1. Credit card applications Commercial banks receive a lot of applications for credit cards. Many of them get rejected for many reasons, like high loan balances, low income levels, or too many inquiries on an individual's credit report, for example. Manually analyzing these applications is mundane, errorprone, and time-consuming (and time is money!). Luckily, this task can be automated with the power of machine learning and pretty much every commercial bank does so nowadays. In this notebook, we will build an automatic credit card approval predictor using machine learning techniques, just like the real banks do. We'll use the Credit Card Approval dataset from the UCI Machine Learning Repository. The structure of this notebook is as follows: -First, we will start off by loading and viewing the dataset. -We will see that the dataset has a mixture of both numerical and non-numerical features, that it contains values from different ranges, plus that it contains a number of missing entries. -We have to preprocess the dataset to ensure the machine learning model we choose can make good predictions. -Next we will do some exploratory data analysis. -Finally, we will build a machine learning model that can predict if an individual's application for a credit card will be accepted. First, loading and viewing the dataset. We find that since this data is confidential, the contributor of the dataset has anonymized the feature names. In [10]: # Import pandas import pandas as pd df = pd.read csv('cc approvals.csv', header=None) df Out[10]: 0 7 8 9 10 11 13 14 15 30.83 0.000 1.25 00202 b 58.67 4.460 3.04 0.500 00280 24.50 1.50 b 27.83 1.540 3.75 00100 20.17 00120 b 5.625 1.71 10.085 1.25 00260 685 b 21.08 h 686 а 22.67 0.750 2.00 00200 00200 687 25.25 13.500 2.00 688 b 17.92 0.205 0.04 00280 c h 8.29 689 b 35.00 3.375 690 rows × 16 columns 2. Inspecting the applications The features of this dataset have been anonymized to protect the privacy. The probable features in a typical credit card application are Gender, Age, Debt, Married, BankCustomer, EducationLevel, Ethnicity, YearsEmployed, PriorDefault, Employed, CreditScore, DriversLicense, Citizen, ZipCode, Income and finally the ApprovalStatus. This gives us a pretty good starting point, and we can map these features with respect to the columns in the output. As we can see from our first glance at the data, the dataset has a mixture of numerical and non-numerical features. This can be fixed with some preprocessing, but before we do that, let's learn about the dataset a bit more to see if there are other dataset issues that need to be df.info() In [11]: <class 'pandas.core.frame.DataFrame'> RangeIndex: 690 entries, 0 to 689 Data columns (total 16 columns): # Column Non-Null Count Dtype 690 non-null 0 0 object 1 1 690 non-null object 2 690 non-null float64 3 3 690 non-null object 4 4 690 non-null object 5 5 690 non-null object 690 non-null 6 object 7 7 690 non-null float64 8 object 8 690 non-null 9 9 690 non-null object 10 10 690 non-null int64 object 11 11 690 non-null 12 12 690 non-null object 13 13 690 non-null object 14 14 690 non-null int64 15 15 690 non-null object dtypes: float64(2), int64(2), object(12) memory usage: 86.4+ KB In [12]: df.describe() Out[12]: 2 7 10 14 **count** 690.000000 690.000000 690.00000 690.000000 4.758725 2.223406 2.40000 1017.385507 mean std 4.978163 3.346513 4.86294 5210.102598 min 0.000000 0.000000 0.00000 0.000000 25% 1.000000 0.165000 0.00000 0.000000 50% 2.750000 1.000000 0.00000 5.000000 75% 7.207500 2.625000 3.00000 395.500000 28.000000 28.500000 67.00000 100000.000000 max In [13]: df.tail(20) Out[13]: 2 3 4 7 8 9 10 11 12 14 15 **670** b 47.17 5.835 u g 00465 150 w v 5.500 **671** b 25.83 12.835 u g cc v 0.500 00000 2 g **672** a 50.25 0.835 u g aa v 0.500 00240 117 **673** ? 29.50 2.000 y p e h 2.000 00256 17 **674** a 37.33 2.500 u g i h 0.210 00260 f f 0 00240 675 a 41.58 1.040 u g aa v 0.665 q h 0.085 **676** a 30.58 10.665 u g f t 12 00129 7.250 u g m v 0.040 00100 **677** b 19.42 1 **678** a 17.92 10.210 u g ff ff 0.000 00000 50 **679** a 20.08 1.250 u g c v 0.000 00000 g 0.290 u g k v 0.290 **680** b 19.50 f f 0 00280 g 00176 537 d h 3.000 f f 0 **681** b 27.83 1.000 y p 3.290 u g i v 0.335 f f 0 t g 00140 **682** b 17.08 683 b 36.42 3 0.750 v 0.585 00240 40.58 3.290 00400 **685** b 21.08 10.085 у р 1.250 00260 0 e h 686 22.67 0.750 u g 2.000 00200 ff 2.000 687 25.25 13.500 р ff 00200 688 17.92 0.205 u aa ٧ 0.040 00280 750 b 35.00 3.375 u g c h 8.290 00000 3. Dealing with missing values Our dataset contains both numeric and non-numeric data (specifically data that are of float64, int64 and object types). The features 2, 7, 10 and 14 contain numeric values (of types float64, float64, int64 and int64 respectively) and all the other features contain non-numeric values. The dataset also contains values from several ranges. Some features have a value range of 0 - 28, some have a range of 2 - 67, and some have a range of 1017 - 100000. Apart from these, we can get useful statistical information (like mean, max, and min) about the features that have numerical values. Finally, the dataset has missing values. The missing values in the dataset are labeled with '?', which can be seen in the last cell's output. Now, let's temporarily replace these missing value question marks with NaN. In [14]: import numpy as np # Replace the '?'s with NaN df = df.replace('?', np.nan) df.tail(20) Out[14]: 1 2 3 4 5 6 8 9 10 11 14 15 00465 670 47.17 5.835 5.500 W 150 671 b 25.83 12.835 u g cc v 0.500 00000 50.25 0.835 u 00240 672 g aa v 0.500 117 29.50 673 NaN 2.000 e h 2.000 00256 17 a 37.33 2.500 u g i h 0.210 00260 246 674 675 a 41.58 1.040 u g aa v 0.665 00240 676 a 30.58 10.665 u g q h 0.085 f t 12 00129 677 b 19.42 7.250 u m v 0.040 00100 678 a 17.92 10.210 ff ff 0.000 00000 50 u 679 a 20.08 1.250 u g c v 0.000 00000 0 680 b 19.50 0.290 00280 u 0.290 681 b 27.83 1.000 у р d h 3.000 00176 b 17.08 3.290 0.335 00140 682 u g 683 b 36.42 0.750 d v 0.585 00240 3 684 b 40.58 3.290 u 00400 m v 3.500 10.085 685 b 21.08 1.250 00260 0 686 a 22.67 0.750 C V 2.000 00200 394 687 a 25.25 13.500 y p ff ff 2.000 00200 688 17.92 0.205 u 00280 750 g aa v 0.040 3.375 u g c h 8.290 f f 689 b 35.00 We replaced all the question marks with NaNs. This will help us in the next missing value treatment that we are going to do Ignoring missing values can affect the performance of a machine learning model heavily. While ignoring the missing values our machine learning model may miss out on information about the dataset that may be useful for its training. Then, there are many models which cannot handle missing values implicitly. So, to avoid this problem, we are going to impute the missing values with a strategy called mean imputation. In [15]: | df.fillna(df.mean(),inplace=True) df.isnull().sum() Out[15]: 0 12 12 0 3 6 6 0 11 0 12 0 13 13 0 15 0 dtype: int64 We have taken care of the missing values present in the numeric columns. There are still some missing values to be imputed for columns 0, 1, 3, 4, 5, 6 and 13. All of these columns contain non-numeric data and this why the mean imputation strategy would not work here. This needs a different treatment. We are going to impute these missing values with the most frequent values as present in the respective columns. In [16]: # Iterate over each column of dataframe for col in df.columns: # Check if the column is of object type if df[col].dtypes == 'object': # Impute with the most frequent value df = df.fillna(df[col].value_counts().index[0]) df.isnull().sum() Out[16]: 0 0 2 0 3 0 5 6 0 7 0 10 11 12 13 14 0 15 0 dtype: int64 The missing values are now successfully handled. 4. Preprocessing the data There is still some minor but essential data preprocessing needed before we proceed towards building our machine learning model. Three main tasks are: 1. Convert the non-numeric data into numeric. 2. Split the data into train and test sets. 3. Scale the feature values to a uniform range. First, we will convert all the non-numeric values into numeric ones. We do this because not only it results in a faster computation but also many machine learning models (like XGBoost) require the data to be in a strictly numeric format. We will do this by using a technique called label encoding. In [17]: # import Labelencoder from sklearn.preprocessing import LabelEncoder le = LabelEncoder() for col in df.columns: # Compare if the dtype is object if df[col].dtypes == 'object': # Use LabelEncoder to do the numeric transformation df[col] = le.fit transform(df[col]) 5. Splitting the dataset into train and test sets We have converted all the non-numeric values to numeric ones. Now, we split our data into train set and test set to prepare our data for two different phases of machine learning modeling: training and testing. Ideally, no information from the test data should be used to scale the training data or should be used to direct the training process of a machine learning model. Hence, we first split the data and then apply the scaling. Features like DriversLicense and ZipCode are not as important as the other features in the dataset for predicting credit card approvals. We should drop them to design our machine learning model with the best set of features. In [18]: # Import train test split from sklearn.model_selection import train test split # Drop the features 11 and 13 and convert the DataFrame to a NumPy array df = df.drop([11,13],axis=1)df = df.valuesX = df[:, 0:13]y = df[:, 13]X train, X test, y train, y test = train test split(X,y,test size=0.3,random state=1) C:\Users\Pranav.LAPTOP-HOVCQVL6\anaconda3\lib\importlib\ bootstrap.py:219: RuntimeWarning: numpy.ufun c size changed, may indicate binary incompatibility. Expected 192 from C header, got 216 from PyObjec return f(*args, **kwds) In [19]: X train Out[19]: array([[0.000e+00, 3.210e+02, 3.350e-01, ..., 0.000e+00, 0.000e+00, 2.197e+03], [1.000e+00, 2.550e+02, 4.915e+00, ..., 0.000e+00, 0.000e+00, 1.442e+03], [1.000e+00, 1.080e+02, 8.350e-01, ..., 0.000e+00, 0.000e+00, 0.000e+00], [0.000e+00, 2.230e+02, 5.000e+00, ..., 0.000e+00, 0.000e+00, 0.000e+00], [0.000e+00, 4.800e+01, 1.835e+00, ..., 5.000e+00, 0.000e+00, [0.000e+00, 7.500e+01, 1.175e+01, ..., 2.000e+00, 0.000e+00,5.510e+02]]) In [20]: y train Out[20]: array([1., 0., 1., 1., 0., 1., 0., 1., 0., 1., 1., 1., 1., 1., 1., 0., 0., 1., 1., 0., 0., 0., 0., 1., 0., 1., 0., 1., 1., 1., 0., 0., 0., 0., 0., 0., 1., 0., 1., 1., 1., 0., 1., 1., 0., 0., 0., 0., 1., 1., 1., 1., 1., 1., 1., 0., 1., 1., 1., 0., 1., 1., 0., 0., 0., 0., 0., 0., 1., 1., 0., 0., 1., 1., 0., 1., 1., 0., 1., 1., 0., 1., 1., 0., 1., 0., 1., 1., 1., 1., 0., 0., 1., 1., 1., 1., 0., 0., 0., 0., 1., 0., 0., 1., 0., 0., 1., 0., 0., 1., 1., 1., 0., 1., 0., 0., 1., 0., 0., 0., 1., 1., 0., 0., 0., 1., 1., 0., 1., 1., 0., 1., 1., 1., 1., 1., 1., 1., 1., 1., 0., 1., 1., 0., 0., 0., 0., 0., 1., 0., 1., 0., 0., 1., 0., 0., 0., 0., 1., 1., 1., 0., 1., 1., 1., 1., 1., 1., 1., 0., 1., 1., 0., 0., 0., 1., 0., 1., 1., 1., 0., 0., 1., 1., 0., 1., 1., 1., 1., 0., 0., 1., 1., 0., 0., 1., 1., 0., 0., 1., 0., 1., 1., 0., 1., 0., 1., 0., 1., 0., 1., 0., 1., 0., 1., 1., 0., 0., 1., 0., 0., 1., 1., 0., 1., 1., 0., 1., 1., 1., 1., 1., 0., 0., 0., 1., 0., 0., 1., 0., 0., 0., 1., 0., 0., 1., 0., 0., 1., 0., 0., 0., 0., 0., 0., 1., 1., 1., 1., 1., 1., 1., 0., 0., 1., 1., 0., 0., 1., 1., 1., 0., 0., 1., 1., 1., 1., 1., 1., 0., 0., 1., 1., 1., 0., 0., 1., 0., 1., 1., 0., 1., 0., 1., 1., 1., 0., 0., 1., 0., 0., 0., 1., 0., 1., 0., 1., 0., 1., 1., 0., 1., 1., 0., 1., 1., 1., 0., 1., 0., 1., 1., 1., 1., 0., 0., 1., 0., 1., 1., 0., 1., 0., 1., 0., 1., 1., 1., 0., 0., 1., 0., 0., 1., 0., 1., 0., 0., 1., 1., 0., 1., 0., 1., 0., 1., 0., 1., 0., 1., 0., 1., 0., 0., 1., 1., 0., 1., 0., 0., 1., 0., 1., 0., 0., 1., 1., 1., 1., 1., 1., 1., 0., 1., 0., 1., 1., 0., 1., 0., 1., 0., 1., 0., 0., 1., 0., 1., 1., 1., 1., 0., 0., 1., 0., 1., 1., 1., 0., 0., 1., 1., 1., 1., 0., 0., 0., 0., 1., 1., 0., 1., 0., 0., 0., 0., 1., 0., 1., 1., 1., 1., 1., 1., 1., 1., 1., 1., 0., 1., 0., 1., 0., 0., 0., 1., 0., 1., 1., 0., 0., 0., 0., 1., 1., 0., 0., 1., 1., 1., 1., 0., 1., 1., 0., 0., 0., 0., 1., 1., 0., 0.]) The data is now split into two separate sets - train and test sets respectively. We are only left with one final preprocessing step of scaling before we can fit a machine learning model to the data. Now, let's try to understand what these scaled values mean in the real world. Let's use CreditScore as an example. The credit score of a person is their creditworthiness based on their credit history. The higher this number, the more financially trustworthy a person is considered to be. So, a CreditScore of 1 is the highest since we're rescaling all the values to the range of 0-1. In [21]: # Import StandardScaler from sklearn.preprocessing import StandardScaler # Instantiate StandardScaler and use it to rescale X_train and X_test scaler = StandardScaler() X_train = scaler.fit_transform(X_train) X_test = scaler.fit_transform(X_test) In [22]: # Import LogisticRegression from sklearn.linear_model import LogisticRegression # Instantiate a LogisticRegression classifier with default parameter values logreg = LogisticRegression() # Fit logreg to the train set logreg.fit(X_train, y_train) C:\Users\Pranav.LAPTOP-HOVCQVL6\anaconda3\lib\importlib_bootstrap.py:219: RuntimeWarning: numpy.ufun c size changed, may indicate binary incompatibility. Expected 192 from C header, got 216 from PyObjec return f(*args, **kwds) Out[22]: LogisticRegression(C=1.0, class weight=None, dual=False, fit intercept=True, intercept_scaling=1, l1_ratio=None, max_iter=100, multi class='auto', n jobs=None, penalty='12', random_state=None, solver='lbfgs', tol=0.0001, verbose=0, warm start=False) 6. Making predictions and evaluating performance We will now evaluate our model on the test set with respect to classification accuracy. But we will also take a look the model's confusion matrix. In the case of predicting credit card applications, it is equally important to see if our machine learning model is able to predict the approval status of the applications as denied that originally got denied. In [23]: # Import confusion matrix from sklearn.metrics import confusion_matrix # Use logreg to predict instances from the test set and store it y_pred = logreg.predict(X_test) # Get the accuracy score of logreg model and print it print("Accuracy of logistic regression classifier: ", logreg.score(X test, y test)) # Print the confusion matrix of the logreg model confusion_matrix(y_test, y_pred) Accuracy of logistic regression classifier: 0.8743961352657005 Out[23]: array([[81, 4], [22, 100]], dtype=int64) Our model was able to yield an accuracy score of almost 88% In []: