90        float64
last_rech_date_ma float64
last_rech_date_da float64
last_rech_amt_ma int64
cnt_ma_rech30    int64
fr_ma_rech30     float64
sumamnt_ma_rech30 float64
medianamnt_ma_rech30 float64
medianmarechprebal30 float64
cnt_ma_rech90    int64
fr_ma_rech90     int64
sumamnt_ma_rech90 int64
medianamnt_ma_rech90 float64
medianmarechprebal90 float64
cnt_da_rech30    float64
fr_da_rech30     float64
cnt_da_rech90    int64
fr_da_rech90     int64
cnt_loans30      int64
amnt_loans30     int64
maxamnt_loans30  float64
medianamnt_loans30 float64
cnt_loans90      float64
amnt_loans90     int64
maxamnt_loans90  int64
medianamnt_loans90 float64
payback30        float64
payback90        float64
pcircle          object
pdate           object
dtype: object
```

- **Data Pre-processing Done**

Verifying the dataset whether it is having NULL values and checking the data types of the provided data. Converted all the float data type to Int data type for modelling and removing the skewness / outliers from computing.

- **Data Inputs- Logic- Output Relationships**

We had considered y variable as Label as it explains whether the customer had made the repayment of loan. X variable consists of other variables which consist of Avg of 30 days and 90 days of parameters which directly or indirectly influences the y variable.

- **Hardware and Software Requirements and Tools Used**

- a. import numpy as np
- b. import pandas as pd
- c. import sklearn
- d. import seaborn as sb
- e. import matplotlib.pyplot as plt
- f. from scipy.stats import zscore
- g. import numpy as np
- h. from sklearn.preprocessing import StandardScaler
- i. from sklearn.preprocessing import power_transform
- j. from sklearn.model_selection import train_test_split
- k. from sklearn.linear_model import LinearRegression
- l. from sklearn.metrics import r2_score
- m. from sklearn.model_selection import GridSearchCV
- n. from sklearn.model_selection import cross_val_score
- o. from sklearn.linear_model import Ridge
- p. from sklearn.tree import DecisionTreeRegressor
- q. from sklearn.ensemble import RandomForestRegressor
- r. from sklearn.ensemble import GradientBoostingRegressor
- s. from sklearn.svm import SVR

Model/s Development and Evaluation

- Identification of possible problem-solving approaches (methods)

Decision Tree is the best possible model in solving the model as it has cross validation and r2 score of 100. All data points are interacting with the predicted model.

- Testing of Identified Approaches (Algorithms)

Lasso Regressor, Ridge Regressor, Decision Tree Regressor, Random Forest Regressor, GridsearchCV

- Run and Evaluate selected models

Regularisation - Lasso

```
1 from sklearn.linear_model import Lasso
2 parameters = {'alpha': [.00001, .0001, .001, .01, .1, 1, 10], 'random_state': list(range(0,10))}
3 ls=Lasso()
4 clf=GridSearchCV(ls,parameters)
5 clf.fit(x_train,y_train)
6 print (clf.best_params_)
```

```
{'alpha': 1e-05, 'random_state': 0}
```

```
1 ls=Lasso(alpha=0.0001,random_state=0)
2 ls.fit(x_train,y_train)
3 ls.score(x_train,y_train)
4 pred_ls=ls.predict(x_test)
5
6 lss=r2_score(y_test,pred_ls)
7 for j in range(2,10):
8     lsscore = cross_val_score(ls,x_t,y,cv=j)
9     lsc=lsscore.mean()
10    print("At CV :-",j)
11    print("Cross Validation Score is :-",lsc*100)
12    print ("R2_score is :-",lss*100)
13    print('\n')
```

```
At CV :- 2
Cross Validation Score is :- 99.99999084559663
R2_score is :- 99.99999083206609
```

```
At CV :- 3
Cross Validation Score is :- 99.99999084584479
R2_score is :- 99.99999083206609
```

```

1 from sklearn.metrics import mean_squared_error,mean_absolute_error
2 print("Error:")
3 print("Mean Absolute Error:",mean_absolute_error(y_test,pred_ls))
4 print("Mean Square Error:",mean_squared_error(y_test,pred_ls))
5 print("Root Mean Squared Error:", np.sqrt(mean_squared_error(y_test,pred_ls)))

```

Error:
Mean Absolute Error: 6.62293035516008e-05
Mean Square Error: 1.0045669535635966e-08
Root Mean Squared Error: 0.0001002280875585081

```

1 plt.figure(figsize=(8,6))
2 plt.scatter(x=y_test,y=pred_ls,color='r')
3 plt.plot(y_test,y_test,color='b')
4 plt.xlabel('Actual Credit Paid in 5 Days',fontsize=14)
5 plt.ylabel('Predicted Credit Paid in 5 Days',fontsize=14)
6 plt.title('Lasso Regression',fontsize=18)
7 plt.show()

```



Data Point is interacting with the predicted line

Ridge Regression

```

1 from sklearn.linear_model import Ridge
2 parameters = {'alpha': [.0001, .001, .01, .1, 1], 'fit_intercept': [True, False], 'normalize': [True, False], 'random_state': [1, 2, 3, 4, 5, 6]}
3 rd=Ridge()
4 clf=GridSearchCV(rd,parameters)
5 clf.fit(x_train,y_train)
6 print (clf.best_params_)

```

```
{'alpha': 0.0001, 'fit_intercept': True, 'normalize': False, 'random_state': 1}
```

```

1 ridge=Ridge(alpha=0.0001,random_state=1,fit_intercept=True,normalize=False)
2 rd.fit(x_train,y_train)
3 rd.score(x_train,y_train)
4 pred_rd=rd.predict(x_test)
5
6 rdd=r2_score(y_test,pred_rd)
7 for j in range(2,10):
8     rdscore = cross_val_score(rd,x_t,y,cv=j)
9     rdc=rdscore.mean()
10    print("At CV :-",j)
11    print("Cross Validation Score is :-",rdc*100)
12    print ("R2_score is :-",rdc*100)
13    print('\n')

```

At CV :- 2
Cross Validation Score is :- 99.9999998983795
R2_score is :- 99.9999998983795

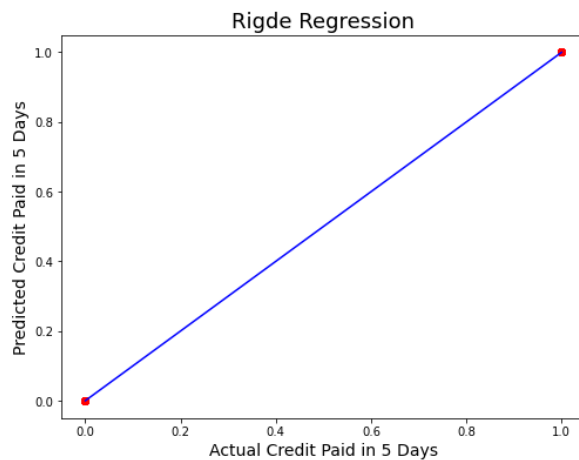
At CV :- 3
Cross Validation Score is :- 99.9999999428626
R2_score is :- 99.9999999428626

At CV :- 4
Cross Validation Score is :- 99.999999954857
R2_score is :- 99.999999954857

```

1 plt.figure(figsize=(8,6))
2 plt.scatter(x=y_test,y=pred_rd,color='r')
3 plt.plot(y_test,y_test,color='b')
4 plt.xlabel('Actual Credit Paid in 5 Days',fontsize=14)
5 plt.ylabel('Predicted Credit Paid in 5 Days',fontsize=14)
6 plt.title('Rigde Regression',fontsize=18)
7 plt.show()

```



Same as Lasso Model all Data points are interacting with the predictd line

Ensembler Technique - Decision Tree

```

1 from sklearn.tree import DecisionTreeRegressor
2 parameters={'criterion':['mse','friedman_mse','mse'],'splitter':['best','random']}
3 dt=DecisionTreeRegressor()
4 clf=GridSearchCV(dt,parameters)
5 clf.fit(x_train,y_train)
6 print (clf.best_params_)

```

```
{'criterion': 'mse', 'splitter': 'best'}
```

```

1 dt=DecisionTreeRegressor(criterion='mse', splitter='best')
2 dt.fit(x_train,y_train)
3 dt.score(x_train,y_train)
4 pred_decision=dt.predict(x_test)
5 dts=r2_score(y_test,pred_decision)
6 print("r2_score:",dts*100)
7 dtscore = cross_val_score(dt,x_t,y,cv=3)
8 dtc=dtscore.mean()
9 print('Cross Val Score:',dtc*100)

```

r2_score: 100.0

Cross Val Score: 100.0

```

1 print("Error:")
2 print("Mean Absolute Error:",round(mean_absolute_error(y_test,pred_decision),2))
3 print("Mean Square Error:",round(mean_squared_error(y_test,pred_decision),2))
4 print("Root Mean Sqaured Error:",round(np.sqrt(mean_squared_error(y_test,pred_decision)),2))

```

```

Error:
Mean Absolute Error: 0.0
Mean Square Error: 0.0
Root Mean Sqaured Error: 0.0

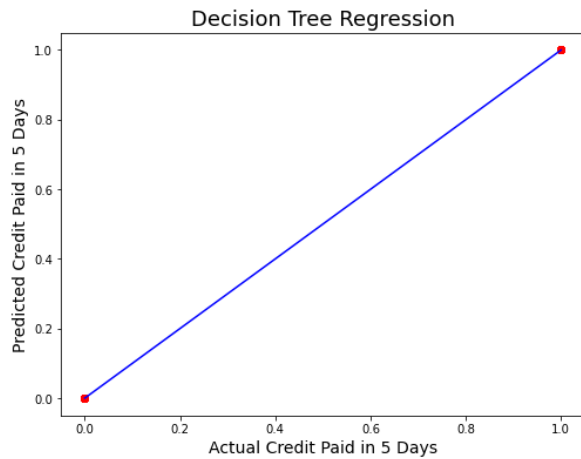
```



```

1 plt.figure(figsize=(8,6))
2 plt.scatter(x=y_test,y=pred_decision,color='r')
3 plt.plot(y_test,y_test,color='b')
4 plt.xlabel('Actual Credit Paid in 5 Days',fontsize=14)
5 plt.ylabel('Predicted Credit Paid in 5 Days',fontsize=14)
6 plt.title('Decision Tree Regression',fontsize=18)
7 plt.show()

```



R2 score and CV Score are 100 , All The Points are interacting with predicted line

Random Forest

```

1 from sklearn.ensemble import RandomForestRegressor
2 parameters={'criterion':['mse','friedman_mse','mse'],'n_estimators':[100,200,300]}
3 rf=RandomForestRegressor()
4 clf=GridSearchCV(rf,parameters)
5 clf.fit(x_train,y_train)
6 print (clf.best_params_)

```

```
{'criterion': 'mse', 'n_estimators': 100}
```

```

1 rf=RandomForestRegressor(criterion='mse', n_estimators= 100)
2 rf.fit(x_train,y_train)
3 rf.score(x_train,y_train)
4 pred_rd=rf.predict(x_test)
5 rfs=r2_score(y_test,pred_rd)
6 print("r2_score:",rfs*100)
7 rfscore = cross_val_score(rf,x_t,y,cv=3)
8 rfc=rfscore.mean()
9 print('Cross Val Score:',rfc*100)

```

```
r2_score: 100.0
```

```
Cross Val Score: 100.0
```

```

1 plt.figure(figsize=(8,6))
2 plt.scatter(x=y_test,y=pred_rd,color='r')
3 plt.plot(y_test,y_test,color='b')
4 plt.xlabel('Actual Credit Paid in 5 Days',fontsize=14)
5 plt.ylabel('Predicted Credit Paid in 5 Days',fontsize=14)
6 plt.title('Random Forest Regressor',fontsize=18)
7 plt.show()

```



In Random Forest Regressor also Data points are intracting with the predicted lines

Gradient Bosster Regressor

```

1 from sklearn.ensemble import GradientBoostingRegressor
2 parameters={'loss':['ls','lad','huber','quantile'],'n_estimators':[50,100,200],'criterion':['friedman_mse','mse']}
3 gbr=GradientBoostingRegressor()
4 clf=GridSearchCV(gbr,parameters)
5 clf.fit(x_train,y_train)
6 print(clf.best_params_)
7

```

```
{'criterion': 'friedman_mse', 'loss': 'ls', 'n_estimators': 200}
```

```

1 gbr=GradientBoostingRegressor(criterion='friedman_mse', loss='ls',n_estimators=100)
2 gbr.fit(x_train,y_train)
3 gbr.score(x_train,y_train)
4 pred_random=gbr.predict(x_test)
5 gbrs=r2_score(y_test,pred_random)
6 print("r2_score:",round(gbrs*100,2))
7 gbscore = cross_val_score(gbr,x_t,y,cv=3)
8 gbrc=gbscore.mean()
9 print('Cross Val Score:',round(gbrc*100,2))
10

```

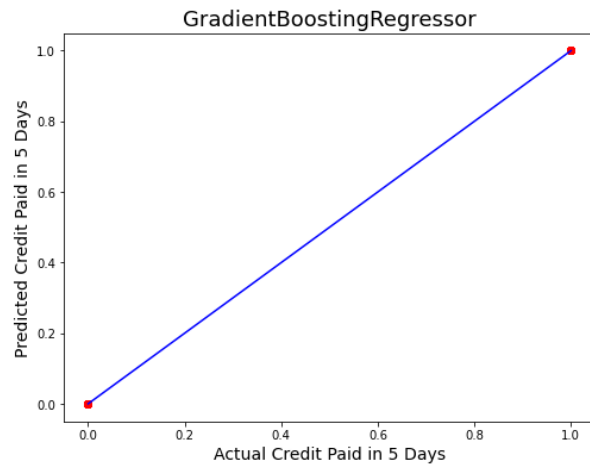
r2_score: 100.0

Cross Val Score: 100.0

```

1 plt.figure(figsize=(8,6))
2 plt.scatter(x=y_test,y=pred_random,color='r')
3 plt.plot(y_test,y_test,color='b')
4 plt.xlabel('Actual Credit Paid in 5 Days',fontsize=14)
5 plt.ylabel('Predicted Credit Paid in 5 Days',fontsize=14)
6 plt.title('GradientBoostingRegressor',fontsize=18)
7 plt.show()
8

```



All Data points in this model also interseting with Predicted Line

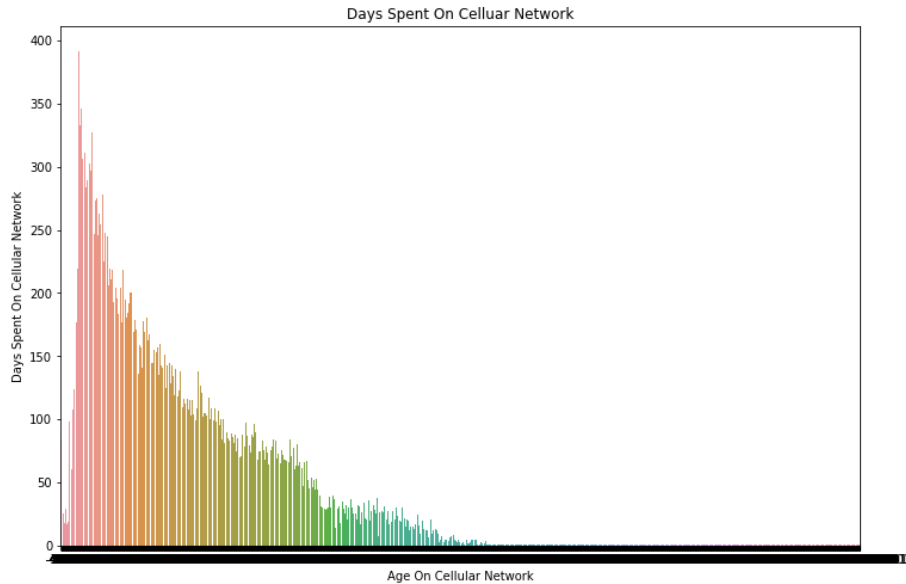
- Key Metrics for success in solving problem under consideration.

Data type conversion, Elimination of outliers and skewed data from dataset.

- Visualizations

Uni Variant Ananlysis

```
: 1 plt.subplots(figsize=(12,8))
2 sb.countplot(x="aon",data=df)
3 plt.title("Days Spent On Celluar Network")
4 plt.xlabel('Age On Cellular Network')
5 plt.ylabel('Days Spent On Cellular Network')
6 plt.show()
```



Highest days spent on network is between bin is in 350 - 400

- Interpretation of the Results

After all models in Linear Regression, I found Decision Tree Regression provide accurate prediction for the test data. Please find the below:

Ensembler Technique - Decision Tree

```
: 1 from sklearn.tree import DecisionTreeRegressor
2 parameters={'criterion':['mse','friedman_mse','mse'],'splitter':['best','random']}
3 dt=DecisionTreeRegressor()
4 clf=GridSearchCV(dt,parameters)
5 clf.fit(x_train,y_train)
6 print (clf.best_params_)
```

```
{'criterion': 'mse', 'splitter': 'best'}
```

```
: 1 dt=DecisionTreeRegressor(criterion='mse', splitter='best')
2 dt.fit(x_train,y_train)
3 dt.score(x_train,y_train)
4 pred_decision=dt.predict(x_test)
5 dts=r2_score(y_test,pred_decision)
6 print("r2_score:",dts*100)
7 dt_score = cross_val_score(dt,x_t,y,cv=3)
8 dtc=dt_score.mean()
9 print('Cross Val Score:',dtc*100)
```

```
r2_score: 100.0
```

```
Cross Val Score: 100.0
```

```
: 1 print("Error:")
2 print("Mean Absolute Error:",round(mean_absolute_error(y_test,pred_decision),2))
3 print("Mean Square Error:",round(mean_squared_error(y_test,pred_decision),2))
4 print("Root Mean Sqaured Error:",round(np.sqrt(mean_squared_error(y_test,pred_decision)),2))
```

```
Error:
```

```
Mean Absolute Error: 0.0
```

```
Mean Square Error: 0.0
```

```
Root Mean Sqaured Error: 0.0
```

```
: 1 plt.figure(figsize=(8,6))
2 plt.scatter(x=y_test,y=pred_decision,color='r')
3 plt.plot(y_test,y_test,color='b')
4 plt.xlabel('Actual Credit Paid in 5 Days',fontsize=14)
5 plt.ylabel('Predicted Credit Paid in 5 Days',fontsize=14)
6 plt.title('Decision Tree Regression',fontsize=18)
7 plt.show()
```



R2 score and CV Score are 100 , All The Points are interacting with predicted line

Testing

```
1 import numpy as np
2 a=np.array(y_test)
3 a
```

array([1, 0, 0, ..., 1, 1, 0], dtype=int64)

```
1 df_com=pd.DataFrame({"Original":a,"Predicted":pred_decision},index=range(len(a)))
2 df_com
```

	Original	Predicted
0	1	1.0
1	0	0.0
2	0	0.0
3	1	1.0
4	1	1.0
...
41914	1	1.0
41915	1	1.0
41916	1	1.0
41917	1	1.0
41918	0	0.0

41919 rows x 2 columns

CONCLUSION

- Key Findings and Conclusions of the Study

Key Finding: Data is having high outliers and data missing in few rows.

Inference: User who are active and maintaining avg balance for 30days and 90days are more likely to repay the loan.

Conclusion: Data need to be more specific in type of recharges like talktime, data, top up recharges or SMS Recharge. Specificity of the data is missing. Considered all the balance maintains is for talk time only.

- Learning Outcomes of the Study in respect of Data Science

Challenges: Handling and treating skewed data / outliers.

Visualisation: It help in understanding the frequency of users who are able to manage 30days / 90days monthly balance, active days on network etc

- Limitations of this work and Scope for Future Work

Data doesn't specifies weather the recharges or balance maintains is of talktime or data balance. If the data is more specific in this instance it would be easier to target audience more accurately.

Scope of future work: It would be more helpful if we can identify the exacta expectation of the user, weather the user is looking for talktime / data / SMS. Introduce packs specific loans can also help in increasing the usage of loan for the credible users,