



# 15 QUESTIONS

## on “Logistic Regression” for Data Science and AI Interviews



# 01

## What is Logistic Regression?

### Explanation:

Logistic Regression is a statistical method to predict a binary outcome (e.g., Yes/No, Pass/Fail) based on one or more independent variables. It's used when the dependent variable is categorical with two classes.

### Detailed Explanation:

- **Logistic Regression** is a classification algorithm used for predicting the probability of a binary outcome. It doesn't predict exact values but gives the probability of an event occurring. For example, it can predict whether a student will pass or fail based on study hours, test preparation, etc.
- It's called logistic regression because it uses the logistic function (**Sigmoid function**) to convert predictions into probabilities.
- The output is a value between 0 and 1, which represents the probability that a given input point belongs to class 1.
- The decision threshold is usually set at 0.5, where values greater than 0.5 predict class 1, and values below predict class 0.

# 02

## How is Logistic Regression different from Linear Regression?

### Easy Explanation:

Linear Regression predicts continuous outcomes, while Logistic Regression predicts probabilities for binary outcomes.

### Details:

- **Linear Regression:** Used for predicting continuous values like house prices, stock prices, etc. The output is a real number, and it assumes a linear relationship between the input variables and the output.
- **Logistic Regression:** It is specifically designed for binary classification. Instead of producing a continuous value, it produces a probability value between 0 and 1 by applying a sigmoid function to the output of the linear equation.
- In Linear Regression, the prediction  $Y$  is calculated as  $Y = b_0 + b_1X_1 + \dots + b_nX_n$
- For Logistic Regression, the prediction  $P(Y=1)$  is obtained by applying the Sigmoid function to this result:

$$P(Y = 1) = \frac{1}{1 + e^{-(b_0 + b_1X_1 + \dots + b_nX_n)}}$$

- **Example:** If you're predicting salary based on years of experience, use Linear Regression. If you're predicting whether someone will get the job based on their experience, use Logistic Regression.

# 03

## What is the Sigmoid Function, and why is it used in Logistic Regression?

### Easy Explanation:

The sigmoid function squashes any real number into a range between 0 and 1, making it ideal for predicting probabilities, ensuring the output is interpretable as a probability.

### Details:

- The **Sigmoid function** is a mathematical function that maps any real-valued number into the range (0, 1). It is often represented as:

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

- where  $z$  is the linear combination of the input features (i.e.,  $b_0 + b_1X_1 + \dots + b_nX_n$ ).
- The function ensures that predictions are between 0 and 1, which makes them interpretable as probabilities. This is critical in classification tasks where the output needs to represent the likelihood of an event.
- Why use it? Without this transformation, the predicted output could be any real number, making it difficult to interpret the result as a probability.

**Example:** A prediction of 0.9 means there's a 90% chance of the event happening (e.g., a customer will buy a product).

# 04

## What is the cost function in Logistic Regression? Why isn't MSE used?

**Easy Explanation:** Logistic Regression uses a special cost function called Log Loss, which is better suited for classification problems. MSE is not ideal because it doesn't handle probabilities well.

### Details:

- In Logistic Regression, the cost function is designed to measure how well the model's predicted probabilities match the true labels. It is defined as Log Loss or Binary Cross-Entropy, which penalizes the model for incorrect predictions based on their confidence.

$$Cost = -\frac{1}{N} \sum_{i=1}^N [y_i \log(\hat{y}_i) + (1 - y_i) \log(1 - \hat{y}_i)]$$

- where  $y_i$  is the true label and  $\hat{y}_i$  is the predicted probability.
- Why not MSE? Mean Squared Error (MSE) works well for regression problems, where the target variable is continuous. However, for classification tasks, MSE is less effective because it doesn't capture the probabilistic nature of the output. Log Loss ensures better optimization by focusing on how confident the model is about each prediction.

**Example:** If the model predicts a probability of 0.95 for class 1 and the true label is also 1, the cost will be very low. If the prediction is 0.2 and the true label is 1, the cost will be high.

# 05

## What are odds and log-odds in Logistic Regression?

### Easy Explanation:

Odds measure the likelihood of an event happening, and log-odds transform odds into a scale that can be modeled with a linear function.

### Detailed Explanation:

- **Odds:** The odds are a ratio of the probability of an event happening to the probability of it not happening.

$$\text{Odds} = \frac{P(Y = 1)}{1 - P(Y = 1)}$$

- If the probability of an event occurring is 0.8, the odds are 4:1 (4 times more likely to occur than not).
- **Log-Odds (Logit):** Logistic Regression works with log-odds, which is the natural logarithm of the odds. This transformation makes the relationship between the predictors and the outcome linear.
  - $\text{Logit} = \ln(\text{Odds}) = b_0 + b_1X_1 + \dots + b_nX_n$
- This ensures that the prediction lies between 0 and 1 after applying the sigmoid function.

**Example:** If the odds are 3:1 (i.e., the event is 3 times more likely to happen than not), the log-odds would be  $\ln(3) \approx 1.1$

# 06

## What are some key assumptions of Logistic Regression?

**Easy Explanation:** Logistic Regression assumes certain conditions for accurate predictions, such as a linear relationship between the predictors and log-odds and independence among observations.

### Assumption Details:

- **Binary outcome:** The dependent variable should be binary (i.e., having only two possible outcomes).
- **No multicollinearity:** Independent variables should not be highly correlated with each other.
- **Linear relationship:** There should be a linear relationship between the independent variables and the log-odds of the dependent variable.
- **Independence:** Observations should be independent of each other (no autocorrelation).
- **Large sample size:** Logistic Regression works best with a large enough sample size to estimate coefficients reliably.

Violating these assumptions can lead to unreliable results, such as biased or inconsistent coefficients.

# 07

## How is model performance evaluated for Logistic Regression?

### Easy Explanation:

Performance metrics like accuracy, precision, recall, and AUC-ROC are used to assess how well the Logistic Regression model is doing.

### Detailed Explanation:

- **Accuracy:** Measures the overall correctness of the model. It's the ratio of correctly predicted instances to total instances.
- **Precision:** Measures how many of the predicted positive cases are actually positive.
  - $\text{Precision} = \frac{TP}{TP+FP}$
- **Recall:** Also known as sensitivity, it measures how many of the actual positive cases the model correctly identified.
  - $\text{Recall} = \frac{TP}{TP+FN}$
- **F1-Score:** Harmonic mean of precision and recall. It's useful when you need a balance between precision and recall.

**AUC-ROC Curve:** The Area Under the Receiver Operating Characteristic Curve (AUC-ROC) measures the ability of the model to discriminate between classes. An AUC of 0.5 indicates a random model, while 1 indicates perfect classification.



# 08

## What is the purpose of the threshold in Logistic Regression?

**Easy Explanation:** The threshold is used to convert the probability output of the logistic model into a class label (0 or 1).

### **Details:**

- Logistic Regression outputs a probability between 0 and 1. By default, a threshold of 0.5 is used: if the predicted probability is greater than or equal to 0.5, the prediction is class 1; otherwise, it is class 0.
- **Threshold tuning:** In some cases, adjusting the threshold can improve model performance, especially with imbalanced datasets. **For example, lowering the threshold can improve recall at the cost of precision.**

**Example:** If a model predicts a probability of 0.7 for a positive event, it would be classified as 1 (positive class) if the threshold is 0.5. If the threshold is set higher (e.g., 0.8), the event would be classified as 0 (negative class), even though the probability is higher than 0.5.

# 09

## How does regularization help in Logistic Regression?

### Easy Explanation:

Regularization helps prevent the model from overfitting by penalizing large coefficients, leading to a simpler and more generalizable model.

### Detailed Explanation:

- Regularization is a technique used to add a penalty term to the cost function in order to prevent overfitting, especially when the model becomes too complex with a large number of features.
- In Logistic Regression, two common types of regularization are used:
  - **L2 Regularization (Ridge Regression):** Adds the squared magnitude of the coefficients to the cost function. It encourages smaller coefficients by penalizing large values.
  - **L1 Regularization (Lasso Regression):** Adds the absolute values of the coefficients to the cost function. It encourages sparsity, meaning it may force some coefficients to become exactly zero, effectively performing feature selection.
- **Why use regularization?** Regularization helps prevent the model from memorizing the training data (overfitting), which leads to poor performance on unseen data. A regularized model is more robust and generalizes better to new data.

**Example:** In a model predicting loan approval, L2 regularization might reduce the influence of less important features (like minor details) on the outcome, improving the model's ability to generalize to new loan applicants.

# 10

## How does Logistic Regression handle multi-class classification?

**Easy Explanation:** Logistic Regression is originally designed for binary classification, but it can be extended to multi-class classification using methods like One-vs-Rest (OvR) or Softmax.

### Details:

- Binary Logistic Regression is designed to predict one of two classes (e.g., 0 or 1). However, in real-world problems, we often have more than two classes.
- Two common techniques are used to extend Logistic Regression to multi-class problems:
- One-vs-Rest (OvR) Approach: For  $k$  classes, it trains  $k$  separate binary classifiers, each distinguishing one class from all the others. The class with the highest predicted probability is selected as the final prediction.
- Example: If you're classifying animals into categories like Dog, Cat, and Rabbit, you would train three separate classifiers: one for Dog vs. Not-Dog, one for Cat vs. Not-Cat, and one for Rabbit vs. Not-Rabbit.
- Softmax Regression (Multinomial Logistic Regression): Instead of creating separate classifiers, it calculates probabilities for all classes simultaneously and selects the class with the highest probability.
- The Softmax function generalizes the sigmoid function for multi-class problems. It takes the raw output of the model (logits) and converts them into probabilities that sum to 1.

$$P(y = k|X) = \frac{e^{z_k}}{\sum_{i=1}^K e^{z_i}}$$

where  $z_k$  is the logit for class  $k$ , and  $K$  is the total number of classes.

**Example:** In a problem where you're predicting types of fruits (Apple, Banana, Orange), Softmax would give you the probability of each fruit class, and the one with the highest probability is chosen.

# 11

## What is the difference between the model coefficients and the odds ratios in Logistic Regression?

### Easy Explanation:

The model coefficients tell us how the predictors affect the log-odds of the outcome, while the odds ratios tell us how the odds of the outcome change with a unit increase in the predictor.

### Detailed Explanation:

- **Model Coefficients:** In Logistic Regression, the coefficients (also called weights) represent the change in the log-odds of the outcome for a one-unit change in the corresponding predictor variable.
- A positive coefficient means the outcome is more likely as the predictor increases. A negative coefficient means the outcome is less likely.
- **Example:** If the coefficient for "years of experience" is 0.5, it suggests that for every 1-year increase in experience, the log-odds of success increase by 0.5.
- **Odds Ratios:** The odds ratio (OR) is the exponential of the coefficient, which represents how the odds of the outcome change with a one-unit increase in the predictor variable.

$$\text{Odds Ratio (OR)} = e^{\beta}$$

- If the odds ratio is greater than 1, the predictor increases the odds of the event happening. If it is less than 1, the predictor decreases the odds.
- **Example:** If the odds ratio for "years of experience" is 1.6, this means that each additional year of experience increases the odds of success by 60%.
- **Summary:** While the coefficient represents the change in log-odds, the odds ratio provides a more interpretable change in the odds, which can be easier to explain to stakeholders.

# 12

## Can Logistic Regression be used for regression tasks?

**Easy Explanation:** Logistic Regression is specifically for classification tasks, but it is not suitable for regression problems where the output is continuous.

### Details:

- Logistic Regression is a classification algorithm, not a regression algorithm, because it predicts categorical outcomes (such as 0 or 1, Yes or No).
- **Why not for regression?** In regression problems, we need to predict continuous values, and techniques like Linear Regression are designed for that. Logistic Regression, on the other hand, is designed to predict probabilities that can then be classified as one of two possible outcomes.
- **Example:** If you're predicting someone's weight (a continuous value), you would use Linear Regression. However, if you're predicting whether a customer will churn (Yes or No), you would use Logistic Regression.
- **Note:** Logistic Regression is often confused with Linear Regression, but while both share similar names, their applications are very different. Logistic Regression is a classification method, while Linear Regression is a regression method for continuous values.

# 13

## How do we interpret the coefficients in Logistic Regression?

### Easy Explanation:

Coefficients tell us how each predictor affects the likelihood of the event. Positive coefficients increase the likelihood, and negative coefficients decrease it.

### Detailed Explanation:

- The coefficients ( $\beta_0, \beta_1, \dots$ ) in Logistic Regression describe the relationship between the predictor variables and the log-odds of the outcome.
- A positive coefficient means that as the predictor variable increases, the probability of the outcome occurring increases (i.e., the event becomes more likely).
- A negative coefficient means that as the predictor variable increases, the probability of the outcome occurring decreases (i.e., the event becomes less likely).

### Interpretation:

- If the coefficient for a predictor (e.g., years of education) is positive, it suggests that increasing education increases the likelihood of the event (e.g., getting hired).
- If the coefficient is negative, it indicates that as the predictor increases, the likelihood of the event decreases (e.g., the likelihood of a disease occurrence decreases with a certain treatment).

**Example:** For a logistic model predicting whether a customer will purchase a product (Yes/No), if the coefficient for the age variable is positive, it means older customers are more likely to buy the product.

# 14

## What are some limitations of Logistic Regression?

**Easy Explanation:** Logistic Regression assumes linearity and may struggle with complex relationships or highly correlated data. It also requires a large amount of data to make reliable predictions.

### Details:

- **Assumes Linearity:** Logistic Regression assumes that the relationship between the predictors and the log-odds of the outcome is linear. If this assumption is violated, the model may not perform well.
- **Multicollinearity:** When predictors are highly correlated with each other, Logistic Regression may produce unreliable coefficient estimates, which could lead to misleading results.
- **Non-Linearity of Features:** If there are complex interactions between variables or non-linear relationships, Logistic Regression may not capture these patterns well. Non-linear models like Decision Trees or Random Forests might be better suited in such cases.
- **Requires Large Sample Size:** Logistic Regression performs better with a larger dataset. With small datasets, the model may overfit, and the results may not be generalizable.
- **Imbalanced Datasets:** Logistic Regression may struggle with imbalanced classes. If one class is heavily underrepresented, the model might be biased towards predicting the majority class. Techniques like class weighting or resampling might help in such cases.

**Example:** In a dataset predicting customer churn, if 90% of customers are non-churners and only 10% churn, Logistic Regression may predict most customers as non-churners, leading to poor performance on the minority class.

# 15

## How do you handle missing values in Logistic Regression?

### Easy Explanation:

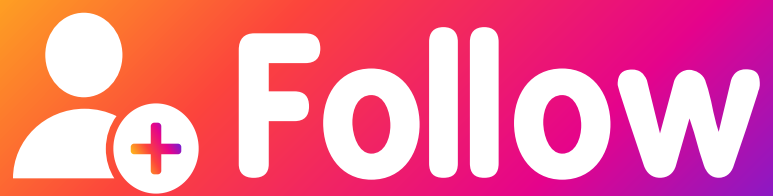
Missing values can be handled by removing rows with missing data, filling in missing values with the mean or median, or using more advanced techniques like imputation.

### Detailed Explanation:

- Handling missing values is a critical step in data preprocessing because Logistic Regression cannot handle missing data directly.
- **Simple Approaches:**
- **Remove rows with missing data:** This is useful if there are very few missing values and removing them won't significantly impact the dataset.
- **Imputation:** Replace missing values with the mean, median, or mode of the column. This is commonly done for numerical features.
- **Advanced Imputation Techniques:** Methods like K-Nearest Neighbors Imputation or Multiple Imputation can be used when the data has complex relationships, and you want a more sophisticated approach.
- **Why not leave missing data?** Logistic Regression requires a complete dataset. If missing values are not handled, the algorithm will fail to work.

**Example:** If the income column in a dataset has missing values, you can replace them with the median income of the dataset (for numerical variables), or you can drop rows where income is missing if it's a small portion of the dataset.





**FOLLOW ALONG**

**For learning more on  
Data Science, AI and  
Generative AI**



**Dinesh Lal**