

### **1. Difference between Linear Regression and Logistic Regression:**

- Linear regression is used for predicting continuous outcomes, whereas logistic regression is used for classification tasks where the outcome is binary (two classes) or multinomial (more than two classes).
- Linear regression predicts the value of a dependent variable based on independent variables by fitting a straight line to the data, while logistic regression predicts the probability of a binary outcome using a logistic function, which produces an S-shaped curve.
- Example scenario: Predicting whether a customer will churn (yes or no) based on their demographics and purchase history is more appropriate for logistic regression than linear regression because the outcome is binary (churned or not churned).

### **2. Cost Function and Optimization in Logistic Regression:**

- The cost function used in logistic regression is the logistic loss function (also called cross-entropy loss or log loss).
- The goal is to minimize the difference between the predicted probabilities and the actual outcomes.
- Optimization is typically done using iterative optimization algorithms such as gradient descent or its variants like stochastic gradient descent.

### **3. Regularization in Logistic Regression:**

- Regularization in logistic regression involves adding a penalty term to the cost function to prevent overfitting.
- It helps in controlling the complexity of the model by shrinking the coefficients towards zero.
- The two most common types of regularization in logistic regression are L1 regularization (Lasso) and L2 regularization (Ridge).

### **4. ROC Curve and Evaluation of Logistic Regression:**

- The ROC (Receiver Operating Characteristic) curve is a graphical plot that illustrates the diagnostic ability of a binary classifier as its discrimination threshold is varied.
- It plots the True Positive Rate (sensitivity) against the False Positive Rate (1-specificity) for different threshold values.
- The area under the ROC curve (AUC-ROC) is a common metric used to evaluate the performance of a logistic regression model, with a higher AUC indicating better performance.

## **5. Feature Selection Techniques in Logistic Regression:**

- Common techniques for feature selection in logistic regression include:
  - L1 regularization (Lasso): It encourages sparsity by shrinking some coefficients to zero, effectively performing feature selection.
  - Recursive Feature Elimination (RFE): It recursively removes features and evaluates their importance until the desired number of features is reached.
  - Information Gain: It measures the reduction in entropy or uncertainty about the target variable after a feature is known.
- These techniques help improve model performance by reducing overfitting, simplifying the model, and removing irrelevant or redundant features.

## **6. Handling Imbalanced Datasets in Logistic Regression:**

- Imbalanced datasets occur when one class dominates the other(s), leading to biased models.
- Strategies for handling imbalanced datasets in logistic regression include:
  - Oversampling the minority class (e.g., using techniques like SMOTE)
  - Undersampling the majority class
  - Using techniques like class weights or cost-sensitive learning to give more importance to minority class samples during training

## **7. Common Issues and Challenges in Logistic Regression:**

- Multicollinearity among independent variables: Addressed by removing highly correlated variables or using dimensionality reduction techniques like PCA.
- Outliers: Detected and treated using robust statistical methods or transformed.
- Missing data: Imputed using techniques like mean imputation, median imputation, or advanced methods like KNN imputation.
- Model interpretability: Addressed by understanding the coefficients' magnitude and direction, as well as using techniques like feature importance analysis