Imports In [3]: # DS Basics import numpy as np import pandas as pd import scipy from math import sqrt import matplotlib.pyplot as plt import seaborn as sns # estimators from sklearn.ensemble import RandomForestRegressor from sklearn.linear_model import LinearRegression from sklearn.svm import SVR from sklearn import linear_model # model metrics from sklearn.metrics import mean_squared_error from sklearn.metrics import r2 score from sklearn.model_selection import cross val score # cross validation from sklearn.model_selection import train test split # classification - pulling in as another way to look at the data from sklearn.metrics import accuracy_score, classification_report from sklearn.ensemble import RandomForestClassifier, GradientBoostingClassifier from sklearn.tree import DecisionTreeClassifier from sklearn.model selection import KFold from sklearn.tree import plot tree # Grahpviz from six import StringIO from IPython.display import Image from sklearn.tree import export_graphviz import pydotplus import graphviz from pydotplus import graph_from_dot_data Pull in the data In [4]: df = pd.read_csv('creditOne_cleaned_formatted.csv', index_col=0) df.rename(columns={'PAY_0':'PAY_1'}, inplace=True) df = pd.get_dummies(df) df.head() Out[4]: ID LIMIT_BAL MARRIAGE AGE PAY_1 PAY_2 PAY_3 PAY_4 PAY_5 PAY_6 ... PAY_AMT5 PAY_AMT6 SEX_female SEX_male 0 1 20000 0 0 2000 1 2 120000 2 26 -1 2 0 0 0 2 ... 0 1 0 90000 0 1000 5000 **2** 3 2 34 0 0 3 50000 1 37 0 0 0 1069 1000 1 0 5 50000 -1 0 0 ... 689 679 0 57 5 rows × 30 columns In [5]: df.info() <class 'pandas.core.frame.DataFrame'> Int64Index: 30000 entries, 0 to 29999 Data columns (total 30 columns): # Non-Null Count Dtype Column 0 ID 30000 non-null int64 1 LIMIT BAL 30000 non-null int64 MARRIAGE 30000 non-null int64 30000 non-null int64 PAY 1 30000 non-null int64 30000 non-null int64 5 PAY 2 PAY 3 30000 non-null int64 30000 non-null int64 7 PAY 4 8 PAY_5 30000 non-null int64 30000 non-null int64 9 PAY 6 30000 non-null int64 10 BILL AMT1 11 BILL AMT2 30000 non-null int64 12 BILL AMT3 30000 non-null int64 13 BILL_AMT4 30000 non-null int64 14 BILL AMT5 30000 non-null int64 15 BILL AMT6 30000 non-null int64 16 PAY AMT1 30000 non-null int64 30000 non-null int64 17 PAY AMT2 18 PAY AMT3 30000 non-null int64 19 PAY AMT4 30000 non-null int64 20 PAY AMT5 30000 non-null int64 21 PAY AMT6 30000 non-null int64 22 SEX_female 30000 non-null uint8 23 SEX male 30000 non-null uint8 24 EDUCATION_graduate school 30000 non-null uint8 25 EDUCATION_high school 30000 non-null uint8 26 EDUCATION_other 30000 non-null uint8
27 EDUCATION_university 30000 non-null uint8
28 DEFAULT_default 30000 non-null uint8
29 DEFAULT_not default 30000 non-null uint8 26 EDUCATION_other dtypes: int64(22), uint8(8) memory usage: 5.5 MB Modeling Business problem: "The bottom line is they need a much better way to understand how much credit to allow someone to use or, at the very least, if someone should be approved or not." In [6]: #features X = df.iloc[:, 3:24]print('Summary of feature sample') X.head(2)Summary of feature sample Out[6]: AGE PAY_1 PAY_2 PAY_3 PAY_4 PAY_5 PAY_6 BILL_AMT1 BILL_AMT2 BILL_AMT3 ... BILL_AMT5 BILL_AMT6 PAY_AMT1 3913 24 3102 689 0 0 0 2 2682 1725 3455 3261 0 26 -1 2682 ... 2 rows × 21 columns In [7]: # dependent variable y=df['LIMIT BAL'] In [8]: model = LinearRegression(n jobs=10) print(cross val score(model, X, y, cv=3)) [0.31050506 0.33093835 0.29553982] **Modeling Format and Model Selection** Import regression algorithms into a list so they can be run serially In [10]: algosClass = []algosClass.append(('Random Forest Regressor', RandomForestRegressor())) algosClass.append(('Linear Regression', LinearRegression())) algosClass.append(('Support Vector Regression', SVR())) In [11]: #regression #results = []#names = []#for name, model in algosClass: result = cross_val_score(model, X,y, cv=3, scoring='r2') names.append(name) # results.append(result) #for i in range(len(names)): print(names[i],results[i].mean()) X = df.iloc[:,3:24]Random Forest Regressor 0.4478359282537678 Linear Regression 0.31232774134169866 Support Vector Regression -0.0503310512284337 Let's try with just those columns that correlated well In [12]: | #X= df.iloc[:,4:10] #results = []#names = []#for name, model in algosClass: # result = cross val score(model, X,y, cv=3, scoring='r2') # names.append(name) # results.append(result) #for i in range(len(names)): print(names[i],results[i].mean()) X = df.iloc[:,4:10]Random Forest Regressor 0.10574276820435224 Linear Regression 0.101961724682885 Support Vector Regression -0.050732632640690355 definitely much worse. It's a good idea to keep more columns in, even if it takes longer. Let's try with them all: In [13]: #X = df.iloc[:,2:30]#results = []#names = []#for name, model in algosClass: # result = cross val score(model, X,y, cv=3, scoring='r2') names.append(name) # results.append(result) #for i in range(len(names)): print(names[i],results[i].mean()) X = df.iloc[:,2:30]Random Forest Regressor 0.467454580873716 Linear Regression 0.35819894266104485 Support Vector Regression -0.05038589493310345 In [14]: X = df.iloc[:,2:30]X.head(2)Out[14]: MARRIAGE AGE PAY_1 PAY_2 PAY_3 PAY_4 PAY_5 PAY_6 BILL_AMT1 BILL_AMT2 ... PAY_AMT5 PAY_AMT6 SEX_female SI 0 24 -2 -2 3913 3102 ... 0 2 0 0 0 2 2682 2000 26 -1 1725 ... 1 2 rows × 28 columns Choose the algorithm In [31]: algo = RandomForestRegressor(n_estimators=100,min_samples_leaf=20) #algo = LinearRegression() In [32]: # ToDo: split data into train/test sections X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.25, random_state=42) print(X_train.shape, X_test.shape) print(y_train.shape, y_test.shape) (22500, 28) (7500, 28) (22500,) (7500,) In [33]: model = algo.fit(X_train,y_train) In [34]: predictions = model.predict(X_test) predRsquared = r2 score(y test, predictions) rmse = sqrt(mean_squared_error(y_test, predictions)) print('R Squared: %.3f' % predRsquared) print('RMSE: %.3f' % rmse) R Squared: 0.485 RMSE: 93578.034 Random Forest (untuned): R Squared: 0.476, RMSE: 94415.243 Random Forest (tuned): R Squared: 0.485, RMSE: 93578.034 Linear Regression: R Squared: 0.370, RMSE: 103500.020 These are pretty close to the cross-validation scores #plt.scatter(y test, predictions, alpha = 0.5) In [36]: plt.scatter(y_test[0:400], predictions[0:400], alpha = 0.5) plt.xlabel('Ground Truth') plt.ylabel('Predictions') plt.plot(y_test[0:400], y_test[0:400], c='red') plt.tight_layout() plt.savefig('model_fitting.png') plt.show(); 600000 500000 400000 300000 200000 100000 300000 400000 500000 100000 200000 600000 Ground Truth Predictions look to run a bit high on the low end of the x-range, and run under at the high end of the x-range maybe consider taking a subset of the data with no defaults? Maybe compare RMS for defaulters vs. non-defaulters? Interpreted from slack: try binning LIMIT_BAL and turning this into a classification model like C1T2 from Task1: The bottom line is they need a much better way to understand how much credit to allow someone to use or, at the very least, if someone should be approved or not Idea: use classification models to look at: · binned limit data default **Classification Modelling - LimitBal** In [142]: #make a copy of the dataframe to play with discretization df2 = df.copy(deep=True) In [143]: | df2['LIMIT BAL'].describe() 30000.000000 Out [143]: count 167484.322667 mean std 129747.661567 min 10000.000000 50000.000000 25% 50% 140000.000000 240000.000000 75% 1000000.000000 max Name: LIMIT_BAL, dtype: float64 In [144]: bins = 9 limit bin = pd.qcut(df2['LIMIT BAL'],bins, labels=range(1, bins+1)) In [145]: | df2['limit bin'] = limit bin df2['limit bin'].head() Out[145]: 0 1 1 4 2 4 3 2 Name: limit bin, dtype: category Categories (9, int64): [1 < 2 < 3 < 4 ... 6 < 7 < 8 < 9]In [146]: # use the limit bin as our dependent variable Yclass=df2['limit bin'] Xclass = df2.iloc[:,2:30]Xclass.head(2) Out[146]: MARRIAGE AGE PAY_1 PAY_2 PAY_3 PAY_4 PAY_5 PAY_6 BILL_AMT1 BILL_AMT2 ... PAY_AMT5 PAY_AMT6 SEX_female SI 0 1 24 3913 3102 ... 0 2 2682 1725 ... 2000 1 2 26 -1 2 0 0 0 1 2 rows × 28 columns In [147]: Xclass_train, Xclass_test, Yclass_train, Yclass_test = train_test_split(Xclass, Yclass, test_size=0.25, random state=42) **Select Classification Models** In [148]: # create empty list and then populate it with the following models models = []models.append(('DT', DecisionTreeClassifier())) models.append(('RF', RandomForestClassifier())) models.append(('GB', GradientBoostingClassifier())) # create empty lists to hold results and model names results = []names = []**Classification Cross Validation** In [149]: **for** name, model **in** models: kfold = KFold(n splits=3, random state=43, shuffle=True) result = cross val score(model, Xclass train, Yclass train, cv=kfold, scoring='accuracy') names.append(name) results.append(result) #msg = '%s: %.4f (%.4f)' % (name, result.mean(), result.std()) # print results for i in range(len(names)): print(names[i], results[i].mean()) DT 0.30017777777778 RF 0.39533333333333333 GB 0.4048444444444444 Gradient Boosting Classifier appears to be the best bet In [158]: gb = GradientBoostingClassifier(max_depth=3, random_state=0) # , learning_rate=0.01, n_estimators=500 # train/fit the mode using amt as dv gbModel = gb.fit(Xclass_train, Yclass_train) # make predicitons with the trained/fit model gbPred = gbModel.predict(Xclass_test) # performance metrics print(accuracy_score(Yclass_test, gbPred)) print(classification_report(Yclass_test, gbPred)) 0.4132 precision recall f1-score support

 0.54
 0.86
 0.66
 1033

 0.56
 0.55
 0.56
 902

 0.45
 0.43
 0.44
 782

 0.42
 0.34
 0.37
 746

 0.36
 0.25
 0.30
 795

 0.28
 0.20
 0.23
 816

 0.27
 0.25
 0.26
 861

 0.29
 0.39
 0.33
 936

 0.44
 0.30
 0.36
 629

 1 2 3 4 5 7 8 0.41 7500 0.39 7500 accuracy 0.40 0.40 0.39 macro avg weighted avg 0.40 0.41 0.40 7500 The binning appears to have decreased the ability to accurately model In [159]: #plt.scatter(y_test, predictions, alpha = 0.5) plt.scatter(Yclass_test[0:200], gbPred[0:200], alpha = 0.5) plt.xlabel('Ground Truth') plt.ylabel('Predictions') #plt.plot(y_test[0:200],y_test[0:200], c='red') plt.show(); 9 • 8 7 6 Predictions 5 4 2 5 Ground Truth We can clearly see some pretty far-off preditctions using this model. Many points were off by a large number of bins **Classification Modelling - Default** In [160]: # use the limit bin as our dependent variable Ydef=df['DEFAULT_default'] Xdef = df.iloc[:,1:27]Xdef.head(2) Out[160]: LIMIT_BAL MARRIAGE AGE PAY_1 PAY_2 PAY_3 PAY_4 PAY_5 PAY_6 BILL_AMT1 ... PAY_AMT2 PAY_AMT3 PAY_AMT4 PAY 0 20000 3913 ... 689 1 120000 2682 ... 1000 1000 1000 2 rows × 26 columns In [161]: Xdef_train, Xdef_test, Ydef_train, Ydef_test = train_test_split(Xdef, test_size=0.25, random_state=42) **Select Classification Models** In [162]: # create empty list and then populate it with the following models models = []models.append(('DT', DecisionTreeClassifier())) models.append(('RF', RandomForestClassifier())) models.append(('GB', GradientBoostingClassifier())) # create empty lists to hold results and model names results = []names = []**Classification Cross Validation** In [163]: for name, model in models: kfold = KFold(n splits=3, random state=43, shuffle=True) result = cross_val_score(model, Xdef_train, Ydef train, cv=kfold, scoring='accuracy') names.append(name) results.append(result) #msq = '%s: %.4f (%.4f)' % (name, result.mean(), result.std()) # print results for i in range(len(names)): print(names[i], results[i].mean()) DT 0.7301333333333333 RF 0.815822222222222 GB 0.8213333333333334 In [166]: gb2 = GradientBoostingClassifier(max_depth=3, random_state=0) # , learning_rate=0.01, n_estimators=500 # train/fit the mode using amt as dv gbModel2 = gb2.fit(Xdef train, Ydef train) # make predicitons with the trained/fit model gbPred2 = gbModel2.predict(Xdef_test) # performance metrics target names=['not defaulted','defaulted'] print(accuracy score(Ydef test, gbPred2)) print(classification_report(Ydef_test, gbPred2, target_names=target_names)) 0.8201333333333334 precision recall f1-score support 0.84 0.95 0.66 0.36 not defaulted 0.89 5873 defaulted 0.46 1627 0.82 7500 accuracy 0.75 0.65 0.68 macro avg 7500 0.80 0.82 weighted avg 0.80 7500 It's close on cross_val, let's look at a decsion tree for visualization purposes In [168]: | dt = DecisionTreeClassifier(splitter='best', max_depth=3, min_samples_split=2) # train/fit the mode using region as dv, and binned by age & amt dtModel3 = dt.fit(Xdef train, Ydef train) # make predicitons with the trained/fit model dtPred3 = dtModel3.predict(Xdef test) # performance metrics print(accuracy score(Ydef test, dtPred3)) print(classification report(Ydef test, dtPred3)) 0.8205333333333333 precision recall f1-score support 0.84 0.95 0.89 5873 0.66 0.36 0.47 1627 0.82 7500 accuracy 0.75 0.65 0.68 macro avg 7500 7500 0.80 0.82 0.80 weighted avg In [173]: target names=['not defaulted','defaulted'] dot data = StringIO() export_graphviz(dtModel3, out file=dot data, filled=True, rounded=True, feature_names=Xdef_train.columns, class_names=target_names, label='all', precision=1, special_characters=True) graph = pydotplus.graph_from_dot_data(dot_data.getvalue()) Image(graph.create_png()) PAY_1 ≤ 1.5 gini = 0.3 samples = 22500 value = [17491, 5009] Out[173]: True PAY_2 ≤ 1.5 gini = 0.3 samples = 20122 value = [16777, 3345] PAY_3 ≤ -0.5 gini = 0.4 samples = 2378 value = [714, 1664] class = not defaulted class = defaulted PAY_6 ≤ 1.0 gini = 0.5 PAY_5 ≤ 1.0 gini = 0.4 samples = 2231 value = [633, 1598] BILL_AMT1 ≤ 2264.0 PAY AMT3 ≤ 573.5 gini = 0.5 samples = 18452 samples = 1670 samples = 147 value = [963, 707] class = not defaulted class = defaulted class = not defaulted class = not defaulted gini = 0.5 gini = 0.5 gini = 0.4 gini = 0.4 samples = 56 gini = 0.3 samples = 996 value = [225, 771] aini = 0.3gini = 0.4 samples = 474 value = [222, 252] samples = 4882 samples = 13570 samples = 1196 samples = 91 samples = 1235 value = [3818, 1064] value = [19, 37] = [408, 827] class = not defaulted class = not defaulted class = not defaulted class = defaulted class = not defaulted class = defaulted class = defaulted class = defaulted In [174]: graph.write_png("Default_DT.png") Out[174]: True In []: