**Loan Prediction** In this model we will predict the eligibility of a person for loan. Importing required libraries In [156]: import pandas as pd import numpy as np import matplotlib.pyplot as plt import seaborn as sns %matplotlib inline Importing our training data In [157]: train=pd.read\_csv('train.csv') In [158]: train.head() Out[158]: Loan\_ID Gender Married Dependents Education Self\_Employed ApplicantIncome CoapplicantIncome LoanAmount Loa **0** LP001002 5849 0.0 NaN Male No 0 Graduate No **1** LP001003 Male Yes Graduate No 4583 1508.0 128.0 **2** LP001005 Graduate 3000 0.0 66.0 Male Yes Yes Not **3** LP001006 2583 2358.0 120.0 Male Yes No Graduate **4** LP001008 Male 0 Graduate 6000 0.0 141.0 **Checking for null values** In [159]: train.isnull().sum() 0 Out[159]: Loan\_ID Gender 13 Married 3 15 Dependents Education 0 Self\_Employed 32 ApplicantIncome 0 0 CoapplicantIncome 22 LoanAmount Loan\_Amount\_Term 14 Credit\_History 50 Property\_Area 0 Loan\_Status 0 dtype: int64 We can see that we have quite a lot null values. Therefore, first we will determine the required features and then deal with the null values. Finding relation between different features and our target variable Relation between gender and loan status. In [160]: sns.countplot(x='Loan\_Status', hue='Gender', data=train) Out[160]: <matplotlib.axes.\_subplots.AxesSubplot at 0x1f53d072e88> Gender Male 300 Female 250 뉟 200 8 150 100 50 Loan\_Status We see that the gender of the applicant does not play a significant role in the process. Relation between relationship status and loan status. In [161]: sns.countplot(x='Loan\_Status', hue='Married', data=train) Out[161]: <matplotlib.axes.\_subplots.AxesSubplot at 0x1f540ceb148> Married 250 200 150 100 50 Loan\_Status We see that married people have a greater chance of getting the loan. Relation between dependents and loan status. In [162]: sns.countplot(x='Loan\_Status', hue='Dependents', data=train) Out[162]: <matplotlib.axes.\_subplots.AxesSubplot at 0x1f540cf5a88> Dependents 0 200 1 3+ 150 100 50 Loan\_Status People with 0 dependents have more chances of getting the loan. Relation between education level of the applicant and loan status. sns.countplot(x='Loan\_Status', hue='Education', data=train) In [163]: Out[163]: <matplotlib.axes.\_subplots.AxesSubplot at 0x1f5414aa948> 350 Education Graduate 300 Not Graduate 250 달 200 8 150 100 50 Loan\_Status Graduates have more chances of getting the loan. Relation between the applicant's employment type and loan status. In [164]: sns.countplot(x='Loan\_Status', hue='Self\_Employed', data=train) Out[164]: <matplotlib.axes.\_subplots.AxesSubplot at 0x1f54150ee48> 350 Self\_Employed 300 Yes 250 는 200 100 50 Loan\_Status People who are not self employed have higher chances of getting the loan. Relation between loan amount term and loan status. In [168]: sns.countplot(x='Loan\_Status', hue='Loan\_Amount\_Term', data=train) Out[168]: <matplotlib.axes.\_subplots.AxesSubplot at 0x1f54283b708> Loan Amount Term 350 12.0 60.0 250 1 200 200 180.0 240.0 150 100 480.0 50 Loan\_Status Relation between credit history and loan status. In [169]: | sns.countplot(x='Loan\_Status', hue='Credit\_History', data=train) Out[169]: <matplotlib.axes.\_subplots.AxesSubplot at 0x1f545c47b08> Credit History 350 0.0 1.0 300 250 200 150 100 50 Loan\_Status Credit history play a significant role in the process. Relation between property area and loan status. In [170]: sns.countplot(x='Loan\_Status', hue='Property\_Area', data=train) Out[170]: <matplotlib.axes.\_subplots.AxesSubplot at 0x1f545c479c8> 175 Property\_Area Urban 150 Rural 125 달 100 75 50 25 Ÿ Loan\_Status Whether the applicant will get the loan or not does not depend on his/her property area. Filtering out required features In [171]: | train.drop(['Gender', 'CoapplicantIncome', 'Loan\_Amount\_Term', 'Property\_Area'], axis=1, inplace= True); In [172]: train.head() Out[172]: Education Self\_Employed ApplicantIncome LoanAmount Credit\_History Loan\_Status Loan\_ID Married Dependents **0** LP001002 No Graduate 5849 NaN 1.0 **1** LP001003 4583 128.0 Ν Yes Graduate No 1.0 Υ **2** LP001005 Graduate Yes 3000 66.0 1.0 Yes 3 LP001006 0 Not Graduate 2583 120.0 1.0 Υ No 4 LP001008 6000 141.0 1.0 Υ No Graduate No Converting categorical data into indicator data married=pd.get\_dummies(train['Married'], drop\_first=True) In [173]: education=pd.get\_dummies(train['Education'], drop\_first=True) self\_employed=pd.get\_dummies(train['Self\_Employed'],drop\_first=True) In [174]: train=pd.concat([train, married, education], axis=1) train.head() Out[174]: Loan\_ID Married Dependents Education Self\_Employed ApplicantIncome LoanAmount Credit\_History Loan\_Status Ye **0** LP001002 Graduate 5849 NaN 1.0 Υ 4583 128.0 1.0 **1** LP001003 Yes Graduate No Ν **2** LP001005 3000 1.0 Graduate Yes 66.0 **3** LP001006 No 2583 120.0 1.0 Yes Graduate 4 LP001008 6000 141.0 1.0 No Graduate No In [175]: train.rename(columns={'Yes':'married'},inplace=True) train.head() Out[175]: Loan\_ID Married Dependents Education Self\_Employed ApplicantIncome LoanAmount Credit\_History Loan\_Status ma **0** LP001002 0 Graduate No 5849 1.0 NaN **1** LP001003 Yes Graduate No 4583 128.0 1.0 **2** LP001005 3000 66.0 1.0 Yes Graduate Yes **3** LP001006 2583 120.0 1.0 No Graduate 4 LP001008 Graduate 6000 141.0 1.0 In [176]: train=pd.concat([train, self\_employed], axis=1) train.head() Out[176]: Loan\_ID Married Dependents Education Self\_Employed ApplicantIncome LoanAmount Credit\_History Loan\_Status ma **0** LP001002 1.0 No Graduate No 5849 NaN **1** LP001003 Yes Graduate No 4583 128.0 1.0 **2** LP001005 Graduate 3000 66.0 1.0 Yes Yes **3** LP001006 No 2583 120.0 1.0 Graduate **4** LP001008 Graduate 6000 141.0 1.0 In [177]: train.rename(columns={'Yes':'self\_employed'},inplace=True) train.head() Out[177]: Loan\_ID Married Dependents Education Self\_Employed ApplicantIncome LoanAmount Credit\_History Loan\_Status ma **0** LP001002 0 Graduate No 5849 NaN 1.0 No **1** LP001003 Graduate No 4583 128.0 1.0 Ν Yes **2** LP001005 Graduate Yes 3000 66.0 1.0 Yes **3** LP001006 2583 120.0 1.0 Yes No Graduate **4** LP001008 6000 141.0 1.0 No 0 Graduate No In [178]: train.drop(['Married', 'Education', 'Self\_Employed'], axis=1, inplace=True) In [179]: train.head() Out[179]: Loan\_ID Dependents ApplicantIncome LoanAmount Credit\_History Loan\_Status married Not Graduate self\_employed **0** LP001002 5849 NaN 1.0 **1** LP001003 1 4583 128.0 1.0 Ν 1 0 0 **2** LP001005 3000 1.0 **3** LP001006 0 2583 120.0 Υ 0 1.0 1 1 **4** LP001008 6000 141.0 1.0 In [180]: train.isnull().sum() Out[180]: Loan\_ID 0 Dependents 15 ApplicantIncome 0 LoanAmount 22 Credit\_History 50 0 Loan\_Status married 0 0 Not Graduate self\_employed dtype: int64 We see that dependents is also categorical data. Therefore, we need to convert it to indicator data. In [181]: | dependents=pd.get\_dummies(train['Dependents'], drop\_first=True) In [182]: | train=pd.concat([train,dependents],axis=1) train.head() Out[182]: Not Loan\_ID Dependents ApplicantIncome LoanAmount Credit\_History Loan\_Status married self\_employed 1 2 Graduate **0** LP001002 0 0 ( 0 5849 1.0 NaN **1** LP001003 1 4583 128.0 1.0 Ν 1 0 0 1 ( **2** LP001005 0 3000 66.0 1.0 1 0 ( **3** LP001006 0 2583 120.0 1.0 Υ 1 0 0 ( 1 **4** LP001008 6000 141.0 1.0 0 0 ( In [183]: train.drop('Dependents', axis=1, inplace=True) train.head() Out[183]: Loan\_ID ApplicantIncome LoanAmount Credit\_History Loan\_Status married Not Graduate self\_employed 1 2 3+ **0** LP001002 5849 0 0 0 0 NaN 1.0 **1** LP001003 4583 128.0 1.0 0 1 0 0 **3** LP001006 120.0 1.0 0 0 0 0 **4** LP001008 6000 141.0 1.0 0 0 0 0 **Checking for null values** In [184]: train.isnull().sum() Out[184]: Loan\_ID ApplicantIncome 0 22 LoanAmount Credit\_History 50 Loan\_Status married Not Graduate self\_employed 1 2 3+ dtype: int64 Filling up null values In [185]: def fill\_loan(cols): amount=cols[0] if pd.isnull(amount): return train['LoanAmount'].mean() return amount In [186]: train['LoanAmount']=train[['LoanAmount']].apply(fill\_loan,axis=1) In [187]: train.isnull().sum() Out[187]: Loan\_ID 0 ApplicantIncome 0 LoanAmount Credit\_History Loan\_Status married Not Graduate self\_employed 1 2 3+ dtype: int64 In [188]: def fill\_credithistory(cols): credit\_history=cols[0] if pd.isnull(credit\_history): return 1 return credit\_history In [189]: | train['Credit\_History']=train[['Credit\_History']].apply(fill\_credithistory,axis=1) In [190]: train.isnull().sum() Out[190]: Loan\_ID ApplicantIncome 0 LoanAmount 0 Credit\_History Loan\_Status married Not Graduate self\_employed 1 2 3+ dtype: int64 We have successfully removed all null values. **Building model** Importing required libraries. In [191]: from sklearn.tree import DecisionTreeClassifier In [192]: classifier=DecisionTreeClassifier() **Splitting our data** Separating our features and target variable. In [193]: x=train.drop(['Loan\_Status', 'Loan\_ID'], axis=1) y=train['Loan\_Status'] Importing required libraries. In [194]: | from sklearn.model\_selection import train\_test\_split In [195]: | x\_train, x\_test, y\_train, y\_test=train\_test\_split(x, y, test\_size=0.2) **Training and testing our model** In [196]: classifier.fit(x\_train,y\_train) Out[196]: DecisionTreeClassifier(ccp\_alpha=0.0, 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='deprecated', random\_state=None, splitter='best') In [197]: y\_hat=classifier.predict(x\_test) In [198]: y\_hat 'N', Out[198]: array(['Y', 'N', 'N', 'Y', 'Y', 'Υ', 'Υ', 'N', 'N', 'Y', 'Y', 'Y', 'Y', 'N', 'N', 'Y', 'Y', 'N', 'Y', 'N', 'Y', 'N', 'Y', 'N', 'Y', 'N', 'Y', **Finding out accuracy score** In [199]: **from sklearn.metrics import** accuracy\_score In [200]: | accuracy=accuracy\_score(y\_test,y\_hat) accuracy Out[200]: 0.7073170731707317 Accuracy is quite high. Importing testing data test=pd.read\_csv('test.csv') In [201]: In [202]: test.head() Out[202]: Loan\_ID Gender Married Dependents Education Self\_Employed ApplicantIncome CoapplicantIncome LoanAmount Loa **0** LP001015 5720 110.0 Male Graduate 1500 **1** LP001022 Male Yes Graduate No 3076 126.0 5000 1800 208.0 **2** LP001031 Male Graduate No Graduate 2340 100.0 **3** LP001035 Male Yes No 2546 **4** LP001051 3276 Male 78.0 No No Graduate **Modifying testing data** In [203]: test.drop(['Gender','CoapplicantIncome','Loan\_Amount\_Term','Property\_Area'],axis=1,inplace=T rue); married=pd.get\_dummies(test['Married'],drop\_first=True) In [204]: education=pd.get\_dummies(test['Education'], drop\_first=True) self\_employed=pd.get\_dummies(test['Self\_Employed'],drop\_first=True) In [205]: test=pd.concat([test,married,education],axis=1) In [206]: test.rename(columns={'Yes':'married'},inplace=True) test.head() Out[206]: Loan\_ID Married Dependents Education Self\_Employed ApplicantIncome LoanAmount Credit\_History married Gradua **0** LP001015 0 Graduate No 5720 110.0 1.0 1 **1** LP001022 No 3076 126.0 1.0 1 Graduate **2** LP001031 Graduate No 5000 208.0 1.0 Yes 1 **3** LP001035 Graduate No 2340 100.0 NaN **4** LP001051 3276 78.0 1.0 Graduate In [207]: test=pd.concat([test,self\_employed],axis=1) In [208]: | test.rename(columns={'Yes':'self\_employed'},inplace=True) test.head() Out[208]: Loan\_ID Married Dependents Education Self\_Employed ApplicantIncome LoanAmount Credit\_History married **0** LP001015 Yes 0 Graduate 5720 110.0 1.0 **1** LP001022 Yes Graduate No 3076 126.0 1.0 1 **2** LP001031 5000 208.0 2 Graduate No 1.0 Yes **3** LP001035 Graduate No 2340 100.0 NaN **4** LP001051 No 3276 78.0 1.0 In [209]: test.drop(['Married', 'Education', 'Self\_Employed'], axis=1, inplace=True) In [210]: dependents=pd.get\_dummies(test['Dependents'], drop\_first=True) In [211]: test=pd.concat([test,dependents],axis=1) In [212]: | test.drop('Dependents', axis=1, inplace=True) test.head() Out[212]: Loan\_ID ApplicantIncome LoanAmount Credit\_History married Not Graduate self\_employed 1 2 3+ 0 0 0 0 **0** LP001015 5720 110.0 1.0 **1** LP001022 3076 126.0 1.0 1 0 1 0 0 **2** LP001031 5000 1 208.0 1.0 0 0 1 0 0 0 1 0 **3** LP001035 2340 100.0 NaN 1 0 **4** LP001051 3276 78.0 1.0 0 0 0 0 In [213]: test.isnull().sum() Out[213]: Loan\_ID 0 0 ApplicantIncome 5 LoanAmount Credit\_History 29 married Not Graduate self\_employed 1 2 3+ dtype: int64 In [214]: def fill\_loan\_test(cols): amount=cols[0] if pd.isnull(amount): return test['LoanAmount'].mean() return amount In [215]: test['LoanAmount']=test[['LoanAmount']].apply(fill\_loan\_test,axis=1) In [216]: | test['Credit\_History']=test[['Credit\_History']].apply(fill\_credithistory,axis=1) In [217]: test.isnull().sum() Out[217]: Loan\_ID ApplicantIncome 0 LoanAmount Credit\_History married Not Graduate self\_employed 2 dtype: int64 **Making predictions on testing data** In [218]: x=test.drop(['Loan\_ID'], axis=1)

In [219]: predictions=classifier.predict(x)

'N', 'Y', 'Y', 'Y',

'Y', 'Y', 'Y',

'Y',

'N',

'Y',

'N',

'Y',

'Y',

'Y',

'Y',

'N',

'N',

'Υ',

'Υ',

'N',

'Y',

'N',

'Υ',

'Υ',

'Y', 'N',

In [221]: predictions.size

In [222]: submission=pd.DataFrame({

Out[221]: 367

'N', 'N',

'Y', 'N',

'Y', 'N', 'N', 'Y', 'Y',

'Y',

'Υ',

'Υ',

'Υ',

'N',

'Υ',

'Υ',

'Y',

'N',

'Y'

'Y',

'Υ',

'Υ',

'Υ',

'N',

'Y',

'N',

'Y',

'Y',

'N',

'Loan\_ID':test['Loan\_ID'],
'Loan\_Status':predictions

'Y',

'Υ',

'Υ',

'Y',

'N',

'Υ',

'N',

'N',

'Υ',

'Y',

'Υ',

'Υ',

'Υ',

'N',

'N',

'N',

'Y',

'Υ',

'Y',

'Y',

'Y', 'Y', 'N'], dtype=object)

submission.to\_csv('submission.csv',index=False)

'Υ',

'Υ',

'Y',

'N',

'N',

'N',

'Y',

'Y',

'Υ',

'N',

'Y',

'N',

'Y',

'Y',

'N',

'Y',

'N',

'N',

'Y',

'Υ',

'N',

'N',

'N',

'Y',

'Υ',

'N',

'Υ',

'Υ',

'Υ',

'Υ',

'N',

'Υ',

'Υ',

'Y',

'Υ',

'Y',

'Υ',

'Y',

'Υ',

'Υ',

'Y',

'N',

'Υ',

'Y',

'N',

'Y',

'N',

'Y',

'Υ',

'N',

'N',

'Υ',

'Y',

'Y',

'Y',

'Υ',

'Y',

'Y',

'Y',

'N',

'Υ',

'Υ',

'Y',

'N',

'Y', 'N', 'Y', 'Y',

'N',

'Y',

'Y',

'N',

'Y',

'Υ',

'N',

'N',

'Y',

'Y',

'Y',

'N',

'Y',

'N',

'Y',

'Υ',

'N',

'N',

'N',

'N',

'Y',

'Υ',

'N',

'Y',

'N',

'Υ',

'Υ',

'Y',

'Υ',

'Υ',

'N',

'Y',

'Y',

'N',

'N',

'N',

'Y', 'Y',

'N', 'Y',

'N',

'Υ',

'Υ',

'Y',

'Υ',

'N',

'Υ',

'Υ',

'Ν',

'Υ',

'Υ',

'N',

'Υ',

'Υ',

'Υ',

'Υ',

'N',

'N',

'Υ',

'Υ',

In [220]: predictions

Out[220]: array(['Y', 'Y', 'Y',