Prediction of LendingClub loan defaulters RIHAD VARIAWA 09-03-2019 The dataset used for this analysis can be found using the following <u>link</u> **About LendingClub:** LendingClub is a US peer-to-peer lending company, headquartered in San Francisco, California. It was the first peer-topeer lender to register its offerings as securities with the Securities and Exchange Commission (SEC), and to offer loan trading on a secondary market. Lending Club is the world's largest peer-to-peer lending platform. The company claims that \$15.98 billion in loans had been originated through its platform up to December 31, 2015. Lending Club enables borrowers to create unsecured personal loans between \$\$1,000 and \$40,000. The standard loan period is three years. Investors can search and browse the loan listings on Lending Club website and select loans that they want to invest in based on the information supplied about the borrower, amount of loan, loan grade, and loan purpose. Investors make money from interest. Lending Club makes money by charging borrowers an origination fee and investors a service fee. **About the Dataset** These files contain complete loan data for all loans issued through the 2007-2015, including the current loan status (Current, Late, Fully Paid, etc.) and latest payment information. The file containing loan data through the "present" contains complete loan data for all loans issued through the previous completed calendar quarter. Additional features include credit scores, number of finance inquiries, address including zip codes, and state, and collections among others. The file is a matrix of about 890 thousand observations and 75 variables. A data dictionary is provided in a separate file. k Purpose of this analysis We will go step by step for building a machine learning algorith for the prediction of loan defaulters based on certain variables present in the dataset. Our main goal is to correctly identifying defaulter's (True positives) so that lending club can decide whether a person is fit for sanctioning a loan or not in the future. **How Lending Club Works?** How Lending Club Works Borrowers apply for loans. Borrowers get funded. Borrowers repay automatically. Investors open an account. Investors build a portfolio. Investors earn & reinvest. **Packages** We will start by importing the packages that will be used throughout the analysis In [0]: # data manipulations and prepocessing import pandas as pd import numpy as np # data visualization import seaborn as sns import matplotlib.pyplot as plt from sklearn.model selection import train test split from sklearn.model_selection import cross_val_score from sklearn.model selection import GridSearchCV from sklearn.preprocessing import StandardScaler from sklearn.metrics import confusion matrix from sklearn.metrics import roc auc score from sklearn.metrics import roc curve from scipy.stats import boxcox In [0]: from google.colab import files files.upload() Choose Files | No file chosen Upload widget is only available when the cell has been executed in the current browser session. Please rerun this cell to enable. Saving loan.csv.zip to loan.csv.zip In [0]: !ls In [0]: !unzip loan.csv.zip In [0]: # read the dataset df = pd.read_csv("loan.csv", low_memory=False) df.head(3)**Data Cleaning** In [0]: # find the count and percentage of values that are missing in the dataframe. df_null = pd.DataFrame({'Count': df.isnull().sum(), 'Percent': 100*df.isnull().sum()/len(df)}) # print columns with null count more than 0 df null[df null['Count'] > 0] As you can see, there are a lot of columns which have huge chunk of data missing. These columns are not necessary for our analysis. The following part will drop any columns where 80% or more data is missing. This will help us clean the Dataset a little bit. In [0]: df1 = df.dropna(axis=1, thresh=int(0.80*len(df))) In [0]: df1.head(5) Now that we explored the whole dataframe easily, we will now select the columns that are necessary for our analysis. In [0]: df_LC = df1.filter(['loan_amnt','term','int_rate','installment','grade','sub_grade','emp_lengt h', 'home ownership', 'annual inc','verification status','purpose','dti','delinq 2yrs','loan stat us']) df LC.dtypes Finding the correlation between variables We will now look at the correlation structure between our variables that we selected above. This will tell us about any dependencies between different variables and help us reduce the dimensionality a little bit more In [0]: plt.figure(figsize=(20,20)) sns.set context("paper", font_scale=1) # find the correllation matrix and changing the categorical data to category for the plot. sns.heatmap(df LC.assign(grade=df LC.grade.astype('category').cat.codes, sub g=df_LC.sub_grade.astype('category').cat.codes, term=df LC.term.astype('category').cat.codes, emp l=df LC.emp length.astype('category').cat.codes, ver =df LC.verification status.astype('category').cat.codes, home=df LC.home ownership.astype('category').cat.codes, purp=df LC.purpose.astype('category').cat.codes).corr(), annot=True, cmap='bwr', vmin=-1, vmax=1, square=True, linewidths=0.5) It can be seen from the plot above that loan amount and installment have a very high correlation amongst each other (0.94). This is intuitive since a person who takes a large sum of loan would require extra time to repay it back. Also, interest rate, sub grade and grade have a very high correlation between them. This is obvious since interest rate is decided by grades once the grades are decided, a subgrade is assigned to that loan (leading to high correlation). Let's drop the three categories alongwith term and verification status(since it doesn't provide any valuable info) for further analysis. In [0]: | df LC.drop(['installment', 'grade', 'sub grade', 'verification status', 'term'] , axis=1, inplace = True) In [0]: # print the count and null values in the dataframe dflc null = pd.DataFrame({'Count': df LC.isnull().sum(), 'Percent': 100*df LC.isnull().sum()/le n(df_LC)}) dflc null[dflc null['Count'] > 0] In [0]: # dropping the null rows since we have sufficient amount of data and there is no need to fill t he null values. df LC.dropna(axis=0) In [0]: # print unique statuses in the loan status column (dependent variable) df_LC['loan_status'].unique() Distribution of the loan status values Let us now see how the values in the status column are distributed. We will plot an histogram of values against count of times the status appears on the dataframe In [0]: m =df LC['loan status'].value counts() m = m.to frame()m.reset_index(inplace=True) m.columns = ['Loan Status', 'Count'] plt.subplots(figsize=(20,8)) sns.barplot(y='Count', x='Loan Status', data=m) plt.xlabel("Length") plt.ylabel("Count") plt.title("Distribution of Loan Status in our Dataset") plt.show() As you can see, we have a lot of loans which are current with fair amount of fully paid loans. other categories (including) default have a really low number. This means the data is imbalanced and we might need to do something about this later in the analysis. For now we will drop all the columns except 'Fully Paid', 'Default' and 'Charged off'. We will also merge 'Charged off' and 'Default' together meaning that anyone who fell into this category defaulted their loan. The following two parts tries to implement this. In [0]: df LC = df LC[df LC.loan status != 'Current'] df LC = df LC[df LC.loan status != 'In Grace Period'] df LC = df LC[df LC.loan status != 'Late (16-30 days)'] df LC = df LC[df LC.loan status != 'Late (31-120 days)'] df_LC = df_LC[df_LC.loan_status != 'Does not meet the credit policy. Status:Fully Paid'] df LC = df LC[df LC.loan status != 'Does not meet the credit policy. Status:Charged Off'] df LC = df LC[df LC.loan status != 'Issued'] In [0]: | df_LC['loan_status'] = df_LC['loan_status'].replace({'Charged Off':'Default'}) df_LC['loan_status'].value_counts() We will now encode the two categories listed above as 0 or 1 for our analysis. This will help us in predicting whether a person defaulted their loan or not. 0 means he deaulted and 1 means he paid off his loan. In [0]: df LC.loan_status=df_LC.loan_status.astype('category').cat.codes df_LC.delinq_2yrs=df_LC.delinq_2yrs.astype('category').cat.codes df LC.head() df LC['loan status'].value counts() #df_LC = pd.get_dummies(df_LC, drop_first=True) In [0]: df LC.dtypes **Transformation** Before training the data, we would first transform the data to account for any skewness in the variable distribution. Various transformation techniques ranging from log transform to power transformation are available. For our analysis, we'll be using Box-cox transformation. It is used to modify the distributional shape of a set of data to be more normally distributed so that tests and confidence limits that require normality can be appropriately used. In [0]: numerical = df_LC.columns[df_LC.dtypes == 'float64'] for i in numerical: **if** df_LC[i].min() > 0: transformed, lamb = boxcox(df_LC.loc[df[i].notnull(), i]) **if** np.abs(1 - lamb) > 0.02: df LC.loc[df[i].notnull(), i] = transformed One Hot Encoding Since we have some categorical variables for the analysis and the machne learning algorithms doesn't take categorical and string variables directly, we have to creat dummy variables for them. We can either encode them using label encoder available for python, but it would be wrong in our analysis since a lot of these variables have multiple categories. Just using weights can cause discrepencies in the algorithm. Instead, we will one hot encode these so that we have a 1 wherever that category turns up and 0 otherwise. This will also create seperate columns for each level of category. Also, we'll be dropping one of the categories so that we have N-1 columns instead of N. In [0]: df_LC = pd.get_dummies(df_LC, drop_first=True) Now splitting the data using scikitlearn's train_test_split and using 60% data for training and 40% for testing. In [0]: traindata, testdata = train_test_split(df_LC, stratify=df_LC['loan_status'],test_size=.4, rando testdata.reset index(drop=True, inplace=True) traindata.reset index(drop=True, inplace=True) We'll now scale the data so that each column has a mean of zero and unit standard deviation. Xunb (unbalanced set) and yunb are the independent and target variable. In [0]: | sc = StandardScaler() Xunb = traindata.drop('loan_status', axis=1) yunb = traindata['loan_status'] numerical = Xunb.columns[(Xunb.dtypes == 'float64') | (Xunb.dtypes == 'int64')].tolist() Xunb[numerical] = sc.fit transform(Xunb[numerical]) In [0]: # preview the shape of training data yunb.shape **Model Selection** We are now ready to build some models. The following would be our approach for building and selecting the best model: 1. Build a model on the imbalance dataset we got from data cleaning. 2. Balance the dataset by using equal amount of default and 'fully paid' loans. Trying the Unbalanced Dataset Let's first try the unbalanced dataset. The function below computes the receiver operating characteristic (ROC) curves for each of the models. This function will be called later in the analysis. In [0]: # roc curve to find a good model that optimizes the trade off between # the False Positive Rate (FPR) and True Positive Rate (TPR) def createROC(models, X, y, Xte, yte): false p, true p = [], [] # false postives and true positives for i in models.keys(): # dict of models models[i].fit(X, y) fp, tp, threshold = roc curve(yte, models[i].predict proba(Xte)[:,1]) # roc curve funct ion true_p.append(tp) false p.append(fp) return true_p, false_p # returning the true postive and false positive Let's try some models on the train dataset With 3 fold cross validation. We are going to use the following 4 machine learning algorithms: 1. Linear Discriminant Analysis 2. Multinomial Naive Bayes 3. Random Forest (tree based model) 4. Logistic Regression from sklearn.discriminant analysis import LinearDiscriminantAnalysis In [0]: from sklearn.ensemble import RandomForestClassifier from sklearn.linear model import LogisticRegression from sklearn.naive bayes import MultinomialNB models = {'LDA': LinearDiscriminantAnalysis(), 'MNB': MultinomialNB(), 'RF': RandomForestClassifier(n estimators=100), 'LR': LogisticRegression(C=1)} unbalset = {} for i in models.keys(): scores = cross_val_score(models[i], Xunb - np.min(Xunb) + 1, yunb, cv=3) unbalset[i] = scores print(i, scores, np.mean(scores)) Looks like Logistic regression provides the best estimate and almost all of the models giving the same results. Because of the issue of collinearity in LDA, we are going to ignore that. Now creating the test set for the analysis and scaling it. In [0]: | Xte = testdata.drop('loan status', axis=1) yte = testdata['loan status'] numerical = Xte.columns[(Xte.dtypes == 'float64') | (Xte.dtypes == 'int64')].tolist() Xte[numerical] = sc.fit_transform(Xte[numerical]) Computing the ROC curves for the models and finding the true positive and false positives. In [0]: tp unbalset, fp unbalset = createROC(models, Xunb - np.min(Xunb) + 1, yunb, Xte - np.min(Xte) + 1, yte) Fitting LR to the test set. In [0]: model = LogisticRegression(C=1) model.fit(Xunb, yunb) predict = model.predict(Xte) # prediction of Xte which can be used to test against yte (testdat a values or true values of y) In [0]: | m = yte.to_frame() m['loan status'].value counts() We will now plot the cross-validation scores, ROC curves and confusion matrix of random forest model. X axis is the true value and Y axis is the predicted value. In [0]: fig, axes = plt.subplots(nrows=1, ncols=3, figsize=(18,5)) ax = pd.DataFrame(unbalset).boxplot(widths=(0.9,0.9,0.9,0.9), grid=False, vert=False, ax=axes[0]) ax.set_ylabel('Classifier') ax.set xlabel('Cross-Validation Score') for i in range(0, len(tp unbalset)): axes[1].plot(fp unbalset[i], tp unbalset[i], lw=1) axes[1].plot([0, 1], [0, 1], '--k', lw=1) axes[1].legend(models.keys()) axes[1].set ylabel('True Positive Rate') axes[1].set_xlabel('False Positive Rate') $axes[1].set_xlim(0,1)$ $axes[1].set_ylim(0,1)$ cm = confusion_matrix(yte, predict).T cm = cm.astype('float')/cm.sum(axis=0) ax = sns.heatmap(cm, annot=True, cmap='Blues', ax=axes[2]); ax.set_xlabel('True Value') ax.set_ylabel('Predicted Value') ax.axis('equal') The cross-validation scores and ROC curves suggest the Logistic Regression is the best model, though the MNB and linear discriminant analysis models are pretty close behind. If we look at the confusion matrix, though, we see a big problem. The model can predict who are going to pay off the loan with a good accuracy of 99% but cannot predict who are going to default. The true positive rate of default (0 predicting 0) is almost 0. Since our main goal is to predict defaulter's, we have to do something about this. The reason this is happening could be because of high imbalance in our dataset and the algorithm is putting everything into 1. We have to chose a new prediction threshold according to the sensitivity and specificity of the model. This will create some balance in predicting the binary outcome. Let's look at the plots below. In [0]: fp, tp, threshold = roc_curve(yte, model.predict_proba(Xte)[:,1]) #getting false and true posit ive from test set fig, ax = plt.subplots(nrows=1, ncols=2, figsize=(16,6)) In [0]: ax[0].plot(threshold, tp + (1 - fp))ax[0].set_xlabel('Threshold') ax[0].set_ylabel('Sensitivity + Specificity') ax[1].plot(threshold, tp, label="tp") ax[1].plot(threshold, 1 - fp, label="1 - fp") ax[1].legend() ax[1].set xlabel('Threshold') ax[1].set ylabel('True Positive & False Positive Rates') In [0]: # finding the optimal threshold for the model function = tp + (1 - fp)index = np.argmax(function) optimal_threshold = threshold[np.argmax(function)] print('optimal threshold:', optimal threshold) The optimal threshold above is where the two graphs meet. 1. Sensitivity (also called the true positive rate, the recall, or probability of detection in some fields) measures the proportion of actual positives that are correctly identified as such (e.g., the percentage of sick people who are correctly identified as having the condition). 2. Specificity (also called the true negative rate) measures the proportion of actual negatives that are correctly identified as such (e.g., the percentage of healthy people who are correctly identified as not having the condition) Now using this threshold for the model: In [0]: predict = model.predict proba(Xte)[:,1] predict = np.where(predict >= optimal threshold, 1, 0) fig, axes = plt.subplots(figsize=(15,6)) cm = confusion matrix(yte, predict).T cm = cm.astype('float')/cm.sum(axis=0) ax = sns.heatmap(cm, annot=True, cmap='Blues'); ax.set xlabel('True Value') ax.set_ylabel('Predicted Value') ax.axis('equal') That's great! The optimum threshold for the classifier have increased out models prediction power of Default (0). Even now the model doesn't provide a lot of prediction power and we have to train the model again using a different algorithm with some tweaks. Part2: Balancing the training dataset and creating a new model Now we will try to use a balanced dataset with equal amount of zeroes and 1's. The following part does the same. y_default = traindata[traindata['loan status'] == 0] In [0]: n_paid = traindata[traindata['loan_status'] == 1].sample(n=len(y_default), random_state=17) # c hosing equal amount of 1's # creating a new dataframe for balanced set data = y default.append(n paid) # creating the independent and dependent array Xbal = data.drop('loan status', axis=1) ybal = data['loan status'] In [0]: # scaling it again numerical = Xbal.columns[(Xbal.dtypes == 'float64') | (Xbal.dtypes == 'int64')].tolist() Xbal[numerical] = sc.fit transform(Xbal[numerical]) Training the model on the balanced set In [0]: from sklearn.discriminant_analysis import LinearDiscriminantAnalysis from sklearn.ensemble import RandomForestClassifier from sklearn.linear_model import LogisticRegression from sklearn.naive_bayes import MultinomialNB models = {'LDA': LinearDiscriminantAnalysis(), 'MNB': MultinomialNB(), 'RF': RandomForestClassifier(n estimators=100), 'LR': LogisticRegression(C=1)} balset = {} for i in models.keys(): scores = cross_val_score(models[i], Xbal - np.min(Xbal) + 1, ybal, scoring='roc auc', cv=3) balset[i] = scores print(i, scores, np.mean(scores)) Even though we almost got the same result as before, This time we are going to select Random Forst method and will try to find the optimal number of trees using the gridsearchcv and try to make the predition based on this and lets see if there is any improvements in predicting 0's In [0]: | model = RandomForestClassifier(n_estimators=100) model.fit(Xbal, ybal) predict = model.predict(Xte) In [0]: predict = model.predict(Xte) fig, axes = plt.subplots(figsize=(8,6)) cm = confusion_matrix(yte, predict).T cm = cm.astype('float')/cm.sum(axis=0) ax = sns.heatmap(cm, annot=True, cmap='Blues'); ax.set xlabel('True Label') ax.set ylabel('Predicted Label') ax.axis('equal') That's a significant improvement over the last model that we built using Logistic regression. Let's find the optimum number of estimators for this model and use that for prediction. This time we are going to use 5 fold cross validation. In [0]: params = {'n estimators': [50, 100, 200, 400, 600, 800]} grid search = GridSearchCV(RandomForestClassifier(), param grid=params, scoring='accuracy', cv=5, n_jobs=-1) grid search.fit(Xbal, ybal) print(grid search.best params) print(grid search.best score) The best model has 600 trees In [0]: #r = pd.DataFrame()#r['x'] = [i for i in params.values()][0]#r['y'] = [i[1] for i in grid search.cv results] #ax = r.plot(x='x', y='y', legend=False, linestyle='-', marker='o', figsize=(8,6)) #ax.set xlabel('n estimators') #ax.set ylabel('5-Fold Cross-Validation Score') As we can see, the score increases as we increase the number of estimators till 600 and then falls for 800 number of estimators. We will use this to estimate to fit the model. In [0]: grid search.best estimator .fit(Xbal, ybal) predict = model.predict(Xte) fig, axes = plt.subplots(figsize=(15,9)) In [0]: cm = confusion matrix(yte, predict).T cm = cm.astype('float')/cm.sum(axis=0) ax = sns.heatmap(cm, annot=True, cmap='Blues'); ax.set xlabel('True Label') ax.set ylabel('Predicted Label') ax.axis('equal') Interestingly, this gives us the same output as the previous model. Even now we have a good accuracy of 71% predicting defaluter's as defaulter's. Since random forest is based on decision trees, we can also plot the variable importance. Variable importance tells us which variable had highest importance when predicting an outcome. In [0]: r = pd.DataFrame(columns=['Feature', 'Importance']) ncomp = 15r['Feature'] = feat_labels = Xbal.columns r['Importance'] = model.feature importances r.set index(r['Feature'], inplace=True) ax = r.sort values('Importance', ascending=False)[:ncomp].plot.bar(width=0.9, legend=False, fig size=(15, 8))ax.set_ylabel('Relative Importance') As you can see, interest rate followed by debt to income ratio, annual income, and loan amount are the most important features in predicting the defaulter's. Lending Club might want to use this as the metric for identifying people defaulting on their loans. Conclusion We have successfully built an machine learning algorithm to predict the people who might default on their loans. This can be further used by LendingClub for their analysis. Also, we might want to look on other techniques or variables to improve the prediction power of the algorthm. One of the drawbacks is just the limited number of people who defaulted on their loan in the 8 years of data (2007-2015) present on the dataset. We can use an updated dataframe which consist next 3 years values (2015-2018) and see how many of the current loans were paid off or defaulted or even charged off. Then these new data points can be used for predicting them or even used to train the model again to improve its accuracy. Since we had a lot of categorical data, we cannot apply PCA for dimensionality reduction. Because of this, we can try some different type of variable selection method like 'MULTIPLE CORRESPONDENCE ANALYSIS' to reduce the dimensionality and select the most important variables from the columns. Since the algorithm puts around 47% of non-defaulters in the default class, we might want to look further into this issue to make the model more robust.