

**Project: - Data Science**

**Group Members: -**

Name	Roll No	GR.No.
Sanket Sonje	322053	17u176
Shardul Pathak	322057	17u666

**Division: - COMP B**

**COURSE FACULTY: Prof. Leena A. Deshpande.**

**Title: - Loan Prediction System**

**1. Problem: -**

- Suppose some company wants to automate the loan eligibility process (real time) based on customer detail provided while filling online application form. These details are Gender, Marital Status, Education, Number of Dependents, Income, Loan Amount, Credit History and others.
- To automate this process, they have given a problem to identify the customers segments, those are eligible for loan amount so that they can specifically target these customers. Here they have provided a data set.

**2. Objectives: -**

- Main Objective: -
  - Automate the Loan Eligibility Process (Real Time).
- Sub-Objectives: -
  - Extract the data from Datahack as a Resource.
  - Analyse the data and fit into the classification model.
  - Perform and Evaluate the model and plot the outcome of the project graphically.

### 3. Data (Variable Description): -

<u>SR.NO.</u>	<u>Variables</u>	<u>Description</u>
1	Loan_ID	Unique Loan ID
2	Gender	Male/ Female
3	Married Applicant	married (Y/N)
4	Dependents	Number of dependents
5	Education	Applicant Education (Graduate/ Under Graduate)
6	Self_Employed	Self-employed (Y/N)
7	ApplicantIncome	Applicant income
8	CoapplicantIncome	Co-applicant income
9	LoanAmount	Loan amount in thousands
10	Loan_Amount_Term	Term of loan in months
11	Credit_History	credit history meets guidelines
12	Property_Area	Urban/ Semi Urban/ Rural
13	Loan_Status	Loan approved (Y/N)

- Rows: 615
- Source: Datahack

### 4. Analyse problem and convert it into data: -

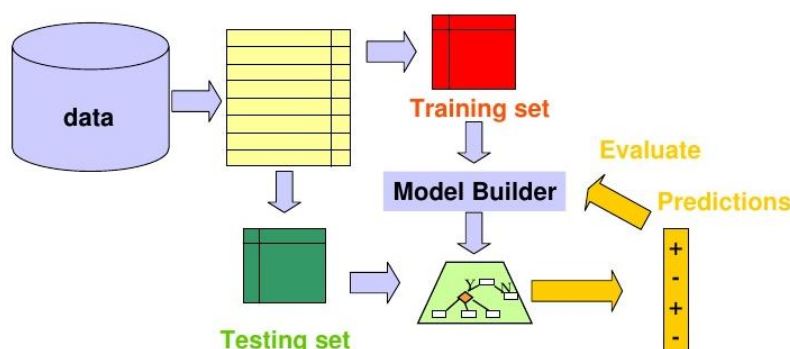
a) How data collected: -

- Datahack (Resource)

b) How data stored:

- Data is stored in the format of csv file.
  1. Test.csv
  2. Train.csv

### 5. Flow of the Project: -



## **6. Flow of the Project: -**

- a) Importing all libraries and storing data sets.
- b) Understanding various features of the datasets.
- c) Understanding the Distribution of Non-Categorical values.
- d) Understanding the Distribution of Categorical values.
- e) Replacing the Null values of Self-Employed column with maximum occurrence of the data.
- f) Nullifying Outliers of Loan Amount and Amount Income.
- g) Data preparation and model building.
- h) Testing and Validating score.

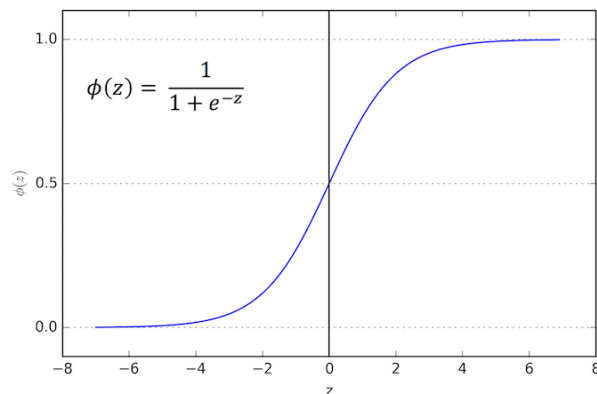
## **7. Logistic Regression: -**

- a) Logistic Regression is used when the dependent variable(target) is categorical.
- b) For example,
  - To predict whether an email is spam (1) or (0)
  - Whether the tumor is malignant (1) or not (0)
- c) In our example, Our Targeted Value is “Loan Status” and it has values Yes or No.
- d) To validate our Accuracy, we calculate the Cross-Validation Score.

## **8. Types of Logistic Regression: -**

- a) Binary Logistic Regression: -
  - The categorical response has only two 2 possible outcomes. Example: Spam or Not. We use this type.
- b) Multinomial Logistic Regression: -
  - Three or more categories without ordering. Example: Predicting which food is preferred more (Veg, Non-Veg, Vegan)
- c) Ordinal Logistic Regression: -
  - Three or more categories with ordering. Example: Movie rating from 1 to 5

## 9. Logistic Regression Graph: -



## 10. Analyze problem and convert it into data: -

- a) How data represented using statistical model: -
  - Data distribution: Binary Logistic Regression Model.
- b) Statistical operations performed: -
  - Min (), Avg (), describe (), Mode (), etc.
- c) Dealing with missing values: -
  - NAN, zero values /variable conversion.

## 11. Model building and training: -

- a) Model used: -
  - Logistic Regression Model
  - Binary output is expected.
- b) Libraries to be used: -
  - SKLEAN, SCIKIT, PANDAS, NUMPY, etc.

## 12. Code: - Importing Libraries: -

```
1 # Importing Library
2 import pandas as pd
3 import numpy as np
4 from sklearn import preprocessing
5 from sklearn.preprocessing import LabelEncoder
6
7 # Reading the training dataset in a dataframe using Pandas
8 df = pd.read_csv("train.csv")
9
10 # Reading the test dataset in a dataframe using Pandas
11 test = pd.read_csv("test.csv")
```

```
1 # First 20 Rows of training Dataset
2 df.head(20)
```

```
1 # Store total number of observation in training dataset
2 df_length = len(df)
3
4 # Store total number of columns in testing data set
5 test_col = len(test.columns)
```

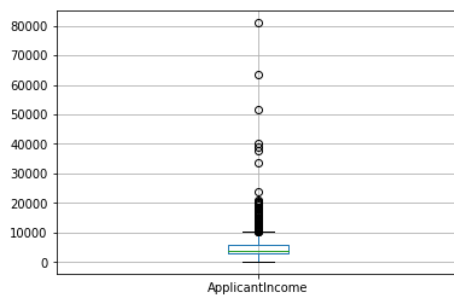
## 13. Understanding Distribution of Numerical Values: -

### a) ApplicationIncome: -

```
1 # Box Plot for understanding the distributions and to observe the outliers.
2
3 %matplotlib inline
4
5 # Histogram of variable ApplicantIncome
6
7 df['ApplicantIncome'].hist()
```

```
1 # Box Plot for variable ApplicantIncome of training data set
2
3 df.boxplot(column='ApplicantIncome')
```

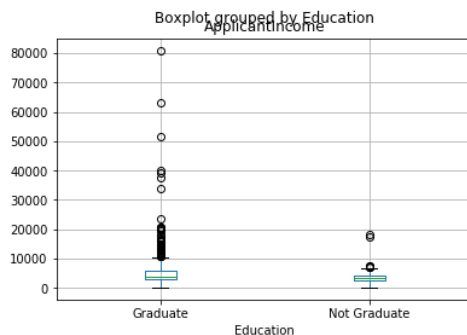
<matplotlib.axes.\_subplots.AxesSubplot at 0x1e8abc75240>



3. The above Box Plot confirms the presence of a lot of outliers/extreme values. This can be attributed to the income disparity in the society.

```
1 # Box Plot for variable ApplicantIncome by variable Education of training data set
2
3 df.boxplot(column='ApplicantIncome', by = 'Education')
```

<matplotlib.axes.\_subplots.AxesSubplot at 0x1e8aba1fc18>



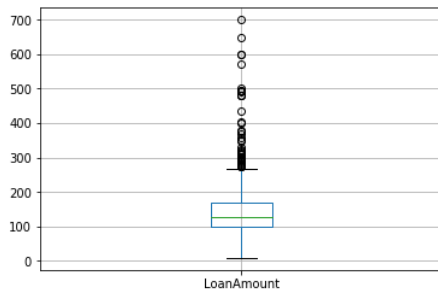
4. We can see that there is no substantial different between the mean income of graduate and non-graduates. But there are a higher number of graduates with very high incomes, which are appearing to be the outliers

## b) LoanAmount: -

```
1 # Histogram of variable LoanAmount
2
3 df['LoanAmount'].hist(bins=50)
```

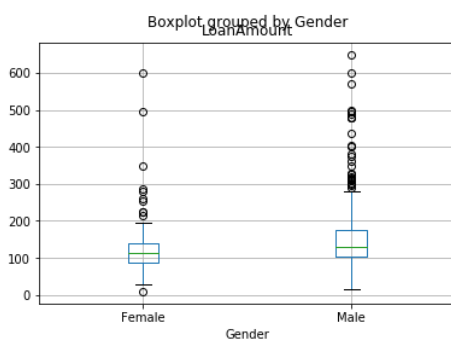
```
1 # Box Plot for variable LoanAmount of training data set
2
3 df.boxplot(column='LoanAmount')
```

<matplotlib.axes.\_subplots.AxesSubplot at 0x1e8abd4f7b8>



```
1 # Box Plot for variable LoanAmount by variable Gender of training data set
2
3 df.boxplot(column='LoanAmount', by = 'Gender')
```

<matplotlib.axes.\_subplots.AxesSubplot at 0x1e8abd96f98>



## 14. Understanding Distribution of Categorical Variables: -

```
1 # Loan approval rates in absolute numbers
2 loan_approval = df['Loan_Status'].value_counts()['Y']
3 print(loan_approval)
```

422

- 422 number of loans were approved.

```
1 # Credit History and Loan Status
2 pd.crosstab(df ['Credit_History'], df ['Loan_Status'], margins=True)
```

Loan_Status	N	Y	All
Credit_History			
0.0	82	7	89
1.0	97	378	475
All	179	385	564

```
1 #Function to output percentage row wise in a cross table
2
3 def percentageConvert(ser):
4     return ser/float(ser[-1])
5
6 df=pd.crosstab(df['Credit_History'], df['Loan_Status'], margins=True).apply(percentageConvert, axis=1)
7 loan_approval_with_credit_1 = df['Y'][1]
8
9 print(loan_approval_with_credit_1*100)
```

79.57894736842105

- 79.58 % of the applicants whose loans were approved have Credit\_History equals to 1.
-

## 15. Data Preparation And Model Building: -

```
1 # Impute missing values for Gender
2 df['Gender'].fillna(df['Gender'].mode()[0],inplace=True)
3
4 # Impute missing values for Married
5 df['Married'].fillna(df['Married'].mode()[0],inplace=True)
6
7 # Impute missing values for Dependents
8 df['Dependents'].fillna(df['Dependents'].mode()[0],inplace=True)
9
10 # Impute missing values for Credit_History
11 df['Credit_History'].fillna(df['Credit_History'].mode()[0],inplace=True)
12
13 # Convert all non-numeric values to number
14 cat=['Gender','Married','Dependents','Education','Self_Employed','Credit_History','Property_Area']
15
16 df.dtypes
```

Loan_ID	object
Gender	object
Married	object
Dependents	object
Education	object
Self_Employed	object
ApplicantIncome	int64
CoapplicantIncome	float64
LoanAmount	float64
Loan_Amount_Term	float64
Credit_History	float64
Property_Area	object
Loan_Status	object
TotalIncome	float64
LoanAmount_log	float64
dtype:	object

---

```
1 #Coverting all non-numerical values to numbers:
2 for var in cat:
3     le = preprocessing.LabelEncoder()
4     df[var]=le.fit_transform(df[var].astype('str'))
5 df.dtypes
```

Loan_ID	object
Gender	int32
Married	int32
Dependents	int32
Education	int32
Self_Employed	int32
ApplicantIncome	int64
CoapplicantIncome	float64
LoanAmount	float64
Loan_Amount_Term	float64
Credit_History	int32
Property_Area	int32
Loan_Status	object
TotalIncome	float64
LoanAmount_log	float64
dtype:	object

---



```

1 #Combining both train and test dataset to work with missing values
2
3 #Create a flag for Train and Test Data set
4 df['Type']='Train'
5 test['Type']='Test'
6 fullData = pd.concat([df,test], axis=0)
7
8 #Look at the available missing values in the dataset
9 fullData.isnull().sum()

```

F:\Anaconda\lib\site-packages\ipykernel\_launcher.py:6: FutureWarning: Sorting because non-concatenation axis is not aligned. A future version of pandas will change to not sort by default.

To accept the future behavior, pass 'sort=False'.

To retain the current behavior and silence the warning, pass 'sort=True'.

```

ApplicantIncome      0
CoapplicantIncome     0
Credit_History       29
Dependents            10
Education             0
Gender               11
LoanAmount            27
LoanAmount_log       389
Loan_Amount_Term      20
Loan_ID              0
Loan_Status          367
Married              0
Property_Area         0
Self_Employed        23
TotalIncome          367
Type                 0
dtype: int64

```

```

1 #Identify categorical and continuous variables
2 ID_col = ['Loan_ID']
3 target_col = ["Loan_Status"]
4 cat_cols = ['Credit_History','Dependents','Gender','Married','Education','Property_Area','Self_Employed']

```

```

1 #Imputing Missing values with mean for continuous variable
2 fullData['LoanAmount'].fillna(fullData['LoanAmount'].mean(), inplace=True)
3 fullData['LoanAmount_log'].fillna(fullData['LoanAmount_log'].mean(), inplace=True)
4 fullData['Loan_Amount_Term'].fillna(fullData['Loan_Amount_Term'].mean(), inplace=True)
5 fullData['ApplicantIncome'].fillna(fullData['ApplicantIncome'].mean(), inplace=True)
6 fullData['CoapplicantIncome'].fillna(fullData['CoapplicantIncome'].mean(), inplace=True)
7
8 #Imputing Missing values with mode for categorical variables
9 fullData['Gender'].fillna(fullData['Gender'].mode()[0], inplace=True)
10 fullData['Married'].fillna(fullData['Married'].mode()[0], inplace=True)
11 fullData['Dependents'].fillna(fullData['Dependents'].mode()[0], inplace=True)
12 fullData['Loan_Amount_Term'].fillna(fullData['Loan_Amount_Term'].mode()[0], inplace=True)
13 fullData['Credit_History'].fillna(fullData['Credit_History'].mode()[0], inplace=True)

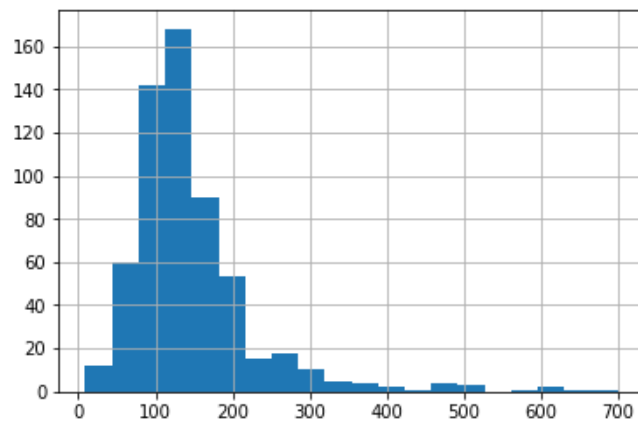
```

## 16. Nullify the Effect of Outliers: -

### a) Loan Amount: -

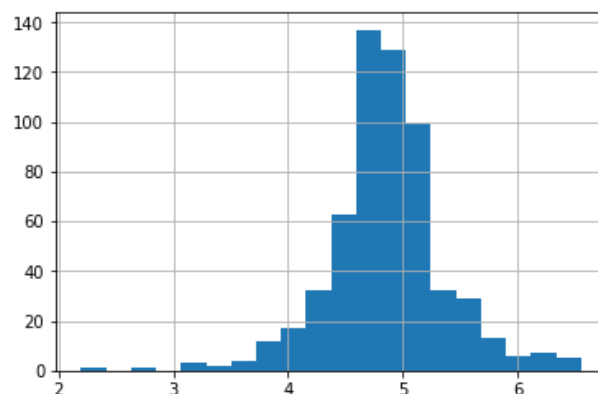
```
1 # Add both ApplicantIncome and CoapplicantIncome to TotalIncome
2 df['TotalIncome'] = df['ApplicantIncome'] + df['CoapplicantIncome']
3
4 # Looking at the distribution of TotalIncome
5 df['LoanAmount'].hist(bins=20)
```

<matplotlib.axes.\_subplots.AxesSubplot at 0x1e8abc5a550>



```
1 # Perform Log transformation of TotalIncome to make it closer to normal
2 df['LoanAmount_log'] = np.log(df['LoanAmount'])
3
4 # Looking at the distribution of TotalIncome_Log
5 df['LoanAmount_log'].hist(bins=20)
```

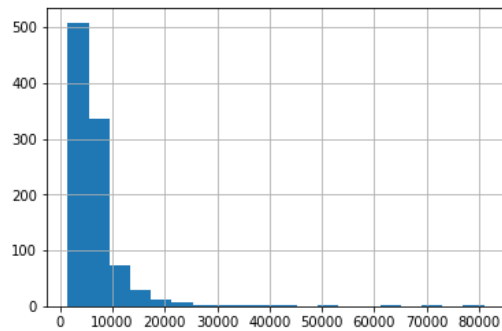
<matplotlib.axes.\_subplots.AxesSubplot at 0x1e8a88e72b0>



## b) Total Income: -

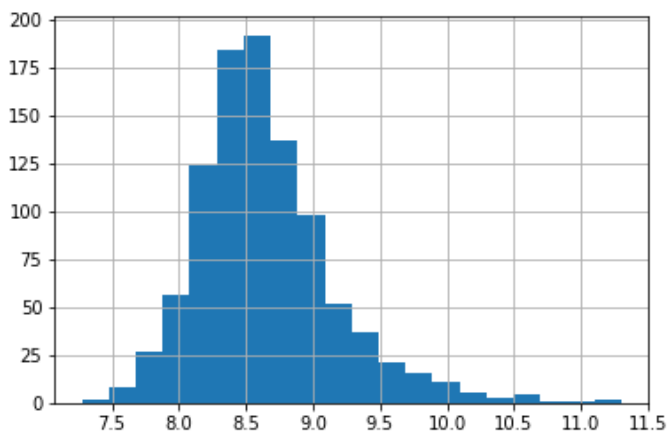
```
1 #Create a new column as Total Income
2
3 fullData['TotalIncome']=fullData['ApplicantIncome'] + fullData['CoapplicantIncome']
4
5 fullData['TotalIncome'].hist(bins=20)
```

<matplotlib.axes.\_subplots.AxesSubplot at 0x1e8abe92a90>



```
1 #Nullify the effect of outliers using log transformation
2 fullData['TotalIncome_log'] = np.log(fullData['TotalIncome'])
3
4 #Histogram for Total Income
5 fullData['TotalIncome_log'].hist(bins=20)
```

<matplotlib.axes.\_subplots.AxesSubplot at 0x1e8abeec2b0>



## 17. Data Preparation And Model Building: -

```
1 #create Label encoders for categorical features
2 for var in cat_cols:
3     number = LabelEncoder()
4     fullData[var] = number.fit_transform(fullData[var].astype('str'))
5
6 train_modified=fullData[fullData['Type']=='Train']
7 test_modified=fullData[fullData['Type']=='Test']
8 train_modified["Loan_Status"] = number.fit_transform(train_modified["Loan_Status"].astype('str'))
```

F:\Anaconda\lib\site-packages\ipykernel\_launcher.py:8: SettingWithCopyWarning:

A value is trying to be set on a copy of a slice from a DataFrame.

Try using .loc[row\_indexer,col\_indexer] = value instead

See the caveats in the documentation: <http://pandas.pydata.org/pandas-docs/stable/indexing.html#indexing-view-versus-copy>

## 18. Classification Model: -

```
1 #Import models from scikit learn module:
2 from sklearn import metrics
3 from sklearn.model_selection import cross_val_predict, KFold
4 from sklearn.model_selection import train_test_split
5
6 #Generic function for making a classification model and accessing performance:
7
8 def classification_model(model, data, predictors, outcome):
9     #Fit the model:
10    model.fit(data[predictors],data[outcome])
11
12    #Make predictions on training set:
13    predictions = model.predict(data[predictors])
14
15    #Print accuracy
16    accuracy = metrics.accuracy_score(predictions,data[outcome])
17    print ("Accuracy : %s" % "{0:.3%}".format(accuracy))
18
19    #Perform k-fold cross-validation with 5 folds
20    kf = KFold(n_splits=5)
21    error = []
22    for train, test in kf.split(data[predictors]):
23        #Filter training data
24        train_predictors = (data[predictors].iloc[train,:])
25
26        #The target we're using to train the algorithm.
27        train_target = data[outcome].iloc[train]
28
29        #Training the algorithm using the predictors and target.
30        model.fit(train_predictors, train_target)
31
32        #Record error from each cross-validation run
33        error.append(model.score(data[predictors].iloc[test:], data[outcome].iloc[test]))
34
35    print ("Cross-Validation Score : %s" % "{0:.3%}".format(accuracy))
36
37    #Fit the model again so that it can be refered outside the function:
38    #model.fit(data[predictors],data[outcome])
```

## 19. Logistic Regression Model and Output: -

```
1 from sklearn.linear_model import LogisticRegression
2
3
4 predictors_Logistic=['Credit_History','Education','Gender']
5
6 x_train = train_modified[list(predictors_Logistic)].values
7 y_train = train_modified["Loan_Status"].values
8
9 x_test=test_modified[list(predictors_Logistic)].values
10
11
12 # Create Logistic regression object
13 model = LogisticRegression()
14
15 # Train the model using the training sets
16 model.fit(x_train, y_train)
17
18 #Predict Output
19 predicted= model.predict(x_test)
20
21 #Reverse encoding for predicted outcome
22 predicted = number.inverse_transform(predicted)
23
24 #Store it to test dataset
25 test_modified['Loan_Status']=predicted
26
27 outcome_var = 'Loan_Status'
28
29 classification_model(model, df,predictors_Logistic,outcome_var)
30
31 test_modified.to_csv("Logistic_Prediction.csv",columns=['Loan_ID','Loan_Status'])
```

F:\Anaconda\lib\site-packages\sklearn\linear\_model\logistic.py:433: FutureWarning: Default solver will be changed to 'lbfgs' in 0.22. Specify a solver to silence this warning.  
FutureWarning)

Accuracy : 80.945%

Cross-Validation Score : 80.945%

## 20. Output: -

Unnamed: 0	Loan_ID	Loan_Status
0	0 LP001015	Y
1	1 LP001022	Y
2	2 LP001031	Y
3	3 LP001035	Y
4	4 LP001051	Y
5	5 LP001054	Y
6	6 LP001055	Y
7	7 LP001056	N
8	8 LP001059	Y
9	9 LP001067	Y
10	10 LP001078	Y
11	11 LP001082	Y
12	12 LP001083	Y
13	13 LP001094	N
14	14 LP001096	Y
15	15 LP001099	Y
16	16 LP001105	Y
17	17 LP001107	Y
18	18 LP001108	Y
19	19 LP001115	Y

## **21. Conclusion And Observations: -**

- a) The chances of getting a loan will be higher for:
  - Applicants having a credit history (we observed this in exploration.)
  - Applicants with higher applicant and co-applicant incomes.
  - Applicants with higher education level.
  - Properties in urban areas with high growth perspectives.

## **22. Analysing and interpreting the results: -**

- a) Test the result on different samples: -
  - Training data
  - Testing dataset
- b) Accuracy: - **80.945%**
- c) Cross – Validation Score: - **80.945%**