

```
In [ ]: #standard imports
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import matplotlib inline

1.Get The Data Ready
```

```
In [2]: #import dataset

loan_default = pd.read_csv('Loan_Default.csv')

#view the data
loan_default.head()
```

	id	year	loan_limit	Gender	approv_in_adv	loan_type	loan_purpose	Credit_Worthiness	open_credit	business_or_commercial	...	credit_type	Credit_Score	applicant_credit_type	co-age	submission_of_application	LTV	Region	
	0	24890	2019	cf	Sex Not Available	ngpre	type1	p1	li	nopc		blvc	---	EXP	758	CIB	25-34	to_inst 98.72814	sou
	1	24891	2019	cf	Male	ngpre	type2	p1	li	nopc		b/c	---	EQUI	552	CIB	55-64	to_inst 98.01985	sou
	2	24892	2019	cf	Male	pre	type1	p1	li	nopc		nb/cv	---	EXP	834	CIB	35-44	not_inst 69.37690	Non
	3	24893	2019	cf	Male	ngpre	type1	p4	li	nopc		nb/cv	---	EXP	587	CIB	45-54	not_inst 69.37690	Non
	4	24894	2019	cf	Joint	pre	type1	p1	li	nopc		nb/cv	---	CRIF	602	EXP	25-34	not_inst 91.88654	Non

5 rows × 34 columns

```
loan_default.shape
```

```
(148670, 34)
```

There are 148,670 Rows and 34 Columns(features) in the dataframe.

The Features are:

- ID = Customer ID of Applicant
- year = Year of Application
- loan_limit = maximum available amount of the loan allowed to be taken
- Gender = sex type

```
In [4]: loan_default.shape

(148678, 34)
```

Out[4]: (148678, 34)

There are 148,670 Rows and 34 Columns(features) in the dataframe.

The Features are;

ID = Customer ID of Applicant

year = Year of Application

loan limit = maximum available amount of the loan allowed to be taken

Gender = sex type

approv_in_adv = is loan pre-approved or not

loan_type = Type of loan

loan_purpose = the reason you want to borrow money

Credit_Worthiness = is how a lender determines that you will default on your debt obligations, or how worthy you are to receive new credit.

open_credit = is a pre-approved loan between a lender and a borrower. It allows the borrower to make repeated withdrawals up to a certain limit.

business_or_commercial = Usage type of the loan amount

loan_amount = The exact loan amount

rate_of_interest = is the amount a lender charges a borrower and is a percentage of the principal—the amount loaned.

Interest_rate_spread = the difference between the interest rate a financial institution pays to depositors and the interest rate it receives from loans

Upfront_charges = Fee paid to a lender by a borrower as consideration for making a new loan

term = the loan's repayment period

Neg_amortization = refers to a situation when a loan borrower makes a payment less than the standard installment set by the bank.

interest_only = amount of interest only without principles

lump_sum_payment = is an amount of money that is paid in one single payment rather than in installments.

property_value = the present worth of future benefits arising from the ownership of the property

construction_type = Collateral construction type

occupancy_type = classifications refer to categorizing structures based on their usage

Secured_by = Type of Collateral

total_units = number of units

income = refers to the amount of money, property, and other transfers of value received over a set period of time

credit_type = type of credit

co-applicant_credit_type = is an additional person involved in the loan application process. Both applicant and co-applicant apply and sign for the loan

age = applicant's age

submission_of_application = Ensure the application is complete or not

LTV = life-time value (LTV) is a prognostication of the net profit

Region = applicant's place

Security_Type = Type of Collateral

status = Loan status (Approved/Declined)

dtir1 = debt-to-income ratio

```
In [17]: #The dataframe is not displaying all of the 34 features.
#To solve this we set number of columns we want to display with pandas settings

pd.set_option("display.max_columns", loan_default.shape[-1])
loan_default.head()
```

	ID	year	loan_limit	Gender	approv_in_adv	loan_type	loan_purpose	Credit_Worthiness	open_credit	business_or_commercial	loan_amount	rate_of_interest	Interest_rate_spread	Upfront_charges	term	Neg_amortization	interest_only
	0	24890	2019	cf	Sex Not Available	ngpre	type1	p1	li	nopc		nb/cv	116500		NaN		NaN
	1	24891	2019	cf	Male	ngpre	type2	p1	li	nopc		b/c	206500		NaN		NaN
	2	24892	2019	cf	Male	pre	type1	p1	li	nopc		nb/cv	405500	4.56	0.2000	595.0	360.0
	3	24893	2019	cf	Male	ngpre	type1	p4	li	nopc		nb/cv	465600	4.25	0.6810	1000.0	360.0
	4	24894	2019	cf	Joint	pre	type1	p1	li	nopc		nb/cv	696500	4.00	0.3042	0.0	360.0

```
In [18]: loan_default.columns
```

Out[18]: Index(['ID', 'year', 'loan_limit', 'gender', 'approv_in_adv', 'loan_type', 'loan_purpose', 'Credit_Worthiness', 'open_credit', 'business_or_commercial', 'loan_amount', 'rate_of_interest', 'Interest_rate_spread', 'Upfront_charges', 'term', 'Neg_amortization', 'interest_only', 'lump_sum_payment', 'property_value', 'construction_type', 'occupancy_type', 'secured_by', 'total_units', 'income', 'credit_type', 'Credit_Score', 'co-applicant_credit_type', 'age', 'submission_of_application', 'LTV', 'Region', 'Security_Type', 'Status', 'dtir1'], dtype='object')

```
In [19]: #converting all features to lower case for easy accessibility.

loan_default.columns = loan_default.columns.str.lower()
loan_default.columns
```

Out[19]: Index(['id', 'year', 'loan_limit', 'gender', 'approv_in_adv', 'loan_type', 'loan_purpose', 'credit_worthiness', 'open_credit', 'business_or_commercial', 'loan_amount', 'rate_of_interest', 'interest_rate_spread', 'upfront_charges', 'term', 'neg_amortization', 'interest_only', 'lump_sum_payment', 'property_value', 'construction_type', 'occupancy_type', 'secured_by', 'total_units', 'income', 'credit_type', 'credit_score', 'co_applicant_credit_type', 'age', 'submission_of_application', 'ltv', 'region', 'security_type', 'status', 'dtir1'], dtype='object')

filtering the data, since there are large number of unuseful columns which can be dropped without a bad effect on the performance of the model.

```
In [20]: #We can drop 'id', 'loan_limit', 'gender', 'approv_in_adv', 'loan_purpose', 'open_credit', 'business_or_commercial' because they don't add any valuable information.
#We can drop 'year' as it includes only one year 2019 so we have the same year.

columns = ['id', 'year', 'loan_limit', 'gender', 'approv_in_adv', 'loan_purpose', 'credit_worthiness', 'open_credit', 'business_or_commercial', 'interest_rate_spread', 'upfront_charges', 'term', 'neg_amortization', 'interest_only', 'lump_sum_payment', 'construction_type', 'occupancy_type', 'secured_by', 'total_units', 'credit_type', 'co_applicant_credit_type', 'submission_of_application', 'ltv', 'region', 'security_type']
loan_default.drop(columns,axis=1, inplace=True)
```

```
In [21]: loan_default.columns
```

Out[21]: Index(['loan_type', 'loan_amount', 'rate_of_interest', 'term', 'property_value', 'income', 'credit_score', 'age', 'status', 'dtir1'], dtype='object')

```
In [22]: loan_default.shape

(148678, 10)
```

```
In [23]: pd.set_option("display.max_columns", loan_default.shape[-1])
loan_default.head()
```

```
loan_default.shape
```

```
(148664, 10)
```

```
CHECKING FOR NAN VALUES
```

```
loan_default.isna().sum()
```

CHECKING FOR DUPLICATES

```
In [25]: loan_default.duplicated().sum()
```

Out[25]: 6

```
In [26]: loan_default.drop_duplicates(inplace=True)
```

```
In [27]: loan_default.duplicated().sum()
```

Out[27]: 0

```
In [28]: loan_default.shape

(148664, 10)
```

CHECKING FOR NAN VALUES

```
In [30]: loan_default.isna().sum()
```

Out[30]: loan_type 0
loan_amount 0
rate_of_interest 36437
term 41
property_value 15896
income 9146
credit_score 209
age 8
status 0
dtir1 24115
dtype: int64

```
In [32]: loan_default.info()
```

```
<class 'pandas.core.frame.DataFrame'>
Int64Index: 148664 entries, 0 to 148669
Data columns (total 10 columns):
# Column Non-Null Count Dtype
---
0 loan_type 148664 non-null object
1 loan_amount 148664 non-null int64
2 rate_of_interest 12227 non-null float64
3 term 148623 non-null float64
4 property_value 133568 non-null float64
5 income 239513 non-null float64
6 credit_score 148664 non-null int64
7 age 148664 non-null object
8 status 148664 non-null int64
9 dtir1 124549 non-null float64
dtypes: float64(5), int64(3), object(2)
memory usage: 12.5+ MB
```

```
In [33]: #filling the null in columns with int and floats.
from sklearn.impute import SimpleImputer
imputer = SimpleImputer()
loan_default[['rate_of_interest', 'term', 'property_value', 'income', 'dtir1']] = imputer.fit_transform(loan_default[['rate_of_interest', 'term', 'property_value', 'income', 'dtir1']])
```

```
In [34]: #filling the null in columns with objects.
from sklearn.impute import SimpleImputer
imputer = SimpleImputer(strategy='most_frequent')
loan_default[['age']] = imputer.fit_transform(loan_default[['age']])
```

```
In [35]: loan_default.isna().sum()
```

Out[35]: loan_type 0
loan_amount 0
rate_of_interest 0
term 0
property_value 0
income 0
credit_score 0
age 0
status 0
dtir1 0
dtype: int64

```
In [36]: #Dealing with non numerical values
#All data must be numerical to be used in the machine learning model

loan_default.info()
```

```
<class 'pandas.core.frame.DataFrame'>
Int64Index: 148664 entries, 0 to 148669
Data columns (total 16 columns):
# Column Non-Null Count Dtype
---
0 loan_type 148664 non-null object
1 loan_amount 148664 non-null int64
2 rate_of_interest 148664 non-null float64
3 term 148664 non-null float64
4 property_value 148664 non-null float64
5 income 148664 non-null float64
6 credit_score 148664 non-null int64
7 age 148664 non-null object
8 status 148664 non-null int64
9 dtir1 148664 non-null float64
dtypes: float64(5), int64(3), object(2)
memory usage: 12.5+ MB
```

```
In [37]: loan_default = pd.get_dummies(loan_default, columns=['loan_type', 'age'], drop_first=True)
```

```
In [38]: loan_default.info()
```

```
<class 'pandas.core.frame.DataFrame'>
Int64Index: 148664 entries, 0 to 148669
Data columns (total 16 columns):
# Column Non-Null Count Dtype
---
0 loan_type 148664 non-null object
1 loan_amount 148664 non-null int64
2 rate_of_interest 148664 non-null float64
3 term 148664 non-null float64
4 property_value 148664 non-null float64
5 income 148664 non-null float64
6 credit_score 148664 non-null int64
7 age_35-44 148664 non-null uint8
8 age_45-54 148664 non-null uint8
9 age_55-64 148664 non-null uint8
10 age_65-74 148664 non-null uint8
11 age_75+ 148664 non-null uint8
dtypes: float64(5), int64(3), uint8(8)
memory usage: 11.3 MB
```

```
In [39]: pd.set_option("display.max_columns", loan_default.shape[-1])
loan_default.head()
```

	loan_amount	rate_of_interest	term	property_value	income	credit_score	dtir1	loan_type_type2	loan_type_type3	age_35-44	age_45-54	age_55-64	age_65-74	age_75+>	age_>74
0	116500	4.045482	360.0	118000.000000	1740.0	758	1 45.000000	0	0	0	0	0	0	0	0
1	206500	4.045482	360.0	497900.200647	4980.0	552	1 37.723932	1	0	0	0	0	1	0	0
2	405500	4.560000	360.0	508000.000000	9480.0	834	0 46.000000	0	0	1	0	0	0	0	0
3	465600	4.250000	360.0	696000.000000	11880.0	587	0 42.000000	0	0	0	1	0	0	0	0
4	696500	4.000000	360.0	758000.000000	10440.0	602	0 39.000000	0	0	0	0	0	0	0	0

148664 rows × 15 columns

```
In [44]: y
```

Out[44]: 0 1
1 1
2 0
3 0
4 0
148665 0
148666 0
148667 0
148668 0
148669 0
Name: status, Length: 148664, dtype: int64

```
In [45]: x.shape
```

Out[45]: (148664, 15)

```
In [46]: y.shape
```

Out[46]: (148664,)

```
In [57]: # Split the data into training and test sets
from sklearn.model_selection import train_test_split
x_train, x_test, y_train, y_test = train_test_split(x, y, test_size=0.3)
```

View the data shapes

```
In [58]: x_train
```

	loan_amount	rate_of_interest	term	property_value	income	credit_score	dtir1	loan_type_type2	loan_type_type3	age_35-44	age_45-54	age_55-64	age_65-74	age_75+>	age_>74
59491	665500	3.990000	360.0	1580000.000000	9060.000000	774	29.000000	0	0	0	1	0	0	0	0
105728	165600	3.750000	180.0	238000.000000	4380.000000	651	44.000000	0	0	0	0	0	0	0	0
71354	136500	3.750000	360.0	2160.000000	2160.000000	775	30.000000	0	0	0	0	1	0	0	0
60811	226500	3.500000	180.0	378000.000000	5160.000000	683	40.000000	0	0	0	0	0	0	0	1
28296	406500	4.250000	360.0	638000.000000	6240.000000	510	47.000000	0	0	0	0	1	0	0	0

...
101798	565500	3.990000	360.0	628000.000000	9180.000000	519	41.000000	0	0	0	0	1	0	0	0
49440	136500	3.990000	360.0	208000.000000	3360.000000	831	45.000000	0	0	0	0	1	0	0	0
34866	276500	4.345000	360.0	138000.000000	6960.000000	569	15.000000	0	0	0	0	0	0	0	0
148668	196500	3.500000	180.0	278000.000000	7140.0	737	29.000000	0	0	0	0	1	0	0	0
148669	406500	4.375000	240.0	558000.000000	7260.0	830	44.000000	0	0	0	1	0	0	0	0

104064 rows × 15 columns

```
In [59]: x_train.shape
```

Out[59]: (104064, 15)

```
In [60]: x_test
```

	loan_amount	rate_of_interest	term	property_value	income	credit_score	dtir1	loan_type_type2	loan_type_type3	age_35-44	age_45-54	age_55-64	age_65-74	age_75+>	age_>74
142373	456500	3.875000	360.0	1758000.000000	7800.000000	774	29.000000	0	0	1	0	0	0	0	0
72668	265500	4.375000	360.0	4990000.000000	9060.000000	513	28.000000	0	0	1	0	0	0	0	0
23085	516500	4.375000	360.0	138000.000000	6960.000000	763	41.000000	0	0	0	0	0	1	0	0
77959	476500	3.375000	180.0	4980000.000000	20340.000000	894	28.000000	0	0	0	1	0	0	0	0
86151	236500	4.045482	180.0	4379002.000000	6900.000000	656	37.723932	0	0	0	0	1	0	0	0

...
74666	426500	3.625000	360.0	4990000.000000	9480.000000	515	39.000000	0	0	0	0	1	0	0	0
105323	726500	3.500000	360.0	1259000.000000	6490.000000	611	40.000000	0	0	0	0	1	0	0	0
1706	236500	4.500000	360.0	3490000.000000	4560.000000	774	59.000000	1	0	0	0	0	1	0	0
136778	546500	4.875000	360.0	4290000.000000	6380.000000	707	38.000000	0	1	0	0	1	0	0	0
113676	546500	3.000000	360.0	6780000.000000	6957.392164	737	37.723932	0	1	0	0	1	0	0	0

44600 rows × 15 columns

```
In [61]: x_test.shape
```

Out[61]: (44600, 15)

```
In [62]: y_train
```

Out[62]: 59491 0
105728 0
71354 0
60811 0
28296 0
101798 0
49440 0
34866 0
148668 1
148669 1
Name: status, Length: 104064, dtype: int64

```
In [63]: y_test.shape
```

Out[63]: (104064,)

```
In [64]: y_test
```

Out[64]: 142373 0
72668 0
23085 0
77959 0
86151 1
74666 0
105323 0
1706 0
136778 0
113676 0
Name: status, Length: 44600, dtype: int64

```
In [65]: y_test.shape
```

Out[65]: (44600,)

```
In [66]: from sklearn.ensemble import RandomForestClassifier
model = RandomForestClassifier()
```

1. Fit the model to the data and use it to make a prediction

A model will (attempt to) learn the patterns in a dataset by calling the fit() function on it and passing it the data.

```
In [68]:
```