Introduction to Logistic Regression

Overview of Logistic Regression

- **Definition**: Logistic regression is a statistical method used for predicting binary outcomes, typically represented as 0 and 1 (e.g., yes/no, success/failure).
- Role in Supervised Learning: Logistic regression falls under the category of classification algorithms and models the relationship between independent variables and a binary dependent variable.

Conceptual Foundations of Logistic Regression

Key Concepts

- Binary Classification: Classifies data points into one of two categories based on predictor variables (e.g., spam or not spam).
- S-shaped Curve: Uses the logistic function to map real-valued numbers into values between 0 and 1, indicating probabilities instead of continuous outcomes.

$$P(Y=1|X) = \frac{1}{1 + e^{-(\beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_n X_n)}}$$
(1)

Where:

- P(Y = 1|X) = Probability of the event occurring (class 1)
- $\blacksquare \beta_0 = Intercept$
- lacksquare $\beta_1, \beta_2, ..., \beta_n = \text{Coefficients of the predictor variables } X_1, X_2, ..., X_n$

Key Points and Applications

Key Points

- **I** Binary Output: Specifically tailored for binary targets.
- **Interpretation of Coefficients**: Represents change in the log odds for a one-unit increase in the predictor variable.
- **Applications**: Used in medical diagnosis, credit scoring, and marketing response prediction.
- 4 Loss Function: Utilizes Maximum Likelihood Estimation (MLE) for coefficient optimization.
- **Extension to Multinomial**: Can extend using techniques like One-vs-Rest or Softmax Regression.

Examples

Healthcare Example: Predicting diabetes based on features like age and weight

What is Logistic Regression? - Definition

Definition of Logistic Regression

Logistic regression is a statistical method used for binary classification in supervised learning. It predicts the probability that a given input point belongs to a particular class (e.g., yes/no, pass/fail).

What is Logistic Regression? - Key Concepts

- Binary Outcomes: Focuses on binary outcomes with two possible classes. Examples:
 - Predicting disease presence: Yes (1) or No (0)
 - Spam detection: Spam (1) or Not Spam (0)
- Logistic Function: Converts predicted values into probabilities:

$$P(Y=1|X) = \frac{1}{1 + e^{-(\beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_n X_n)}}$$
(2)

Odds and Odds Ratio: Estimates the log odds of probability:

log odds =
$$\beta_0 + \beta_1 X_1 + \beta_2 X_2 + ... + \beta_n X_n$$
 (3)



6/1

Logistic vs Linear Regression

Differences Between Logistic and Linear Regression

- Output Format:
 - Linear Regression: Continuous values (e.g., prices)
 - Logistic Regression: Probabilities (0 to 1)
- Error Measurement:
 - Linear Regression: Mean Squared Error (MSE)
 - Logistic Regression: Likelihood estimation for binary cross-entropy loss
- Assumption:
 - Linear Regression: Linear relationship
 - Logistic Regression: Logistic relationship ideal for binary outcomes



Mathematical Foundation - Overview

The logistic function is a mathematical representation central to logistic regression. It predicts probabilities of binary outcomes (e.g., success/failure, yes/no) and constrains these probabilities between 0 and 1.

- Unlike linear regression, which predicts continuous outcomes,
- Logistic regression is specifically for binary outcomes.
- Ensures predicted probabilities are within the range [0, 1].

The Logistic Function

The logistic function is defined as:

$$f(z) = \frac{1}{1 + e^{-z}} \tag{4}$$

Where:

- *e* is Euler's number (approximately 2.71828).
- z is a linear combination of input features:

$$z = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_n X_n \tag{5}$$

- lacksquare β_0 : intercept (bias)
- $\beta_1, \beta_2, \dots, \beta_n$: coefficients for features X_1, X_2, \dots, X_n

The function f(z) outputs a value between 0 and 1, interpreted as the probability of the positive class (e.g., P(Y=1|X)).

Example of Logistic Function

Consider predicting if a student will pass (1) or fail (0) based on hours studied (X). The logistic regression model yields:

$$z = -4 + 1.2 \cdot \text{Hours Studied} \tag{6}$$

Example Calculation for 5 hours studied:

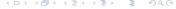
Calculate z:

$$z = -4 + 1.2 \cdot 5 = 2$$

2 Calculate P(Y = 1|X):

$$P(Y = 1|X) = f(2) = \frac{1}{1 + e^{-2}} \approx 0.8808$$

This indicates a high likelihood that the student will pass.



10 / 1

Logistic Function Equation

Overview

The logistic function equation transforms a linear combination of inputs into a probability, crucial for binary classification problems in supervised learning.

Understanding the Logistic Function

- The logistic function outputs values between 0 and 1.
- It is used for predicting binary outcomes (e.g., success/failure).
- Ideal for supervised learning applications, particularly in logistic regression.

The Logistic Function Equation

$$P(Y = 1|X) = \frac{1}{1 + e^{-z}}$$
 where $z = \beta_0 + \beta_1 X_1 + \dots + \beta_n X_n$ (7)

- P(Y = 1|X): Predicted probability that outcome Y is 1 given X.
- \bullet e^{-z} : Exponential function ensuring output is between 0 and 1.
- z: Linear combination of inputs weighted by coefficients.
 - β_0 : Intercept (log odds when all X are zero).
 - β_1, \ldots, β_n : Coefficients influencing the log odds.



13 / 1

Example Interpretation

- Consider:
 - $\beta_0 = -4$
 - $\beta_1 = 0.5$ (influence of X_1)
 - $X_1 = 5$
- Calculate z:

$$z = -4 + 0.5 \times 5 = -4 + 2.5 = -1.5$$
 (8)

■ Plug z into the logistic function:

$$P(Y=1|X) = \frac{1}{1+e^{1.5}} \approx 0.18 \tag{9}$$

■ Interpretation: 18% probability that the outcome Y = 1 given $X_1 = 5$.



Key Points to Emphasize

- The logistic function converts linear combinations into probabilities.
- lacktriangle Coefficients eta are optimized using Maximum Likelihood Estimation (MLE).
- Understanding z affects interpretations of the logistic model's outputs.

Next Steps

Upcoming Content

Prepare for the next slide focused on the cost function used in logistic regression, which automates the optimization of coefficients!

Cost Function in Logistic Regression - Overview

- A cost function quantifies how well the model predicts outcomes versus actual values.
- The objective in logistic regression is to minimize the cost function to enhance prediction accuracy.

Log-Loss Function in Logistic Regression

- The cost function in logistic regression is known as the log-loss function or cross-entropy loss.
- It quantifies the performance of a classification model with outputs as probabilities between 0 and 1.

Mathematical Formulation of Log-Loss

Log-Loss for a Single Data Point

$$Log Loss(y, \hat{y}) = -(y log(\hat{y}) + (1 - y) log(1 - \hat{y}))$$
(10)

- Where:
 - y = Actual label (0 or 1)
 - \hat{y} = Predicted probability of the positive class

July 10, 2025

19 / 1

Total Cost Function

Total Cost for All Predictions

For a dataset with *m* samples:

$$J(\beta) = -\frac{1}{m} \sum_{i=1}^{m} (y_i \log(\hat{y}_i) + (1 - y_i) \log(1 - \hat{y}_i))$$
 (11)



 $July\ 10,\ 2025 \\ 20\ /\ 1$

Advantages of Log-Loss

- Sensitive to Differences: Log-loss penalizes confident mistakes more heavily.
- **Probabilistic Interpretation**: It presents a smooth landscape for optimization, simplifying the minimization process.
- As predicted probabilities diverge from the true labels, log-loss can grow towards infinity.

Practical Application: Python Code

def log loss(y true, y pred):

v true = np.array([0, 1, 1, 0])

y pred = np.array([0.1, 0.9, 0.8, 0.2])

print("Log_Loss:", log loss(y true, y pred))

import numpy as np

```
epsilon = 1e-15  # To prevent log(0)
  y_pred = np.clip(y_pred, epsilon, 1-epsilon)
# Clip predictions
  return -np.mean(y_true * np.log(y_pred) + (1 - y_true) * np.log(
# Example usage:
```

Conclusion

- The log-loss function is vital for evaluating logistic regression models.
- Minimizing total cost is essential for effective model training.
- A solid understanding of the cost function provides insight into logistic regression in binary classification.

Understanding Gradient Descent for Logistic Regression

1. Introduction to Gradient Descent

- Objective: Minimize the cost function (log-loss) in logistic regression by finding optimal parameters (weights).
- **Gradient Descent**: An iterative optimization algorithm that adjusts parameters to minimize the cost.

How Gradient Descent Works

- Initialization:
 - Start with random weights, w.
- Compute the Cost:
 - Log-loss function:

$$J(w) = -\frac{1}{m} \sum_{i=1}^{m} \left[y^{(i)} \log(h_w(x^{(i)})) + (1 - y^{(i)}) \log(1 - h_w(x^{(i)})) \right]$$
(12)

where:

- $\mathbf{m} = \mathsf{number}$ of training examples
- $h_w(x) = hypothesis function (sigmoid function)$
- $y^{(i)} = \text{actual label for example } i$
- Calculate the Gradient:
 - Gradient of cost function:

$$\nabla J(w) = \frac{1}{m} \sum_{i=1}^{m} (h_w(x^{(i)}) - y^{(i)}) x^{(i)}$$
(13)



Key Points and Conclusion

Key Points to Emphasize

- Convergence: Continues until cost function changes are negligible.
- Learning Rate Selection:
 - \blacksquare Small α : Slow convergence.
 - Large α : Risks overshooting.
- Types of Gradient Descent:
 - Batch Gradient Descent: Uses all data points.
 - Stochastic Gradient Descent: Updates for each training example.
 - Mini-batch Gradient Descent: Combines both approaches.

Conclusion

Gradient descent is crucial for training logistic regression models, with its understanding essential for effective classification algorithm development.

Evaluation Metrics Overview

- Evaluation metrics are crucial for assessing the performance of logistic regression models.
- Common metrics include:
 - Accuracy
 - Precision
 - Recall
 - F1-Score
- Each metric provides insight into different aspects of model performance.

1. Accuracy

Definition

Accuracy is the proportion of correctly predicted instances among the total instances.

Formula

$$Accuracy = \frac{TP + TN}{TP + TN + FP + FN}$$
 (15)

Where:

- TP = True Positives
- TN = True Negatives
- FP = False Positives
- FN = False Negatives

Example

2. Precision and Recall

2. Precision

- **Definition**: Precision indicates the accuracy of positive predictions.
- Formula:

$$Precision = \frac{TP}{TP + FP}$$
 (16)

29 / 1

Example

From our previous example, if we have 80 positive predictions (70 TP and 10 FP):

Precision =
$$\frac{70}{70 + 10} = \frac{70}{80} = 0.875$$
 or 87.5%

3. Recall

Definition: Recall measures the model's ability to identify relevant instances (true

4. F1-Score

Definition

The F1-score is the harmonic mean of precision and recall.

Formula

$$F1-Score = 2 \times \frac{Precision \times Recall}{Precision + Recall}$$
 (18)

Example

Using the values computed earlier:

F1-Score =
$$2 \times \frac{0.875 \times 0.778}{0.875 + 0.778} \approx 0.823$$
 or 82.3%

Key Points

Implementing Logistic Regression - Overview

What is Logistic Regression?

Logistic Regression is a statistical model used for binary classification problems. This guide focuses on implementing it using Python's scikit-learn library.

Step 1: Import Necessary Libraries

Before you begin, ensure you have installed the required libraries. You can install them using pip as shown below:

```
# Install with pip if necessary
# !pip install numpy pandas scikit—learn
import numpy as np
import pandas as pd
from sklearn.model_selection import train_test_split
from sklearn.linear_model import LogisticRegression
from sklearn.metrics import accuracy_score, confusion_matrix, classif
```

Step 2: Load Your Dataset

Start by loading your dataset. For illustration, the popular Iris dataset is used.

```
from sklearn.datasets import load_iris
data = load_iris()
X = data.data[data.target != 0]  # Two classes for binary classifica
y = data.target[data.target != 0]
```

Step 3: Pre-process the Data

- Ensure your data is clean.
- Handle missing values, scale features, and transform categorical variables if necessary.

Step 4: Split the Dataset

Divide your dataset into training and testing sets to evaluate the model's performance.

 $X_{train}, X_{train}, y_{train}, y_{train}, test_split(X, y, test_size = train_test_split(X, y, test_size = train_test_$

Step 5: Create the Logistic Regression Model

Initialize the Logistic Regression model:

```
model = LogisticRegression()
model.fit(X train, y train)
```

Step 6: Make Predictions

Use the trained model to make predictions on the test set.

$$y_pred = model.predict(X_test)$$

July 10, 2025

Step 7: Evaluate the Model

Assess the model's performance using:

- Accuracy Score
- Confusion Matrix
- Classification Report

```
accuracy = accuracy_score(y_test, y_pred)
conf_matrix = confusion_matrix(y_test, y_pred)
class_report = classification_report(y_test, y_pred)

print("Accuracy:", accuracy)
print("Confusion_Matrix:\n", conf_matrix)
print("Classification_Report:\n", class_report)
```

Key Points to Emphasize

- Data Preparation: Quality input data significantly influences model performance.
- Model Evaluation: Always use appropriate metrics to ensure reliability.
- Benefits of Logistic Regression: Easy interpretation and good performance on binary outcomes.

Conclusion

You now have a foundation to implement Logistic Regression for binary classification tasks. Next, we will explore the assumptions underlying logistic regression models to deepen our understanding.

Logistic Regression Assumptions - Overview

Overview of Logistic Regression:

Logistic regression is a statistical method used for binary classification problems, predicting the probability that an instance belongs to a particular category. It is particularly popular for its simplicity and interpretability.

Logistic Regression Assumptions - Key Assumptions

Key Assumptions of Logistic Regression:

- Binary Outcome:
 - Designed for binary outcomes (0/1, True/False).
- Linearity Between Features and Log-Odds:

Assumes a linear relationship between independent variables and the log-odds of the dependent variable:

$$\mathsf{Log}\mathsf{-}\mathsf{Odds} = \mathsf{log}\left(\frac{p}{1-p}\right) \tag{19}$$

represented as:

$$Log-Odds = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \ldots + \beta_n X_n$$
 (20)

- **3** Independence of Observations:
 - Assumes that observations are independent of one another.
- 4 No Multicollinearity:

Independent variables must not be highly correlated.



Logistic Regression Assumptions - Key Points and Practical Considerations

Key Points to Emphasize:

- Understanding these assumptions is crucial for correct application and interpretation of logistic regression results.
- Violations can lead to model misfit and inaccurate predictions.

Practical Consideration:

- Visualize relationships with scatterplots or residual plots to check for linearity.
- Check for multicollinearity using statistical tests or variance inflation factors (VIF).

Logistic Regression Assumptions - Conclusion and Additional Resources

Conclusion:

By keeping these assumptions in mind, we can ensure appropriate application of logistic regression, leading to robust insights.

Additional Resources:

- Use visualization libraries like Matplotlib or Seaborn in Python for diagnostic plots.
- Perform exploratory data analysis (EDA) to validate assumptions before modeling.

Logistic Regression Assumptions - Code Snippet

```
Code Snippet (Python):
from statsmodels api import Logit
import pandas as pd
# Example DataFrame
data = pd.DataFrame({
    'hours studied': [1, 2, 3, 4, 5],
    'passed exam': [0.0, 1, 1, 1]
})
# Fitting the logistic regression model
model = Logit (data ['passed exam'], data ['hours studied'])
result = model.fit()
```

Applications of Logistic Regression - Introduction

Overview

Logistic regression is a powerful statistical method used for binary classification problems, where the outcome variable is dichotomous. This technique is widely utilized to make predictions about binary outcomes based on multiple predictor variables.

Applications of Logistic Regression - Industries

Healthcare

- Disease Prediction: Predicting disease likelihood based on factors such as age and BMI.
- Clinical Decision-Making: Identifying patients at risk for surgical complications.

2 Finance

- Credit Scoring: Assessing creditworthiness of loan applicants.
- Fraud Detection: Identifying potentially fraudulent transactions.

Marketing

- Customer Retention: Determining factors contributing to customer churn.
- Targeted Advertising: Predicting customer responses to marketing campaigns.



July 10, 2025

Key Points and Formula Overview

Key Points

- Versatility: Applicable in various domains.
- Interpretability: Coefficients provide insights into predictor relationships.
- Decision-Making Tool: Aids in informed decisions through quantified probabilities.

Logistic Regression Model

The logistic regression model predicts the probability P(Y = 1|X) as follows:

$$P(Y=1|X) = \frac{1}{1 + e^{-(\beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_n X_n)}}$$
(21)

where:

• e is the base of the natural logarithm,

July 10, 2025

48 / 1

Multiclass Logistic Regression

Logistic regression is widely used for binary classification, but many real-world problems involve multiple categories.

- Multiclass classification can be handled using:
 - One-vs-Rest (OvR)
 - Softmax Regression

1. One-vs-Rest (OvR) Approach

Concept

The One-vs-Rest method creates a separate binary classifier for each class in the dataset.

- \blacksquare Train K classifiers for K classes.
- 2 Each classifier C_k predicts if a sample belongs to class k.
- 3 For a new input, the class with the highest predicted probability is chosen.

Example

For classes A, B, and C:

- C_A: Predict A vs. (B, C)
- C_B: Predict B vs. (A, C)
- C_C: Predict C vs. (A, B)

2. Softmax Regression (Multinomial Logistic Regression)

Concept

Softmax regression generalizes logistic regression for multi-class predictions in a single model.

How it works

Given a feature vector x:

$$P(y = i | \mathbf{x}) = \frac{e^{\theta_i^T \mathbf{x}}}{\sum_{j=1}^{K} e^{\theta_j^T \mathbf{x}}}$$

where θ_i are the weights for class i.

Example

For classes A, B, and C, with logits:

■ For A: 2

Case Study Example

In this slide, we will explore a practical scenario where logistic regression is applied to solve a real-world problem related to predicting customer defaults on loans.

Introduction to the Case Study

Case Study: Predicting Customer Default on Loans

A financial institution aims to predict whether a customer will default on a loan. Logistic regression will be utilized to handle this binary outcome.

Data Collection

- Input Features:
 - Credit Score
 - Annual Income
 - Loan Amount
 - Employment Status
 - Age
- Binary Outcome:
 - Default (1) or No Default (0)

Logistic Regression Model Overview

Logistic regression is used for binary target variables and models the relationship between input features and event probabilities.

The logistic function transforms a linear combination of inputs into a probability:

$$P(Y=1|X) = \frac{1}{1 + e^{-(\beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_n X_n)}}$$
(22)

where:

- $\beta_0 = Intercept$
- lacksquare $\beta_1, \beta_2, ..., \beta_n = \text{Coefficients for features}$
- $X_1, X_2, ..., X_n =$ Input features



July 10, 2025

Implementation Steps

- Data Preparation:
 - Preprocess data (handle missing values, encode categorical variables).
- 2 Model Training:
 - Split data into training and testing sets.
 - Fit the model using Python's 'scikit-learn':

```
from sklearn.model_selection import train_test_split
from sklearn.linear_model import LogisticRegression
```

```
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size
model = LogisticRegression()
model.fit(X train, y train)
```

- **3** Model Evaluation:
 - Use accuracy, precision, recall, and the confusion matrix:

```
from sklearn.metrics import classification_report, confusion_matrix
```

Key Takeaways

- Predictive Analysis: Logistic regression assesses probabilities of customer default based on financial metrics.
- Interpretability: Coefficients reveal the impact of each feature.
- Decision Making: Enables informed decisions on loan approvals.

Conclusion

Logistic regression is a powerful tool in finance for predicting binary outcomes. By effectively implementing and interpreting this model, institutions can enhance risk assessment and improve financial strategies.

Challenges and Limitations - Overview

Objective

Understand the principal challenges and limitations encountered when using logistic regression for predictive modeling.

Challenges and Limitations - Assumptions

- Assumptions of Logistic Regression
 - Logistic regression assumes a linear relationship between independent variables and log odds.
 - Violation of this assumption leads to inaccurate predictions.

Example

If the true relationship is quadratic or exponential, logistic regression may struggle to fit the data.

Challenges and Limitations - Multicollinearity and Data Imbalance

Multicollinearity

- Occurs when independent variables are highly correlated.
- Inflates variance of coefficient estimates, making tests unreliable.

Example

Including both "age" and "years of experience" might introduce multicollinearity.

3 Data Imbalance

- Logistic regression may predict the majority class if classes are imbalanced.
- Accuracy may be high while sensitivity is poor.

Solution

Use oversampling, undersampling, or cost-sensitive learning techniques.



Challenges and Limitations - Overfitting and Non-linearity

4 Overfitting

- Model captures noise as signal leading to poor performance on new data.
- Occurs when too many predictors or interactions are included.

Prevention

Use regularization techniques (e.g., L1 or L2 penalties).

5 Non-linearity in Relationships

- Logistic regression assumes linearity in log-odds; non-linear true relationships lead to suboptimal predictions.
- Consider polynomial or interaction terms to fit non-linear relationships.

Logistic Regression Formula

The logistic regression model is defined as:

$$P(Y=1|X) = \frac{1}{1 + e^{-(\beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_n X_n)}}$$
(23)

Where:

- $lackbr{\blacksquare} P(Y=1|X) = ext{Probability that } Y ext{ equals 1 given inputs } X_1, X_2, \dots, X_n.$
- lacksquare β_0 = Intercept.
- lacksquare β_i = Coefficients for predictors.



July 10, 2025

Key Points and Conclusion

- Logistic regression is simple and efficient but has notable limitations.
- Understanding data characteristics and model assumptions is crucial for effectiveness.
- Addressing challenges like multicollinearity, data imbalance, and overfitting can significantly enhance performance.

Conclusion

Being aware of the challenges associated with logistic regression leads to better model diagnostics and improved predictions.

Future Directions - Overview

Overview

Logistic regression has proven to be a powerful tool in statistical modeling and machine learning. As technology evolves, so do methodologies and applications of logistic regression.

- Exploring emerging trends and future directions.
- Emphasizing advancements that enhance utility and efficiency.

Future Directions - Integration with Advanced ML Techniques

- 1. Integration with Advanced Machine Learning Techniques
 - **Hybrid Models**: Combining logistic regression with ensemble methods for improved accuracy.
 - Examples:
 - Initial feature selection with logistic regression.
 - Using Random Forest for capturing complex feature interactions.

Future Directions - Feature Engineering and Explainability

2. Feature Engineering and Selection

- Automated Feature Engineering: Trends towards automating feature creation with techniques like deep learning embeddings.
- Key Point: Robust feature selection improves model performance and decreases overfitting.

4. Explainability and Transparency

- Model Interpretability: Enhancing interpretability with tools like SHAP.
- Key Point: Crucial for applications like healthcare and finance.

July 10, 2025

Future Directions - Addressing Imbalanced Data

3. Addressing Imbalanced Data

- Improved Techniques: Focus on handling class imbalance through methods like SMOTE and cost-sensitive learning.
- Illustration: Utilizing SMOTE in fraud detection to train reliable logistic regression models.

Future Directions - Application in Novel Domains

5. Application in Novel Domains

- **Expansion into New Areas:** Logistic regression's application in various fields.
- Examples:
 - **Healthcare**: Predicting disease presence based on risk factors.
 - Social Sciences: Analyzing voting behavior and survey responses.

Future Directions - Summary

Summary

The future of logistic regression is promising, with advancements focusing on:

- Accuracy
- Interpretability
- Application in diverse fields

Key Takeaway: Ongoing advancements will enrich logistic regression's capabilities, making it a versatile tool in predictive analytics.

Summary and Key Takeaways - Overview of Logistic Regression

Overview

Logistic Regression is a powerful statistical method used for binary classification problems, helping to predict one of two outcomes based on input variables. It is particularly effective when the dependent variable is categorical (e.g., pass/fail, yes/no).

Summary and Key Takeaways - Key Concepts

Logistic Function:

- Maps predicted values between 0 and 1 using the logistic (sigmoid) function.
- Formula:

$$P(Y=1|X) = \frac{1}{1+e^{-z}}$$
 (24)

• Where $z = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + ... + \beta_n X_n$.

Odds and Odds Ratio:

- **Odds**: Ratio of the probability of an event occurring to it not occurring.
- Odds Ratio: Indicates change in odds for a one-unit increase in the predictor variable.

Odds Ratio =
$$e^{\beta_i}$$
 (25)



72 / 1

July 10, 2025

Summary and Key Takeaways - Interpretation of Coefficients

Interpretation of Coefficients

- Positive coefficients ($\beta_i > 0$) increase the odds of the outcome.
- Negative coefficients (β_i < 0) decrease the odds.

Model Evaluation

- Confusion Matrix: For performance metrics (accuracy, precision, recall, F1-score).
- ROC Curve and AUC: Trade-off between sensitivity and specificity, with AUC indicating overall model performance.

July 10, 2025

Summary and Key Takeaways - Examples and Conclusions

Examples Illustrating Logistic Regression

- **Example 1**: Predicting whether a student passes (1) or fails (0) based on study hours and attendance.
- Example 2: Identifying whether an email is spam (1) or not spam (0) based on features like word frequency and sender information.

Key Points to Emphasize

- Logistic Regression can adapt to non-linear relationships through transformations.
- Check for multicollinearity among predictors to avoid distortion.
- Assumes independence of errors and requires a sufficient sample size.

Conclusion

Logistic regression is a foundational tool in machine learning for binary classification tasks.

Q&A Session - Overview

Overview

In this open floor session, we invite you to ask questions and engage in discussions about Logistic Regression, a fundamental supervised learning algorithm used for binary classification tasks.

Q&A Session - Key Concepts Recap

- Logistic Regression Purpose: Predicts the probability of a given input belonging to a certain class (0 or 1).
- Sigmoid Function: The model outputs a probability using the sigmoid function:

$$P(Y=1|X) = \frac{1}{1 + e^{-(\beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_n X_n)}}$$
(26)

■ **Cost Function**: It uses the log-likelihood function:

$$Cost = -\frac{1}{m} \sum_{i=1}^{m} [y^{(i)} \log(h(x^{(i)})) + (1 - y^{(i)}) \log(1 - h(x^{(i)}))]$$
 (27)



July 10, 2025

76 / 1

Q&A Session - Applications and Discussion Prompts

Applications of Logistic Regression

- Medical Diagnosis: Predicting disease presence based on test results.
- Customer Churn Prediction: Identifying potential service cancellations.
- Credit Scoring: Evaluating borrower default risks.

Discussion Prompts

- What scenarios might logistic regression perform poorly in?
- Can you provide examples where logistic regression offers valuable insights?
- How could feature selection affect model performance?



July 10, 2025

Q&A Session - Interactive Component

model = LogisticRegression()
model.fit(X train. v train)

Interactive Component

Practical Exercise: If time permits, we'll review a dataset and perform logistic regression using Python's Scikit-learn library.

```
from sklearn.linear_model import LogisticRegression
from sklearn.model_selection import train_test_split
from sklearn.metrics import confusion_matrix, accuracy_score

# Sample Data Split
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size = # Model Training
```