

Logistic Regression - Theory

What is Classification?

Predict Categories of dependent Variable (output)
Using independent (input) Variable(s).

↳ Will it rain tomorrow?

↳ Will the flight arrive in time?

Log. Regression :- Type of classification technique
whose underlying concept is based on
Linear Regression to predict Categories.

Eg:- given the past data, should we give
loan to this customer?

Use Cases of Logistic Regression

Finance \Rightarrow Can Credit be approved?

HealthCare \Rightarrow Do am I Diabetic?

Engineering \Rightarrow ^{will} The machine work after 6 months?

Benefits of Logistic Regression

- Dependent Variable doesn't need to be correlated with independent variables.
- No Normality assumption of dependent Variable
- There is no homogeneity of variance assumption.
- Normal distribution of errors is not assumed.
- Tree based & other algorithms may have higher accuracy, but Logistic Regression is used preferably in Probability Based Solution.

Why not use Linear Regression for Classification?

$\hat{y} = \beta_1 x + \beta_0$ - This equation has solution to the gm.

If we use Linear Regression the line incorrectly divides a single Category into two.

→ No Linear Regression can't Predict Value of Categorical Variable.

→ The line drawn doesn't capture any information about outcome variable.

Challenges with LR:-

↳ output becomes greater than 1 & less than 0.

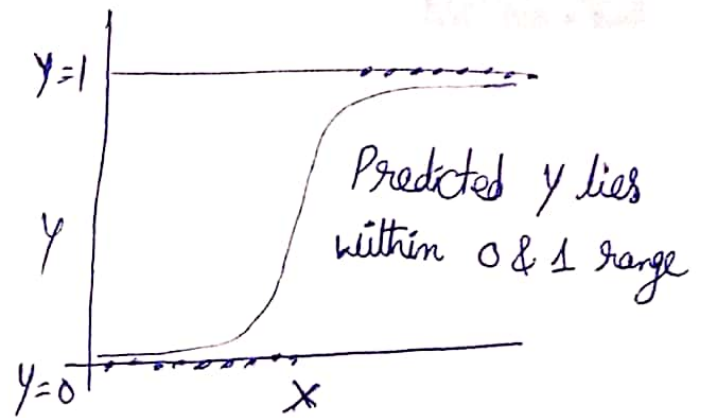
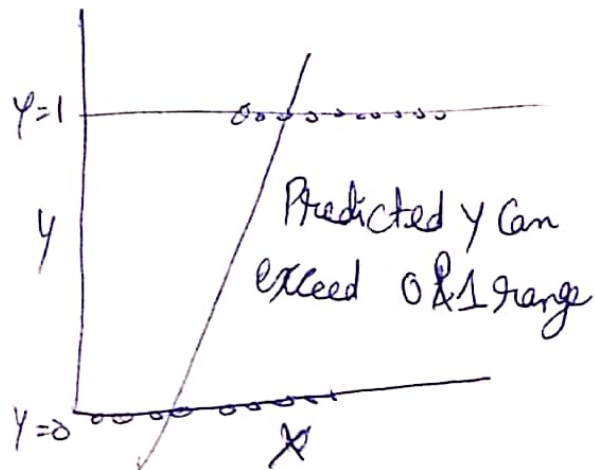
↳ This violates the rule of probability.

$$P(y) = \beta_0 + \beta_1 x$$

→ Logistic Regression helps us when output variable is categorical.

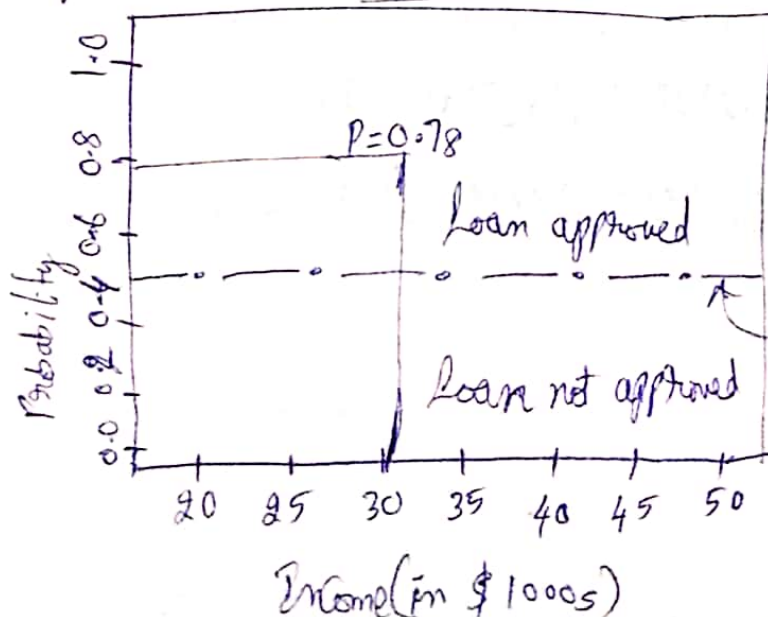
→ Transforms Linear Function to a function that

Can predict Categories with probabilities.



Logistic Regression always outputs a probability value to assign observations to classes.

Assigning observations to classes (Yes(1)/No(0)) depends on threshold probability.



1 (approved), because $0.78 > 0.5$
 $0.78 > 0.5$ (default threshold probability)

Odds Ratio

Logistic Regression can be written in terms of Odds Ratio

$$\text{Odds Ratio} = \frac{\text{Probability of event occurring}}{\text{Probability of event not occurring}}$$

$$\text{Odds Ratio} = \frac{P}{1-P}$$

applying odds ratio on loan approved example,

odds of loan being approved is

$$\frac{0.78}{0.22} = \underline{3.54}$$

Logit Function / Sigmoid Function

Transformation of linear function into non linear function that predicts categories.

$$\text{logit}(P) = \ln \frac{P}{1-P} = \beta_0 + \beta_1 x$$

Let's Break logit function into Probability

$$\ln \left(\frac{P(y)}{1-P(y)} \right) = \beta_0 + \beta_1 x$$

$$P(y) = \frac{e^{\beta_0 + \beta_1 x}}{1 + e^{\beta_0 + \beta_1 x}}$$

This is also called sigmoid function that is used to Predict probability to assign observations to Categories.

$$0 \leq \frac{e^{\beta_0 + \beta_1 x}}{1 + e^{\beta_0 + \beta_1 x}} < 1$$

↳ Equation above gives probability b/w 0 & 1

↳ As a result, it fulfills the Criteria of assigning observations to classes with probability.

Classification metrics

metrics used to evaluate Classification

Algo

Confusion matrix, Accuracy, Precision and Recall

Model	CATS	True Cat True Negative (TN) 25	False Cat False Negative (FN) 10
	Dogs	False dog False positive (FP) 5	True dog True positive (TP) 60
		← Precision measurement Cells →	
		Cats	Dogs

Recall measurement Cells

Reality \Rightarrow Total Animal = 100
Total Dogs = 70
Total Cats = 30

Model Prediction \Rightarrow objective \rightarrow identify Dogs

True positive = 60

\Rightarrow Dogs correctly identified

False Negative = 10

\Rightarrow Dogs not correctly identified

True Negative = 25

\Rightarrow Cats correctly identified.

False Positive = 5

\Rightarrow Cats not correctly identified.

Model Performance

To identify dogs

$$\text{Accuracy} = (TP + TN) / \text{Total} = 85\%$$

→ How many Cats and Dogs are correctly identified by model.

$$\text{Precision} = TP / (TP + FP) = 92.3\%$$

→ When model identifies an animal as Dog, how many times model is right.

$$\text{Recall} = TP / (TP + FN) = 85.7\%$$

→ out of all Dogs (Total dogs), how many dogs model is able to identify.

$$\text{F-1 score} = \frac{(\text{Precision} \times \text{Recall}) \times 2}{(\text{Precision} + \text{Recall})} \Rightarrow \text{mean of Precision \& Recall}$$

harmonic mean $\frac{2xy}{x+y}$

Close the F score to 1, the better the model

Logistic Regression

① Import packages (load the data)

② EDA \rightarrow Boxplot, Hist, Cor, Pairplot,

③ ~~EDA~~ Data Cleaning

\rightarrow missing values
 \rightarrow outlier treatment

Step 1 Substeps

④ Feature Engineering.

Converting Categorical to Numerical for Column

`default_dummies = pd.get_dummies(Cred-cred_df.default, prefix='default',
drop_first=True)`

`cred_df = pd.concat([cred_df, default_dummies], axis=1)`

`student_dummies = pd.get_dummies(cred_df.student, prefix='student',
drop_first=True)`

`cred_df = pd.concat([cred_df, student_dummies], axis=1)`

Label encoding is used for more than two-classes

Remember the Confusion matrix; very
very imp

Step = 1

~~Labels~~ & ~~Splitting~~ features & labels
Splitting/create

X = - - - - -

Y = - - - - -

Step = 2

Splitting the data into train & test

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, random_state
= 0)

Step 3 Creating Instance of the model

```
from sklearn.linear_model import LogisticRegression
```

from sklearn.metrics.

log reg = Logistic Regression ()

Step 4 Fitting the model

~~Logit = Logistic Regression~~

$$\text{log reg. fit}(x_{\text{train}}, y_{\text{train}})$$

Step 5 Predict the model

$$y_{\text{pred-test}} = \text{logreg.predict}(x_{\text{test}})$$

Step 6 = Evaluation of model

Print ("Test accuracy: ", metrics.accuracy_score
(y_test, y_pred_test))

Print ("Train accuracy = ", metrics.accuracy_score(~~Y_test~~
(Y_train, Y_pred_train)
after 6 months!

Creating Confusion matrix

Conf = metrics.confusion_matrix(y_test, y_pred_test)

~~Conf~~ Classification Report

cr = metrics.classification_report(y_test, y_pred_test)