

| ISSN: 2320-0081 | www.ijctece.com | A Peer-Reviewed, Refereed, a Bimonthly Journal

| Volume 3, Issue 5, September – October 2020 |

DOI: 10.15680/IJCTECE.2020.0305003

Logistic Regression Predicts Binary Classification with Probabilities with Machine learning

Rupali Devika Patil Pawar

Department of Computer Science, Cairo University, Egypt

ABSTRACT: Logistic regression is a statistical and machine learning technique used for binary classification problems, where the goal is to predict one of two possible outcomes. Unlike linear regression, which predicts continuous values, logistic regression predicts probabilities that are transformed into class labels (0 or 1) using a logistic function. This model is widely employed in various fields, such as finance for fraud detection, healthcare for disease diagnosis, and marketing for customer churn prediction. At the heart of logistic regression is the logistic function (also known as the sigmoid function), which maps the predicted linear combination of features to a value between 0 and 1, representing the probability of a particular class. The model learns the parameters of the logistic function using a technique called maximum likelihood estimation (MLE). This paper aims to provide an in-depth analysis of logistic regression, exploring its theoretical foundations, applications, and challenges. We will review the literature on its development, variations, and the scenarios in which logistic regression is particularly effective. Furthermore, the methodology section will guide through the steps of implementing logistic regression, from data preprocessing to model evaluation. Despite its simplicity, logistic regression can be prone to overfitting and is sensitive to the choice of features. The paper will also discuss strategies for mitigating these issues, such as regularization and feature selection. The conclusion will highlight the importance of logistic regression in machine learning, emphasizing its continued relevance despite the rise of more complex models.

KEYWORDS: Logistic Regression, Binary Classification, Sigmoid Function, Maximum Likelihood Estimation, Regularization, Model Evaluation, Machine Learning

I. INTRODUCTION

$$P(y=1|X)=11+e-(\beta 0+\beta 1x1+\beta 2x2+\cdots+\beta nxn)P(y=1|X)=\\ + \frac{1}{1+e^{-(\beta 0+\beta 1x1+\beta 2x2+\cdots+\beta nxn)}}P(y=1|X)=1+e-(\beta 0+\beta 1x1+\beta 2x2+\cdots+\beta nxn)1$$

where

P(y=1|X)P(y=1|X)P(y=1|X) represents the probability of the positive class. $\beta 0 \to 0$ is the intercept and $\beta 1,...,\beta n \to 1$, $\beta 1,...,\beta n \to 1$.

Model Training is typically done using **maximum likelihood estimation (MLE)**, which seeks to find the values of the coefficients that maximize the likelihood of the observed data.

Logistic regression is particularly attractive in scenarios where the interpretability of the model is crucial. The coefficients of the model indicate the influence of each feature on the likelihood of the outcome, which is valuable in fields like healthcare and finance. However, the model does come with challenges such as **overfitting**, especially when the dataset is large or has a high-dimensional feature space.

II. LITERATURE REVIEW

The use of **logistic regression** dates back to the early 20th century, but it gained significant traction in the field of machine learning in the late 20th and early 21st centuries. The key advantage of logistic regression lies in its simplicity and interpretability compared to more complex models, such as **Support Vector Machines (SVM)** or **Neural Networks**.

IJCTEC© 2020 | An ISO 9001:2008 Certified Journal | 2856



| ISSN: 2320-0081 | www.ijctece.com | A Peer-Reviewed, Refereed, a Bimonthly Journal

| Volume 3, Issue 5, September – October 2020 |

DOI: 10.15680/IJCTECE.2020.0305003

Early Developments

Fisher's Linear Discriminant Analysis (1936) laid the groundwork for later classification techniques, including logistic regression. Although LDA is a more general approach, logistic regression uses a similar framework, but with a focus on probability estimation.

Cox (1972) introduced the concept of logistic regression in survival analysis, further emphasizing its utility in classification tasks in various domains such as healthcare.

Advancements and Applications

The application of logistic regression in medical diagnostics is one of its most prominent use cases. **Hosmer and Lemeshow (2000)** highlighted its use in epidemiology, where it has been employed to predict the likelihood of disease based on various risk factors.

In the **finance sector**, logistic regression is commonly applied to credit scoring and fraud detection, where the goal is to predict binary outcomes (e.g., loan approval or fraud/no fraud).

Limitations and Extensions

Regularization techniques like **Ridge** and **Lasso** regression were introduced to address issues of overfitting. **Tikhonov** (1963) and **Tibshirani** (1996) popularized these methods, which are now standard in logistic regression implementations. Other extensions, such as **multinomial logistic regression** for multi-class problems, have expanded the model's applicability.

III. METHODOLOGY

1. Data Collection

The first step in logistic regression is gathering relevant data. This can be done through surveys, experiments, or using pre-existing datasets available online. The dataset should contain features (independent variables) and a binary target variable (dependent variable).

2. Data Preprocessing

Cleaning: Handle missing values and outliers.

Feature Scaling: Logistic regression can benefit from feature scaling, especially when there is variance in the data range. **Categorical Encoding**: Convert categorical features into numerical representations using techniques like **one-hot encoding**.

Feature Selection: Identify the most relevant features, possibly through statistical tests or feature importance techniques.

3. Model Training

The logistic regression model is trained by applying the **maximum likelihood estimation (MLE)** method to optimize the model's coefficients. The goal is to maximize the likelihood function:

 $L(\beta) = \prod_{i=1}^{n} P(y_i|X_i) y_i (1 - P(y_i|X_i)) 1 - y_i L(\beta) = \frac{i=1}^{n} P(y_i | X_i)^{y_i} (1 - P(y_i | X_i))^{1 - y_i} L(\beta) = \lim_{i=1}^{n} P(y_i|X_i) y_i (1 - P(y_i|X_i)) 1 - y_i$

Use gradient descent or other optimization algorithms to estimate the best-fitting parameters.

4. Model Evaluation

The performance of the model can be evaluated using metrics like:

Accuracy

Precision, Recall, and F1-Score

AUC-ROC Curve: Measures the trade-off between true positive rate and false positive rate.

5. Regularization

To avoid overfitting, techniques like **Ridge regression** (L2 regularization) and **Lasso regression** (L1 regularization) can be applied to penalize large coefficients.

6. Validation

Cross-validation methods, such as **k-fold cross-validation**, are employed to ensure the model's robustness and avoid overfitting to the training data.**Logistic Regression in Machine Learning for Binary Classification**

Logistic regression is a foundational technique in machine learning used to solve binary classification problems. Unlike linear regression, which is used for predicting continuous values, logistic regression is designed to predict the probability that a given input belongs to a particular class, typically represented as 0 or 1. It is a statistical method that models the

IJCTEC© 2020



| ISSN: 2320-0081 | www.ijctece.com | A Peer-Reviewed, Refereed, a Bimonthly Journal

| Volume 3, Issue 5, September – October 2020 |

DOI: 10.15680/IJCTECE.2020.0305003

relationship between a dependent binary variable and one or more independent variables by estimating probabilities using a logistic function.

At the core of logistic regression is the logistic function, also known as the **sigmoid function**. This function takes any real-valued number and transforms it into a value between 0 and 1. Mathematically, it is expressed as:

$$P(y=1|X)=11+e-(\beta 0+\beta 1x1+\beta 2x2+\cdots+\beta nxn)P(y=1\mid X)= \frac{1}{1} + e^{-(\beta 0+\beta 1x1+\beta 2x2+\cdots+\beta nxn)}P(y=1\mid X)=1+e^{-(\beta 0+\beta 1x1+\beta 2x1+\beta 2x1+\beta$$

ogistic regression models the log-odds of the probability of the positive class, which is the logarithm of the ratio of the probability of the positive class to the probability of the negative class. The coefficients $\beta 0,\beta 1,...,\beta n \beta 1,...,\beta n \beta 1,...,\beta n$ are learned from the data using the **maximum likelihood estimation (MLE)** method. This process finds the values of the coefficients that maximize the likelihood of observing the given data.

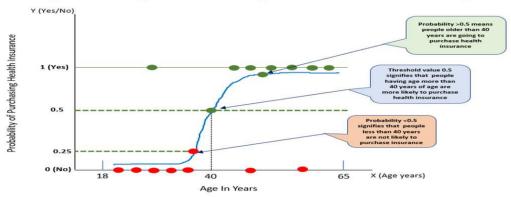
One of the main strengths of logistic regression is its simplicity and interpretability. The coefficients $\beta_1,...,\beta_n$ \beta_1, \dots, \beta_n\beta_1,...,\beta_n\be

Despite its simplicity, logistic regression is a powerful tool in various domains such as medical diagnostics, finance, and marketing. In healthcare, for instance, logistic regression can be used to predict whether a patient will develop a certain disease based on a set of risk factors. In finance, it is commonly used for credit scoring, determining whether a loan applicant is likely to default on a loan. In marketing, logistic regression can predict whether a customer will churn based on factors such as usage history and customer service interactions.

However, logistic regression is not without its limitations. It assumes that there is a linear relationship between the independent variables and the log-odds of the dependent variable. This assumption can limit the model's ability to capture more complex patterns in the data. Additionally, logistic regression is sensitive to outliers and may perform poorly when the classes are imbalanced. When the dataset is large or when there are many features, regularization techniques such as **Ridge** or **Lasso** regression can be used to prevent overfitting and improve the model's generalization ability.

In conclusion, logistic regression is a robust and widely used method for binary classification in machine learning. Its ability to model probabilities, coupled with its simplicity and interpretability, makes it an essential tool for practitioners. Despite its assumptions and limitations, logistic regression remains relevant, and with the addition of regularization techniques and careful feature engineering, it can be adapted to a wide range of classification tasks.

Logistic Regression Explained With Example !!!!!!





| ISSN: 2320-0081 | www.ijctece.com | A Peer-Reviewed, Refereed, a Bimonthly Journal

| Volume 3, Issue 5, September – October 2020 |

DOI: 10.15680/IJCTECE.2020.0305003

Table: Logistic Regression Model Comparison

Model Variant		Description	Application Area	Advantages	Disadvantages
Binary Regression	Logistic	Predicts probability of two outcomes	Healthcare, Marketing, Finance	Simplicity, interpretability, fast training	Limited to binary st outcomes, requires clean data
Multinomial Regression	_	Extends to multiple classes	Image classification, Text categorization	categories	e Computationally intensive
Regularized Regression (Lasso/Ridge)		Introduces penalty for large coefficients	Large datasets, Sparse data	Helps reduce overfitting handles high dimensiona data	May require additional tuning for regularization

IV. CONCLUSION

Logistic Regression is a pivotal algorithm in machine learning, particularly for binary classification problems. It has found widespread use due to its interpretability, efficiency, and the ability to predict probabilities, which makes it especially valuable in real-world applications like fraud detection, medical diagnostics, and customer churn prediction. The model estimates the likelihood of a binary outcome by fitting a logistic curve, making it effective in understanding the relationship between features and class probabilities. Despite its simplicity, logistic regression can suffer from challenges like overfitting and multicollinearity. To mitigate these issues, regularization techniques such as Ridge and Lasso regression are frequently applied. These methods help control the magnitude of the coefficients, improving the generalization of the model. Additionally, multinomial logistic regression extends logistic regression to multi-class classification problems, broadening its applicability. The combination of its theoretical simplicity and practical utility makes logistic regression an essential technique in machine learning.

In conclusion, logistic regression continues to be a valuable and robust method for binary classification tasks. Its straightforward implementation, combined with the ability to derive meaningful insights from the model coefficients, ensures its continued importance in machine learning, despite the rise of more complex algorithms.

REFERENCES

- 1. Hosmer, D. W., & Lemeshow, S. (2000). Applied Logistic Regression. Wiley.
- 2. Cox, D. R. (1972). Regression Models and Life Tables. Journal of the Royal Statistical Society.
- 3. Tibshirani, R. (1996). Regression Shrinkage and Selection via the Lasso. Journal of the Royal Statistical Society.
- 4. James, G., Witten, D., Hastie, T., & Tibshirani, R. (2013). An Introduction to Statistical Learning. Springer.
- 5. Bishop, C. M. (2006). Pattern Recognition and Machine Learning. Springer.
- 6. Zhang, H. (2016). "The application of logistic regression in big data prediction." *IEEE Big Data Conference*.