%matplotlib inline In [271... df = pd.read csv('loan data.csv') In [272... df.head() Out[272... credit.policy purpose int.rate installment log.annual.inc dti fico days.with.cr.line revol.bal revol.util inq.last.6mths delinq.2yrs pub.rec not.fully.paid 1 debt\_consolidation 11.350407 19.48 737 52.1 0.1189 829.10 5639.958333 28854 0 0 0 76.7 credit\_card 0.1071 228.22 11.082143 14.29 707 0 2760.000000 33623 25.6 1 debt consolidation 366.86 10.373491 11.63 682 3511 0 0.1357 4710.000000 0 1 debt\_consolidation 0.1008 162.34 11.350407 8.10 712 2699.958333 33667 73.2 0 0 4740 39.5 credit\_card 0.1426 102.92 11.299732 14.97 667 0 4 4066.000000 In [273... df.info() <class 'pandas.core.frame.DataFrame'> RangeIndex: 9578 entries, 0 to 9577 Data columns (total 14 columns): Column Non-Null Count Dtype credit.policy 9578 non-null int64 9578 non-null object purpose int.rate 9578 non-null float64 installment float64 9578 non-null 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 ing.last.6mths 9578 non-null int64 11 deling.2yrs 9578 non-null int64 12 pub.rec 9578 non-null int64 13 not.fully.paid 9578 non-null int64 dtypes: float64(6), int64(7), object(1) memory usage: 1.0+ MB In [274... df.describe() int.rate installment log.annual.inc Out[274... credit.policy dti fico days.with.cr.line revol.util inq.last.6mths delinq.2yrs revol.bal pub.r€ count 9578.000000 9578.000000 9578.000000 9578.000000 9578.000000 9578.000000 9578.000000 9.578000e+03 9578.000000 9578.000000 9578.000000 9578.00000 0.804970 0.122640 319.089413 10.932117 12.606679 710.846314 4560.767197 1.691396e+04 46.799236 1.577469 0.163708 0.06212 mean 0.546215 0.396245 0.026847 207.071301 0.614813 37.970537 2496.930377 3.375619e+04 2.200245 0.26212 std 6.883970 29.014417 7.547502 178.958333 0.000000e+00 0.00000 min 0.000000 0.060000 15.670000 0.000000 612.000000 0.000000 0.000000 0.000000 682.000000 0.00000 25% 1.000000 0.103900 163.770000 10.558414 7.212500 2820.000000 3.187000e+03 22.600000 0.000000 0.000000 4139.958333 8.596000e+03 1.000000 0.122100 707.000000 0.00000 **50%** 268.950000 10.928884 12.665000 46.300000 1.000000 0.000000 432.762500 1.000000 0.140700 11.291293 737.000000 5730.000000 1.824950e+04 70.900000 2.000000 0.000000 0.00000 **75%** 17.950000 29.960000 0.216400 940.140000 14.528354 827.000000 5.00000 1.000000 17639.958330 1.207359e+06 119.000000 33.000000 13.000000 max FICO Score Distribution by Loan Payment Status sns.histplot(x='fico', data=df, hue='not.fully.paid', bins=20) <Axes: xlabel='fico', ylabel='Count'> not.fully.paid 1200 \_\_\_\_1 1000 800 600 400 200 700 650 750 800 fico Loan Purpose Count by Payment Status plt.figure(figsize=(12,5)) sns.countplot(x='purpose', data=df, hue='not.fully.paid') <Axes: xlabel='purpose', ylabel='count'> 3500 not.fully.paid 3000 2500 -2000 1500 1000 500 debt\_consolidation all\_other home\_improvement small\_business credit\_card major\_purchase educational purpose Relationship Between FICO Score and Interest Rate by Payment Status sns.jointplot(x='fico', y='int.rate', data=df, hue='not.fully.paid') <seaborn.axisgrid.JointGrid at 0x1aa433937f0> not.fully.paid 0.225 0.200 0.175 o.150 0.125 0.100 0.075 0.050 650 700 750 800 850 600 fico Convert Categorical Variable 'purpose' into Dummy Variables In [278... data = pd.get dummies(df,columns=['purpose'],drop\_first=True) data.info() In [279... <class 'pandas.core.frame.DataFrame'> RangeIndex: 9578 entries, 0 to 9577 Data columns (total 19 columns): Column Non-Null Count Dtype 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 int64 revol.bal 9578 non-null revol.util 9578 non-null float64 inq.last.6mths 9578 non-null int64 9578 non-null delinq.2yrs int64 11 pub.rec 9578 non-null int64 not.fully.paid 9578 non-null int64 purpose credit card 9578 non-null bool purpose debt consolidation 9578 non-null bool 9578 non-null purpose educational bool purpose home improvement 9578 non-null bool purpose major purchase 9578 non-null bool purpose small business 9578 non-null bool dtypes: bool(6), float64(6), int64(7) memory usage: 1.0 MB Split Data into Training and Test Sets In [280... X = data.drop('not.fully.paid', axis=1) y = data['not.fully.paid'] from sklearn.model selection import train test split X\_train, X\_test, y\_train, y\_test = train\_test\_split(X, y, test\_size=0.30) Train Decision Tree Model and Make Predictions In [281... from sklearn.tree import DecisionTreeClassifier, plot tree Dtree = DecisionTreeClassifier() Dtree fit(X train, y train) predictions = Dtree.predict(X test) Evaluate Model Performance In [282... from sklearn.metrics import classification report, confusion matrix print(classification report(y test, predictions)) pd.DataFrame(confusion matrix(y test, predictions)) precision recall f1-score support 0.84 2389 0.85 0.83 0.23 0.25 0.24 485 2874 0.73 accuracy 0.54 2874 0.54 0.54 macro avg 0.74 0.74 2874 weighted avg 0.73 Out[282... **0** 1985 404 364 121 The Decision Tree model performed moderately in predicting whether a loan would be fully paid or not. Its precision and recall are acceptable, but it struggled with some incorrect predictions, as seen in the confusion matrix. It tends to overfit the training data, which can lead to less reliable performance on unseen data. Train Random Forest Classifier In [283... from sklearn.ensemble import RandomForestClassifier rfc = RandomForestClassifier(n estimators=600) rfc.fit(X train, y train) rfc predict = rfc.predict(X test) Evaluate Random Forest Classifier from sklearn.metrics import classification report, confusion matrix print(classification\_report(y\_test, rfc\_predict)) pd.DataFrame(confusion\_matrix(y\_test, rfc\_predict)) precision recall f1-score support 0.83 1.00 0.91 2389 0.53 0.02 0.04 485 0.83 2874 accuracy 0.47 0.51 2874 0.68 macro avg 0.78 weighted avg 0.83 0.76 2874 0 1 Out[284... **0** 2381 8 **1** 476 9 The Random Forest Classifier outperformed the Decision Tree by a significant margin. It achieved higher precision, recall, and F1-score, showing it was better at predicting loan repayment status. The RFC model's performance is more consistent, with fewer errors in the confusion matrix, indicating that it generalizes well to new data. Summary In conclusion, the Random Forest Classifier is the better model for this loan prediction task, providing more accurate and reliable results than the Decision Tree. The Decision Tree model is simpler but less effective, while RFC offers improved performance due to its ensemble nature.

In [270...

import numpy as np

import pandas as pd

import seaborn as sns

import matplotlib.pyplot as plt