

## PREDICTING LENDERS USING MACHINE LEARNING

This project explores publicly available data from [LendingClub.com \(www.lendingclub.com\)](http://www.lendingclub.com). Lending Club connects people who need money (borrowers) with people who have money (investors). I created a model to predict if an investor will want to connect with a borrow based on whether or not the borrower paid back their loan in full. Full data can be viewed [here \(https://www.lendingclub.com/info/download-data.action\)](https://www.lendingclub.com/info/download-data.action).

Here are what the columns represent:

- credit.policy: 1 if the customer meets the credit underwriting criteria of LendingClub.com, and 0 otherwise.
- purpose: The purpose of the loan (takes values "credit\_card", "debt\_consolidation", "educational", "major\_purchase", "small\_business", and "all\_other").
- int.rate: The interest rate of the loan, as a proportion (a rate of 11% would be stored as 0.11). Borrowers judged by LendingClub.com to be more risky are assigned higher interest rates.
- installment: The monthly installments owed by the borrower if the loan is funded.
- log.annual.inc: The natural log of the self-reported annual income of the borrower.
- dti: The debt-to-income ratio of the borrower (amount of debt divided by annual income).
- fico: The FICO credit score of the borrower.
- days.with.cr.line: The number of days the borrower has had a credit line.
- revol.bal: The borrower's revolving balance (amount unpaid at the end of the credit card billing cycle).
- revol.util: The borrower's revolving line utilization rate (the amount of the credit line used relative to total credit available).
- inq.last.6mths: The borrower's number of inquiries by creditors in the last 6 months.
- delinq.2yrs: The number of times the borrower had been 30+ days past due on a payment in the past 2 years.
- pub.rec: The borrower's number of derogatory public records (bankruptcy filings, tax liens, or judgments).

A Random Forest Classification model will be used and will highlight my skills in:

- Machine learning algorithms
- Creating data visualizations
- Statistical modeling

-Sarah Aristil, Marketing Analyst Applicant saraharistil@outlook.com 305-924-4143

```
In [1]: import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
%matplotlib inline
```

```
In [2]: loans = pd.read_csv('loan_data.csv')
```

### Exploring data

```
In [3]: loans.head()
```

Out[3]:

	credit.policy	purpose	int.rate	installment	log.annual.inc	dti	fico	days.with.cr.line	revol.bal	revol.util	inq.last.6mths	delinq.2yrs	p
0	1	debt_consolidation	0.1189	829.10	11.350407	19.48	737	5639.958333	28854	52.1	0	0	
1	1	credit_card	0.1071	228.22	11.082143	14.29	707	2760.000000	33623	76.7	0	0	
2	1	debt_consolidation	0.1357	366.86	10.373491	11.63	682	4710.000000	3511	25.6	1	0	
3	1	debt_consolidation	0.1008	162.34	11.350407	8.10	712	2699.958333	33667	73.2	1	0	
4	1	credit_card	0.1426	102.92	11.299732	14.97	667	4066.000000	4740	39.5	0	1	

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

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 9578 entries, 0 to 9577
Data columns (total 14 columns):
credit.policy      9578 non-null int64
purpose           9578 non-null object
int.rate          9578 non-null float64
installment       9578 non-null float64
log.annual.inc    9578 non-null float64
dti               9578 non-null float64
fico              9578 non-null int64
days.with.cr.line 9578 non-null float64
revol.bal         9578 non-null int64
revol.util        9578 non-null float64
inq.last.6mths    9578 non-null int64
delinq.2yrs       9578 non-null int64
pub.rec           9578 non-null int64
not.fully.paid    9578 non-null int64
dtypes: float64(6), int64(7), object(1)
memory usage: 1.0+ MB
```

In [5]: `loans.describe()`

Out[5]:

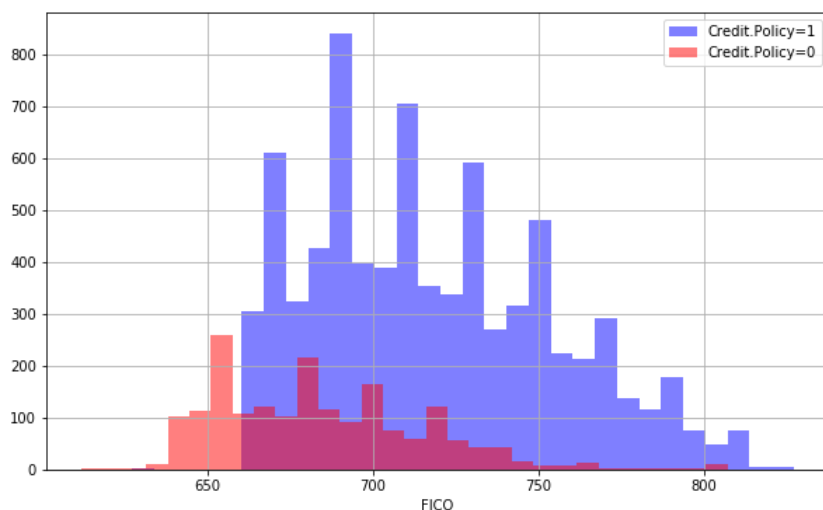
	credit.policy	int.rate	installment	log.annual.inc	dti	fico	days.with.cr.line	revol.bal	revol.util	inq.last.6mths
count	9578.000000	9578.000000	9578.000000	9578.000000	9578.000000	9578.000000	9578.000000	9.578000e+03	9578.000000	9578.000000
mean	0.804970	0.122640	319.089413	10.932117	12.606679	710.846314	4560.767197	1.691396e+04	46.799236	1.577469
std	0.396245	0.026847	207.071301	0.614813	6.883970	37.970537	2496.930377	3.375619e+04	29.014417	2.200245
min	0.000000	0.060000	15.670000	7.547502	0.000000	612.000000	178.958333	0.000000e+00	0.000000	0.000000
25%	1.000000	0.103900	163.770000	10.558414	7.212500	682.000000	2820.000000	3.187000e+03	22.600000	0.000000
50%	1.000000	0.122100	268.950000	10.928884	12.665000	707.000000	4139.958333	8.596000e+03	46.300000	1.000000
75%	1.000000	0.140700	432.762500	11.291293	17.950000	737.000000	5730.000000	1.824950e+04	70.900000	2.000000
max	1.000000	0.216400	940.140000	14.528354	29.960000	827.000000	17639.958330	1.207359e+06	119.000000	33.000000

## Data Analysis

In [6]: *#Creating a histogram of two FICO distributions, one for each credit.policy outcome*

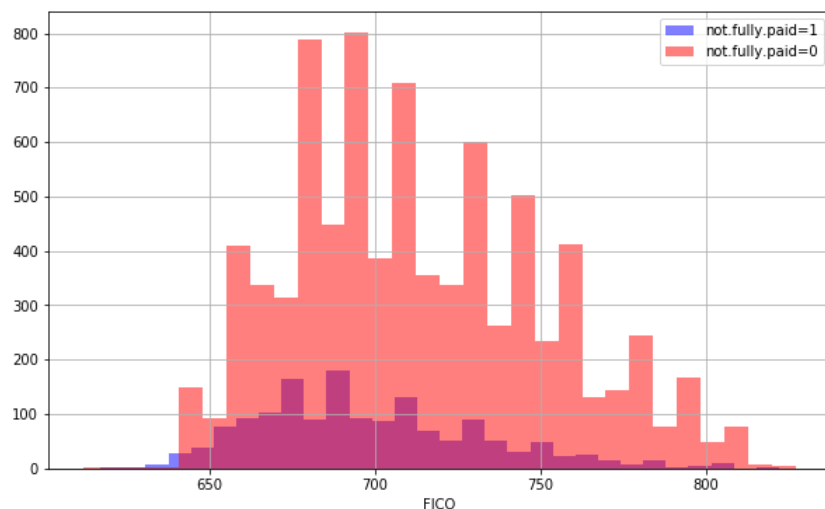
```
plt.figure(figsize=(10,6))
loans[loans['credit.policy']==1]['fico'].hist(alpha=0.5,color='blue',
                                              bins=30,label='Credit.Policy=1')
loans[loans['credit.policy']==0]['fico'].hist(alpha=0.5,color='red',
                                              bins=30,label='Credit.Policy=0')
plt.legend()
plt.xlabel('FICO')
```

Out[6]: Text(0.5, 0, 'FICO')



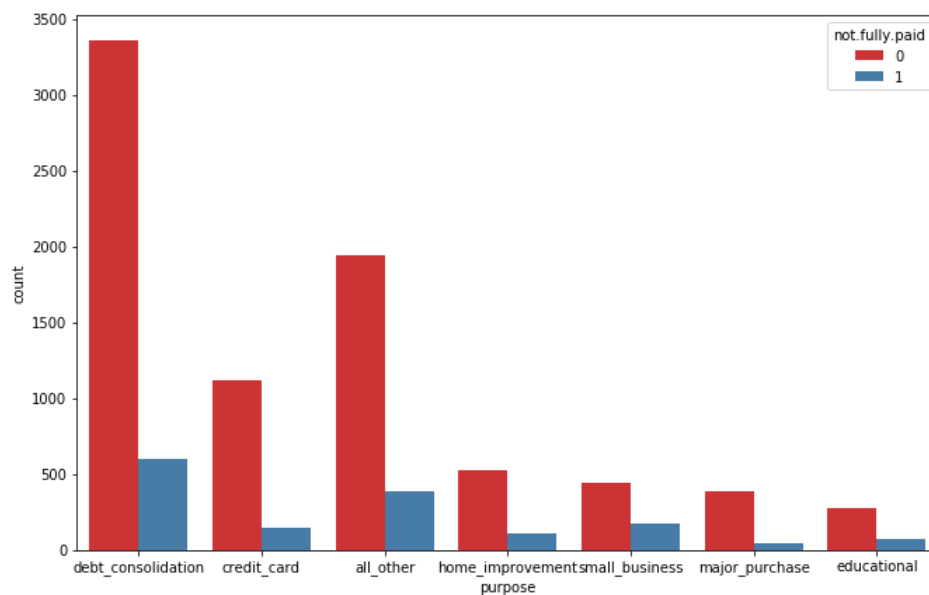
```
In [7]: # Creating a similar figure, using the not.fully.paid column
plt.figure(figsize=(10,6))
loans[loans['not.fully.paid']==1]['fico'].hist(alpha=0.5,color='blue',
                                                bins=30,label='not.fully.paid=1')
loans[loans['not.fully.paid']==0]['fico'].hist(alpha=0.5,color='red',
                                                bins=30,label='not.fully.paid=0')
plt.legend()
plt.xlabel('FICO')
```

Out[7]: Text(0.5, 0, 'FICO')



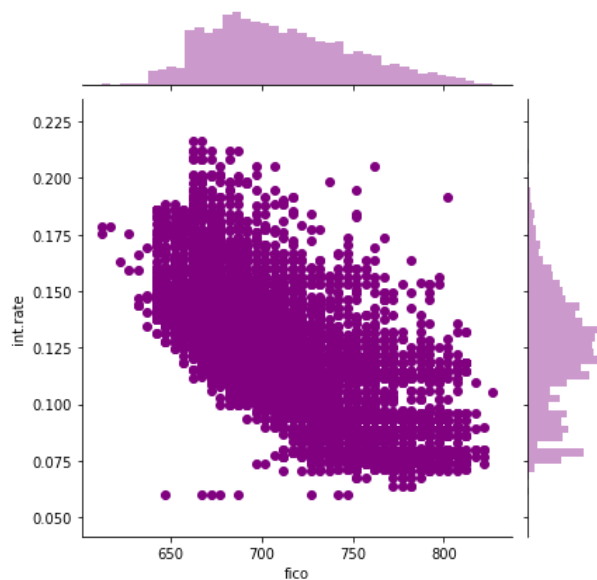
```
In [8]: #Showing the counts of loans by purpose, with the color hue defined by not.fully.paid
plt.figure(figsize=(11,7))
sns.countplot(x='purpose',hue='not.fully.paid',data=loans,palette='Set1')
```

Out[8]: <matplotlib.axes.\_subplots.AxesSubplot at 0x1144258d0>



```
In [9]: #Viewing the trend between FICO score and interest rate via jointplot
sns.jointplot(x='fico',y='int.rate',data=loans,color='purple')
```

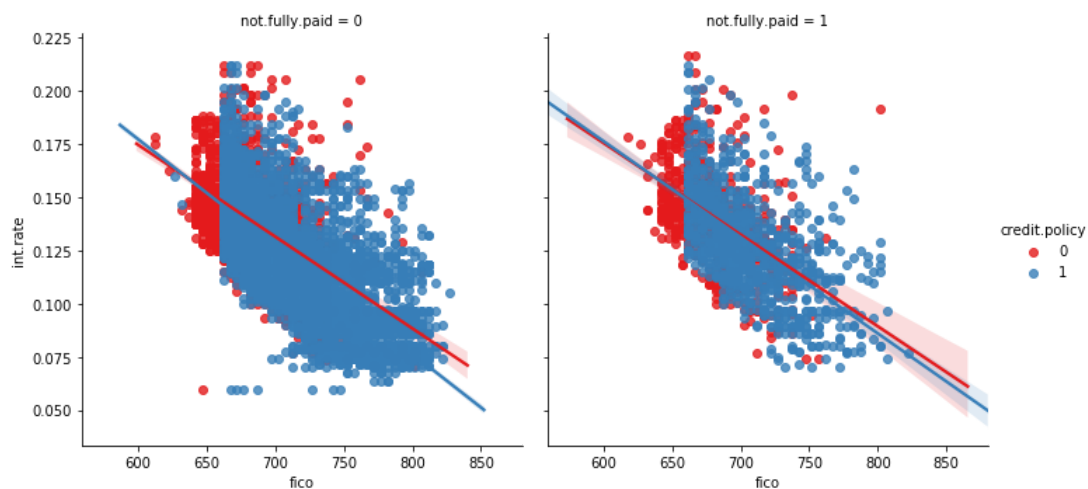
```
Out[9]: <seaborn.axisgrid.JointGrid at 0x117470898>
```



```
In [10]: #Creating lmplots to see if the trend differs between not.fully.paid and credit.policy
plt.figure(figsize=(11,7))
sns.lmplot(y='int.rate',x='fico',data=loans,hue='credit.policy',
          col='not.fully.paid',palette='Set1')
```

```
Out[10]: <seaborn.axisgrid.FacetGrid at 0x1175fab0>
```

<Figure size 792x504 with 0 Axes>



## Random Forest Classification Model

```
In [11]: loans.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 9578 entries, 0 to 9577
Data columns (total 14 columns):
credit.policy      9578 non-null int64
purpose           9578 non-null object
int.rate          9578 non-null float64
installment       9578 non-null float64
log.annual.inc    9578 non-null float64
dti               9578 non-null float64
fico              9578 non-null int64
days.with.cr.line 9578 non-null float64
revol.bal         9578 non-null int64
revol.util        9578 non-null float64
inq.last.6mths    9578 non-null int64
delinq.2yrs       9578 non-null int64
pub.rec           9578 non-null int64
not.fully.paid    9578 non-null int64
dtypes: float64(6), int64(7), object(1)
memory usage: 1.0+ MB
```

The **purpose** column is categorical. This will be transformed using dummy variables so sklearn will be able to understand those data.

```
In [12]: cat_feats = ['purpose']
```

```
In [13]: final_data = pd.get_dummies(loans,columns=cat_feats,drop_first=True)
```

```
In [14]: final_data.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 9578 entries, 0 to 9577
Data columns (total 19 columns):
credit.policy      9578 non-null int64
int.rate          9578 non-null float64
installment       9578 non-null float64
log.annual.inc    9578 non-null float64
dti               9578 non-null float64
fico              9578 non-null int64
days.with.cr.line 9578 non-null float64
revol.bal         9578 non-null int64
revol.util        9578 non-null float64
inq.last.6mths    9578 non-null int64
delinq.2yrs       9578 non-null int64
pub.rec           9578 non-null int64
not.fully.paid    9578 non-null int64
purpose_credit_card 9578 non-null uint8
purpose_debt_consolidation 9578 non-null uint8
purpose_educational 9578 non-null uint8
purpose_home_improvement 9578 non-null uint8
purpose_major_purchase 9578 non-null uint8
purpose_small_business 9578 non-null uint8
dtypes: float64(6), int64(7), uint8(6)
memory usage: 1.0 MB
```

### Creating Train Test Split for Model

```
In [15]: from sklearn.model_selection import train_test_split
```

```
In [16]: X = final_data.drop('not.fully.paid',axis=1)
y = final_data['not.fully.paid']
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.30, random_state=101)
```

```
In [17]: from sklearn.tree import DecisionTreeClassifier
```

```
In [18]: dtree = DecisionTreeClassifier()
```

```
In [19]: dtree.fit(X_train,y_train)
```

```
Out[19]: DecisionTreeClassifier(class_weight=None, criterion='gini', max_depth=None,
                                max_features=None, max_leaf_nodes=None,
                                min_impurity_decrease=0.0, min_impurity_split=None,
                                min_samples_leaf=1, min_samples_split=2,
                                min_weight_fraction_leaf=0.0, presort=False, random_state=None,
                                splitter='best')
```

### Predicting and Evaluating the Decision Tree

```
In [20]: #Creating a predictions from the test set, a classification report and a confusion matrix
         predictions = dtree.predict(X_test)
```

```
In [21]: from sklearn.metrics import classification_report,confusion_matrix
```

```
In [22]: print(classification_report(y_test,predictions))
```

	precision	recall	f1-score	support
0	0.85	0.82	0.84	2431
1	0.18	0.22	0.20	443
micro avg	0.73	0.73	0.73	2874
macro avg	0.52	0.52	0.52	2874
weighted avg	0.75	0.73	0.74	2874

```
In [23]: print(confusion_matrix(y_test,predictions))
```

```
[[2003  428]
 [ 346   97]]
```

### Training the Random Forest model

```
In [24]: from sklearn.ensemble import RandomForestClassifier
```

```
In [25]: rfc = RandomForestClassifier(n_estimators=600)
```

```
In [26]: rfc.fit(X_train,y_train)
```

```
Out[26]: RandomForestClassifier(bootstrap=True, class_weight=None, criterion='gini',
                                max_depth=None, max_features='auto', max_leaf_nodes=None,
                                min_impurity_decrease=0.0, min_impurity_split=None,
                                min_samples_leaf=1, min_samples_split=2,
                                min_weight_fraction_leaf=0.0, n_estimators=600, n_jobs=None,
                                oob_score=False, random_state=None, verbose=0,
                                warm_start=False)
```

### Predictions and Evaluation

```
In [27]: #Predict off the y_test values and the class of not.fully.paid for the X_test data.**
         predictions = rfc.predict(X_test)
```

```
In [28]: #Classification report
         from sklearn.metrics import classification_report,confusion_matrix
```

```
In [29]: print(classification_report(y_test,predictions))
```

	precision	recall	f1-score	support
0	0.85	1.00	0.92	2431
1	0.50	0.02	0.04	443
micro avg	0.85	0.85	0.85	2874
macro avg	0.67	0.51	0.48	2874
weighted avg	0.79	0.85	0.78	2874

```
In [30]: print(confusion_matrix(y_test, predictions))
```

```
[[ 2422    9]
 [  434    9]]
```

Looking at the classification report, neither did very well, more feature engineering is needed for a more accurate prediction model.