

Binary Classification with Logistic Regression

TM Quest

Overview

Overview

What Will we Learn in This Module?

■ Logistic Regression

- What is **binary classification**?
- How to train a logistic regression model?
- What is the difference between **predicting classes** and **predicting probabilities**?

■ Evaluating the Model

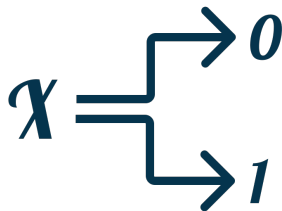
- How to evaluate a logistic regression model?
- What is **accuracy score**?

■ Bonus: What is an **Estimator** in Scikit-Learn?

Binary Classification & Logistic Regression

What is Binary Classification?

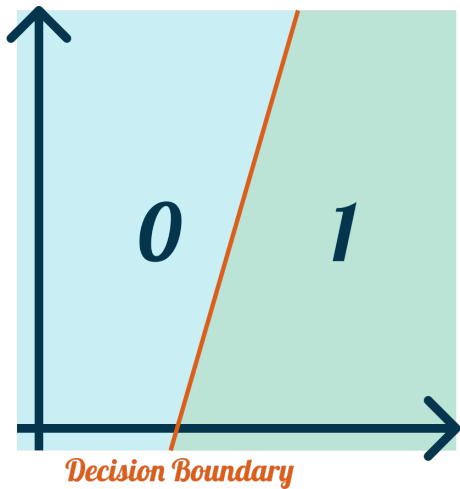
- The target values y are one of two categories.
- The categories are often encoded as 0 and 1.



Example

- Given internet browsing data as the features, predict the sex (male/female) of the user.
- Given historical stock data as the features, should you sell or keep your stock?
- Given data about a tumor as features, decide if it is malignant or benign.

Logistic Regression in the Feature Space



Logistic Regression Idea

In logistic regression one tries to find a **separating function** $p(x)$ such that:

- $p(x_0) \leq 0.5$, then x_0 should belong to category 0.
- $p(x_0) > 0.5$, then x_0 should belong to category 1.

The function $p(x)$ is (for a single feature) on the form

$$p(x) = \frac{1}{1 + e^{-(ax+b)}},$$

where training the model finds the best choices for a and b .

Logistic Regression Example

Example

Say that we have a single feature (size of tumor) and we want to predict whether the tumor is **malignant** or **benign**. If we find out (by training) that $a = 2$ and $b = -5$ then we can solve

$$p(x) = \frac{1}{1 + e^{2x-5}} \leq 0.5$$

and get $x \leq 2.5$. Hence if x (the size of the tumor) is less than 2.5, then we predict that the tumor is benign. If $x > 2.5$, then the tumor is malignant.

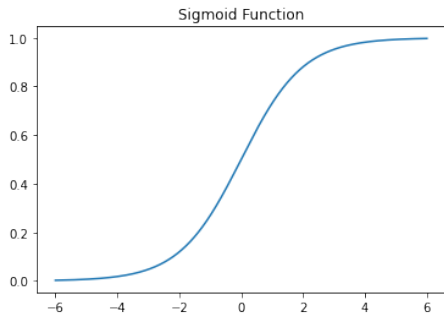
Takeaway

Training a logistic regression model finds the parameters a and b so that we can predict which category an observation is in.

Relationship with Linear Regression

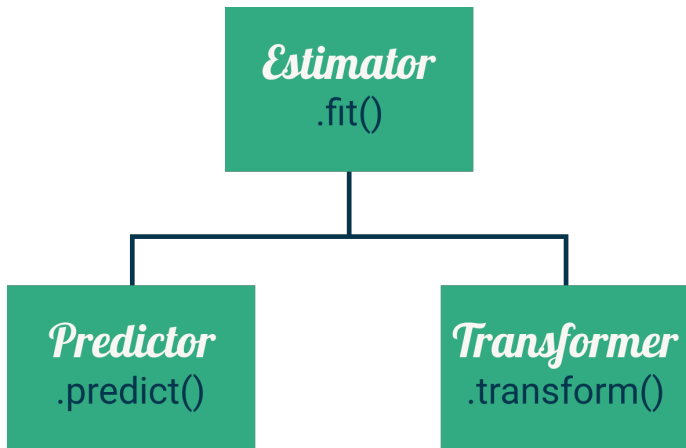
From a mathematical standpoint, the logistic regression function $p(x)$ is just linear regression $\mathbf{a} \cdot \mathbf{x} + b$ composed with the *sigmoid function*

$$\sigma(x) = \frac{1}{1 + e^{-x}}.$$



Estimators & Predictors

Estimators & Predictors



Accuracy Score

Accuracy Score

Definition

The **accuracy score** of a binary classifier is given by

$$\frac{\text{Number of Correctly Classified Observations}}{\text{Total Observations}}.$$

Example

If your model manages to guess 10 of 15 observations correctly, we get that the accuracy score is

$$\frac{10}{15} \approx 0.667 = 66.7\%.$$