

Peer Lending Project

Stage 3 - Data Preparation/Modeling

Group C

In this stage we prepared our data for modeling using methods to deal with missing values, outliers and encoding of categorical variables as well as checking the option of scaling the data. We started creating our models for proof of concept.

The steps we took:

1. Changes from the previous stage:

- a. We returned two attributes that we dropped due to irrelevance - "earliest_cr_line" and "issue_d", we used those to create a new attribute "cr_hist" - the length of a person's credit history at the time the loan is issued. We will check if the new attribute improves our model.
- b. We decided to drop the "emp_title" attribute since it has too many unique values (approx 80,000) and will not help our model.

2. Dealing with missing values:

We treated the attributes by their N/A Fraction:

- Less than 1% - dropped those N/As, since it's negligent.
- More than 1% - the traditional approach : replacing with the median, which is not biased by the values at the far end of the distribution, therefore, the median is a good representation of the majority of the values in the variable.

We've decided on those methods in order not to lose many instances in our initial modeling stage. In the next step we will try to use a KNN-imputation to replace missing values (also the ones we dropped for now) and see if it improves our model (See Appendix 1 for the N/A fractions).

3. Dealing with categorical attributes:

We had to change the categorical attributes since Machine Learning algorithms require that attributes will be numeric. We used few different methods to encoding those attributes:

- a. "loan_status" & "initial_list_status": _____

These attributes have only 2 possible values, therefore, they will be encoded with the values as "0" and "1" (Appendix 2.1).

- b. Encoding of categorical attributes with more than 2 unique values:

- **Hierarchical attributes** - "grade", "purpose" (after manipulation).
Encoded by their order (Appendix 2.1, 2.6 for "grade" and Appendix 2.1, 2.3, 2.4, 2.5 for "purpose").
- **Non Hierarchical attributes** - "addr_state", "home_ownership", "verification_status"
One-hot encoding: Some attributes had many possible values so we grouped them to decrease the number of values (Appendix 2.1 and Appendix 2.2 for "home_ownership").

As we move forward if an encoded attribute will not contribute to our model, it will be dropped.

4. Dealing with outliers:

Before detecting outliers we want to mention that the distribution of the realized return between fully paid and defaulted loans - 99% of fully paid loans have realized return of more than 2% whereas only 5% of the defaulted loans yield more than 2%. This is worth mentioning since the company currently yields around 2% yearly and we want to detect the loans with a better return.

We understand that the majority of loans we want to detect and avoid are defaulted, with only 5% of them with a better realized return than the company generates now. This means that if we do a good job detecting fully paid loans, 99% of them will yield more than 2%.

We used both the empirical¹ and the IQR² rule to detect the absolute number and fraction of outliers for each attribute (Appendix 3.1).

We saw that when using the IQR rule we have much more outliers, hence removing outliers based on this criteria will result in a massive loss of data. Therefore, we will focus on the empirical rule results.

To understand the nature of the outliers better, we checked the division between fully paid and defaulted loans for every attribute's outliers (Appendix 3.2). Our target column - realized_return, has the highest number of outliers, however, all of them are defaulted loans - as we mentioned above these are the loans we want to detect, removing them will skew the data even more than it already is (only 13% defaulted). Therefore we will not ignore or replace these outliers.

Outliers are a sensitive subject. On one hand, they might bias our future model but on the other hand, they might contain relevant information. For now, we will not remove or change outliers, we will check each attribute's outliers separately, and treat them only if it improves our model before making any hasty decisions.

5. Scaling

Scaling the data is an important step since some models require the data to be scaled so it will not assign unproportional weights for attributes. For now we are testing our models with and without scaling, at this time we will use a standard scaler, when we proceed we will try different scaling methods and see if it improves our models.

6. Modeling:

Approach A - Regression model for predicting realized return:

We used linear regression and Gradient Boosting regression to predict the realized return for each loan. We then compared the results to our baseline model - investing only based on grade (appendix 4.1). Using grade we can choose loans that fit our desired risk. Loans graded 'A' has the highest average return and the lowest standard deviation - thus being the safest investment. In the second extreme, loans graded 'F' and 'G' offer average return close to zero but high standard deviation - meaning a chance for high return. Our goal is either to predict loans with the same realized return and lower standard deviation or a higher return with the same standard deviation.

Key evaluation metrics:

RMSE - it indicates the absolute fit of the model to the data – how close the observed data points are to the model's predicted values

R² - the proportion of the variance in the dependent variable that is predictable (explained) from the independent variables.

¹ 3 SD from the mean

² below Q1 - 1.5 IQR or above Q3 + 1.5 IQR

Both models performed poorly. The RMSE is too large, meaning that many loans with predicted positive return can have negative return, exactly the type of error we want to avoid - hence we can't trust the predictions. We will attempt to improve our model and consider dropping it if we can't reach better results.

Approach B - Classification model on loan status:

We used 3 different classification methods to determine whether a loan is fully paid or defaulted. We then created an expected return equation as follows:

$$E(\text{loans}) = P(\text{fully paid}) * E(\text{return fully paid}) + P(\text{default}) * E(\text{return default})$$

Where:

$P(\text{fully paid})$ - The proportion of loans predicted as fully paid that were actually fully paid

$P(\text{default})$ - The proportion of loans predicted as fully paid that were actually default (1 - precision)

$E(\text{return fully paid})$ - The average realized return of fully paid loans calculated from the training set

$E(\text{return default})$ - The average realized return of defaulted loans calculated from the training set

In our equation, we multiply $P(\text{fully paid})$ with the mean realized return of fully paid loans and deduct ("punish") $P(\text{default})$ multiplied by the mean realized return of defaulted loans to estimate the average return for loans our model predicts as fully paid. Using this method, we don't predict the realized return per loan, rather we show the average expected realized return for investing in all loans our model predicts as fully paid. As we mentioned above, 99% of fully paid loans yield more than 2%. Thus, predicting correctly the maturity of the loans can yield a higher return. We decided to try this different approach since the regression models performed poorly, this might suggest that our data is not suitable for regression.

We will use cross validation and other metrics to raise our confidence in the estimate of the precision our model produces since it's crucial for our expected return equation. (Appendix 4.2)

Key evaluation metric:

Recall - The ratio of correct positive predictions to the total positives examples. Defaulted loans classified as fully paid generate losses where fully paid loans classified as default prevent gain. One way to increase recall is to lower the threshold for classifying loans as default, however this will also result in loss of fully paid loans with positive return. We will try different thresholds and choose the one with the highest return suitable to a specific budget constraint.

AUC- The Area Under the Curve is the measure of the ability of a classifier to distinguish between classes (loan status) better than a random model.

Appendices

Appendix 1 - missing values Fraction

Columns	Missing value fraction
total_bal_il	0.0166
max_bal_bc	0.0166
all_util	0.0226
revol_util	0.0667
realized_return	0.1514
bc_open_to_buy	1.1464
percent_bc_gt_75	1.1551
bc_util	1.1969
il_util	13.9307
emp_length	14.8805

Appendix 2: Dealing with Categorical attributes

Appendix 2.1 - Encoding description for categorical attributes

Attribute	Encoding Description
"loan_status"	Can get the values Fully Paid\Charged Off
"initial_list_status"	Can get the values w\f
"verification_status"	We will merge "Source Verified" and "Verified" to "Verified" since we believe they have the same meaning. we will encode it as 0 and 1
"grade"	<u>Hierarchical attribute</u> . Since the different grades are ordered where A is the highest, we will replace G to A with 1 to 7
'purpose'	The values "wedding" and "renewable_energy" appear 1 and 204 times respectively in our data. Since these won't be helpful for our model, we decided to remove the rows containing these values. For each of the remaining values, we checked their default rate as well as the distribution of their realized return. The results showed that the realized return distribution is more or less the same for each

	purpose while the default rate can be divided into three subgroups, "under 16.36%", "between 16.36% and 19.5%" and "over 19.5%". We decided to group the values into these 3 subgroups, labeled 1,2 and 3 (Appendix 2.3). After this process, "purpose" is a <u>hierarchical</u> attribute
"addr_state"	Contains values for 50 states so we will group the states to regions and use <u>one hot encoding</u>
"home_ownership"	Has 4 unique values, the value "ANY" has 58 instances, it is negligible so we will drop these instances and use <u>one hot encoding</u> for the column
"sub_grade"	There is an overlap of information with the rankings column. In addition it contains 35 different variables. We will currently drop the column and return it if necessary according to the model's performance.

Appendix 2.2 - Home Ownership Attribute values count:

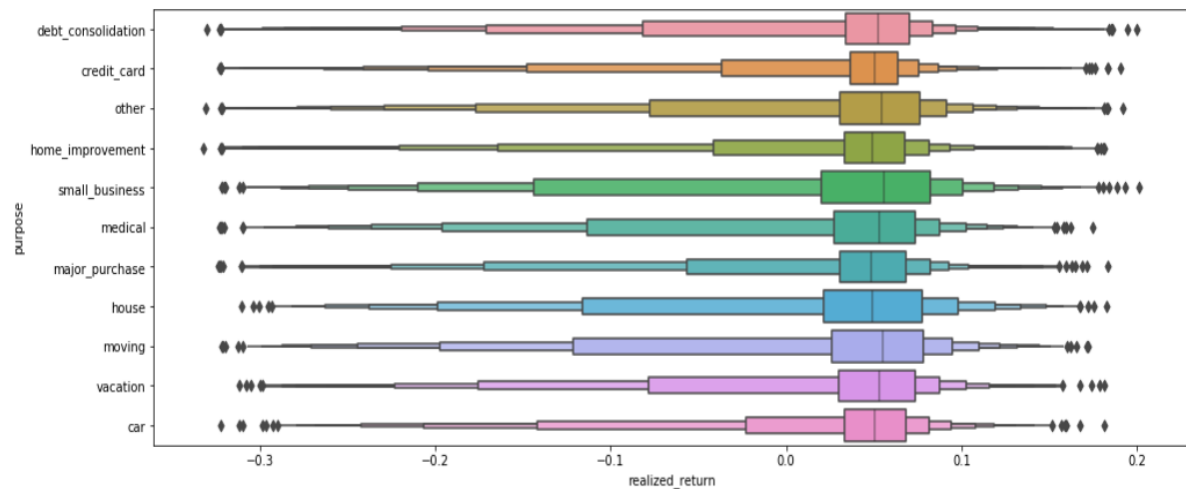
Home ownership	Value count
ANY	58
MORTGAGE	121648
OWN	33539
RENT	109574

Appendix 2.3 - Purpose Attribute values count:

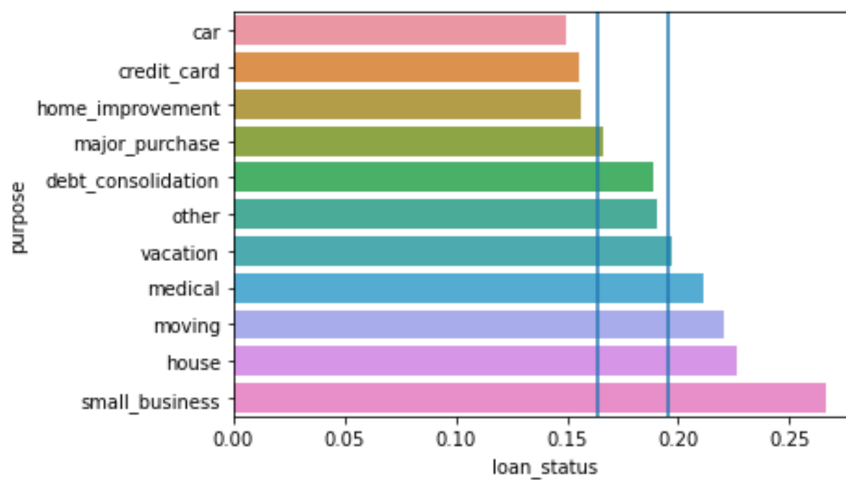
Purpose	Value count
car	3305
credit_card	57098
debt_consolidation	146266
home_improvement	19116
house	1259
major_purchase	6749
medical	3618
moving	2328
other	19380
Renewable_energy	204
small_business	3029
vacation	2408
wedding	1

Appendix 2.4 - Realized return distribution for each purpose

We can see the realized return distribution is more or a less the same for each purpose.



Appendix 2.5 - Loan status for each purpose



car, credit_card, home_improvement, major_purchase → encode 1

debt_consolidation, other, vacation → encode 2

medical, moving, house, small_business → encode 3

Appendix 2.6 - Grades realized return distribution

Grade	Mean	Median	Std	Min	Max
G	-0.0039	0.0332	0.1541	-0.3197	0.1942
F	0.0081	0.0579	0.1453	-0.3305	0.2010
E	0.0158	0.0677	0.1306	-0.3215	0.1813
D	0.0189	0.0683	0.1138	-0.3214	0.1575
C	0.0248	0.0659	0.0936	-0.3319	0.1301
B	0.0301	0.0558	0.0714	-0.3303	0.1222
A	0.0320	0.0410	0.0461	-0.3227	0.1007

Appendix 3 - Outliers

Appendix 3.1: Method comparison

<u>Columns</u>	<u>IQR Outliers</u>	<u>Prop IQR Outliers</u>	<u>Gaussian Outliers</u>	<u>Prop Gauss Outliers</u>
<u>realized_return</u>	<u>44529</u>	<u>0.168315971</u>	<u>7546</u>	<u>0.028523262</u>
<u>bc_open_to_buy</u>	<u>22389</u>	<u>0.084628585</u>	<u>5726</u>	<u>0.021643811</u>
<u>pct_tl_nvr_dlt</u>	<u>15417</u>	<u>0.058274997</u>	<u>5161</u>	<u>0.019508157</u>

<u>total_bc_limit</u>	<u>16660</u>	<u>0.062973435</u>	<u>5147</u>	<u>0.019455238</u>
<u>total_bal_il</u>	<u>16640</u>	<u>0.062897836</u>	<u>5042</u>	<u>0.019058347</u>
<u>total_bal_ex_mort</u>	<u>15864</u>	<u>0.05996462</u>	<u>4876</u>	<u>0.01843088</u>
<u>num_bc_sats</u>	<u>13071</u>	<u>0.049407309</u>	<u>4785</u>	<u>0.018086908</u>
<u>avg_cur_bal</u>	<u>15722</u>	<u>0.059427872</u>	<u>4765</u>	<u>0.01801131</u>
<u>mort_acc</u>	<u>9679</u>	<u>0.036585827</u>	<u>4761</u>	<u>0.01799619</u>
<u>total_il_high_credit_limit</u>	<u>13743</u>	<u>0.051947414</u>	<u>4655</u>	<u>0.017595519</u>
<u>delinq_2yrs</u>	<u>55811</u>	<u>0.210961006</u>	<u>4327</u>	<u>0.016355705</u>
<u>total_rev_hi_lim</u>	<u>16444</u>	<u>0.062156972</u>	<u>4218</u>	<u>0.015943694</u>
<u>max_bal_bc</u>	<u>16918</u>	<u>0.063948654</u>	<u>4125</u>	<u>0.015592162</u>
<u>fico_range_high</u>	<u>9096</u>	<u>0.034382135</u>	<u>4090</u>	<u>0.015459865</u>
<u>fico_range_low</u>	<u>9096</u>	<u>0.034382135</u>	<u>4090</u>	<u>0.015459865</u>
<u>tot_cur_bal</u>	<u>11104</u>	<u>0.04197221</u>	<u>4068</u>	<u>0.015376707</u>
<u>num_actv_bc_tl</u>	<u>6088</u>	<u>0.023012141</u>	<u>3812</u>	<u>0.014409048</u>
<u>tot_hi_cred_lim</u>	<u>10836</u>	<u>0.040959192</u>	<u>3709</u>	<u>0.014019716</u>
<u>num_bc_tl</u>	<u>6121</u>	<u>0.023136878</u>	<u>3702</u>	<u>0.013993257</u>
<u>pub_rec</u>	<u>51963</u>	<u>0.196415882</u>	<u>3676</u>	<u>0.013894979</u>
<u>num_rev_tl_bal_gt_0</u>	<u>7165</u>	<u>0.027083113</u>	<u>3672</u>	<u>0.013879859</u>

<u>open_acc</u>	<u>6849</u>	<u>0.025888659</u>	<u>3661</u>	<u>0.01383828</u>
<u>revol_bal</u>	<u>18613</u>	<u>0.070355615</u>	<u>3648</u>	<u>0.013789141</u>
<u>num_op_rev_tl</u>	<u>12068</u>	<u>0.045616051</u>	<u>3643</u>	<u>0.013770241</u>
<u>num_sats</u>	<u>10208</u>	<u>0.038585403</u>	<u>3576</u>	<u>0.013516987</u>
<u>num_rev_accts</u>	<u>7244</u>	<u>0.027381726</u>	<u>3493</u>	<u>0.013203254</u>
<u>il_util</u>	<u>11673</u>	<u>0.044122983</u>	<u>3299</u>	<u>0.01246995</u>
<u>pub_rec_bankruptcies</u>	<u>35323</u>	<u>0.133518045</u>	<u>2841</u>	<u>0.010738747</u>
<u>int_rate</u>	<u>7003</u>	<u>0.026470766</u>	<u>2758</u>	<u>0.010425014</u>
<u>grade</u>	<u>11279</u>	<u>0.042633696</u>	<u>2575</u>	<u>0.009733289</u>
<u>cr_hist</u>	<u>8087</u>	<u>0.030568197</u>	<u>2525</u>	<u>0.009544293</u>
<u>tot_coll_amt</u>	<u>47489</u>	<u>0.179504528</u>	<u>2455</u>	<u>0.009279699</u>
<u>installment</u>	<u>9675</u>	<u>0.036570707</u>	<u>1954</u>	<u>0.00738596</u>
<u>annual_inc</u>	<u>13103</u>	<u>0.049528266</u>	<u>1881</u>	<u>0.007110026</u>
<u>loan_amnt</u>	<u>10069</u>	<u>0.038059995</u>	<u>1555</u>	<u>0.005877773</u>
<u>dti</u>	<u>906</u>	<u>0.003424606</u>	<u>729</u>	<u>0.00275556</u>
<u>all_util</u>	<u>1599</u>	<u>0.006044089</u>	<u>327</u>	<u>0.001236033</u>
<u>bc_util</u>	<u>18</u>	<u>6.80E-05</u>	<u>23</u>	<u>8.69E-05</u>
<u>revol_util</u>	<u>19</u>	<u>7.18E-05</u>	<u>19</u>	<u>7.18E-05</u>

<u>RENT</u>	<u>0</u>	<u>0</u>	<u>0</u>	<u>0</u>
<u>Southwest</u>	<u>33017</u>	<u>0.124801554</u>	<u>0</u>	<u>0</u>
<u>South East</u>	<u>63249</u>	<u>0.239076037</u>	<u>0</u>	<u>0</u>
<u>North East</u>	<u>60317</u>	<u>0.227993317</u>	<u>0</u>	<u>0</u>
<u>Midwest</u>	<u>47919</u>	<u>0.181129893</u>	<u>0</u>	<u>0</u>
<u>1</u>	<u>0</u>	<u>0</u>	<u>0</u>	<u>0</u>
<u>0</u>	<u>0</u>	<u>0</u>	<u>0</u>	<u>0</u>
<u>MORTGAGE</u>	<u>0</u>	<u>0</u>	<u>0</u>	<u>0</u>
<u>OWN</u>	<u>33517</u>	<u>0.126691513</u>	<u>0</u>	<u>0</u>
<u>percent_bc_gt_75</u>	<u>0</u>	<u>0</u>	<u>0</u>	<u>0</u>
<u>emp_length</u>	<u>0</u>	<u>0</u>	<u>0</u>	<u>0</u>
<u>loan_status</u>	<u>47310</u>	<u>0.178827923</u>	<u>0</u>	<u>0</u>
<u>purpose</u>	<u>0</u>	<u>0</u>	<u>0</u>	<u>0</u>
<u>initial_list_status</u>	<u>0</u>	<u>0</u>	<u>0</u>	<u>0</u>
<u>West</u>	<u>60054</u>	<u>0.226999199</u>	<u>0</u>	<u>0</u>

Appendix 3.2: Outliers ratio of fully paid/default

<u>Columns</u>	<u>Fully_paid</u>	<u>Default</u>	<u>frac</u>
<u>realized_return</u>	<u>0</u>	<u>7546</u>	<u>1</u>

<u>int_rate</u>	<u>1482</u>	<u>1276</u>	<u>0.462654097</u>
<u>grade</u>	<u>1411</u>	<u>1164</u>	<u>0.452038835</u>
<u>installment</u>	<u>1348</u>	<u>606</u>	<u>0.31013306</u>
<u>all_util</u>	<u>236</u>	<u>91</u>	<u>0.278287462</u>
<u>dti</u>	<u>534</u>	<u>195</u>	<u>0.267489712</u>
<u>revol_util</u>	<u>14</u>	<u>5</u>	<u>0.263157895</u>
<u>num_rev_tl_bal_gt_0</u>	<u>2798</u>	<u>874</u>	<u>0.238017429</u>
<u>pub_rec_bankruptcies</u>	<u>2212</u>	<u>629</u>	<u>0.221400915</u>
<u>num_actv_bc_tl</u>	<u>2982</u>	<u>830</u>	<u>0.217733473</u>
<u>bc_util</u>	<u>18</u>	<u>5</u>	<u>0.217391304</u>
<u>delinq_2yrs</u>	<u>3393</u>	<u>934</u>	<u>0.21585394</u>
<u>num_op_rev_tl</u>	<u>2886</u>	<u>757</u>	<u>0.207795773</u>
<u>pub_rec</u>	<u>2926</u>	<u>750</u>	<u>0.204026115</u>
<u>num_sats</u>	<u>2867</u>	<u>709</u>	<u>0.198266219</u>
<u>open_acc</u>	<u>2936</u>	<u>725</u>	<u>0.198033324</u>
<u>pct_tl_nvr_dlq</u>	<u>4156</u>	<u>1005</u>	<u>0.194729704</u>
<u>il_util</u>	<u>2668</u>	<u>631</u>	<u>0.191270082</u>
<u>cr_hist</u>	<u>2044</u>	<u>481</u>	<u>0.19049505</u>

<u>num_bc_sats</u>	<u>3899</u>	<u>886</u>	<u>0.185161964</u>
<u>num_rev_accts</u>	<u>2854</u>	<u>639</u>	<u>0.182937303</u>
<u>num_bc_tl</u>	<u>3041</u>	<u>661</u>	<u>0.178552134</u>
<u>total_bal_il</u>	<u>4242</u>	<u>800</u>	<u>0.158667196</u>
<u>tot_coll_amt</u>	<u>2069</u>	<u>386</u>	<u>0.157230143</u>
<u>total_il_high_credit_limit</u>	<u>3972</u>	<u>683</u>	<u>0.146723953</u>
<u>total_bal_ex_mort</u>	<u>4161</u>	<u>715</u>	<u>0.146636587</u>
<u>loan_amnt</u>	<u>1349</u>	<u>206</u>	<u>0.132475884</u>
<u>max_bal_bc</u>	<u>3651</u>	<u>474</u>	<u>0.114909091</u>
<u>annual_inc</u>	<u>1669</u>	<u>212</u>	<u>0.112706007</u>
<u>revol_bal</u>	<u>3238</u>	<u>410</u>	<u>0.112390351</u>
<u>mort_acc</u>	<u>4228</u>	<u>533</u>	<u>0.111951271</u>
<u>tot_cur_bal</u>	<u>3648</u>	<u>420</u>	<u>0.103244838</u>
<u>tot_hi_cred_lim</u>	<u>3339</u>	<u>370</u>	<u>0.099757347</u>
<u>total_rev_hi_lim</u>	<u>3798</u>	<u>420</u>	<u>0.099573257</u>
<u>avg_cur_bal</u>	<u>4317</u>	<u>448</u>	<u>0.094018888</u>
<u>total_bc_limit</u>	<u>4673</u>	<u>474</u>	<u>0.092092481</u>
<u>bc_open_to_buy</u>	<u>5324</u>	<u>402</u>	<u>0.070206078</u>

<u>fico_range_low</u>	<u>3887</u>	<u>203</u>	<u>0.049633252</u>
<u>fico_range_high</u>	<u>3887</u>	<u>203</u>	<u>0.049633252</u>

Appendix 4 -Modeling

4.1 Approach A – Regression model on realized return

MODEL	R Squared	RMSE
Baseline model linear regression only by grade	0.75327	0.04130
Linear Regression	0.01189	0.08282
KNN	0.18461	0.09068
Gradient Boosting Regressor	0.02078	0.08244
Linear Regression - scaled	0.01188	0.08282
KNN - scaled	0.16875	0.09007
Gradient Boosting Regressor - scaled	0.02076	0.08245

4.2 Approach B – classification model on loan status

Baseline model - predicting by grade: confusion matrix

	Fully Paid	Charged Off
Fully Paid	43424	26
Charged Off	9455	7

Model recall: 0.00074

Model precision score: 0.21212

Model roc_auc_score: 0.50007

Model accuracy score: 0.82082

Logistic regression: confusion matrix

	Fully Paid	Charged Off
Fully Paid	43432	18
Charged Off	9460	2

Model recall: 0.00021

Model precision score: 0.1000

Model roc_auc_score: 0.49990

Model accuracy score: 0.82087

KNN: confusion matrix

	Fully Paid	Charged Off
Fully Paid	41299	2151
Charged Off	8820	642

Model recall: 0.06785

Model precision score: 0.22986

Model roc_auc_score: 0.50917

Model accuracy score: 0.79266

Random forest: confusion matrix

	Fully Paid	Charged Off
Fully Paid	43216	234
Charged Off	9222	240

Model recall: 0.02536

Model precision score: 0.50633

Model roc_auc_score: 0.50999

Model accuracy score: 0.82129

Logistic regression scaled: confusion matrix

	Fully Paid	Charged Off
Fully Paid	43103	347
Charged Off	9116	346

Model recall: 0.03657

Model precision score: 0.49928

Model roc_auc_score: 0.51429

Model accuracy score: 0.82116

KNN scaled: confusion matrix

	Fully Paid	Charged Off
Fully Paid	41137	2313
Charged Off	8507	955

Model recall: 0.10093

Model precision score: 0.29223

Model roc_auc_score: 0.52385

Model accuracy score: 0.79551

Random forest scaled: confusion matrix

	Fully Paid	Charged Off
Fully Paid	43201	294
Charged Off	9179	238

Model recall: 0.02527

Model precision score: 0.44736

Model roc_auc_score: 0.50999

Model accuracy score: 0.82096