



John Smith, Ph.D.

Department of Computer Science  
University Name

Email: [email@university.edu](mailto:email@university.edu)  
Website: [www.university.edu](http://www.university.edu)

July 7, 2025

# Introduction to Supervised Learning

## What is Supervised Learning?

Supervised learning is a type of machine learning where an algorithm learns from labeled training data to make predictions on unseen data. Each training example consists of an input-output pair, where the input is the feature set and the output is the label.

## Key Characteristics

- **Labeled Data:** Each training sample includes known outputs (labels).
- **Predictive Modeling:** The goal is to create a model that can predict the output for new, unseen data based on the learned relationships.

# Importance of Supervised Learning in Data Mining

- 1 **Decision Making:** Aids in informed decision-making by predicting outcomes based on historical data.
  - Example: Predicting whether a customer will buy a product based on age, gender, and previous purchasing behavior.
- 2 **Classification and Regression:** Solves both classification problems (e.g., spam detection) and regression problems (e.g., sales forecasting).
- 3 **Versatility:** Wide range of applications in finance (credit scoring), healthcare (disease diagnosis), and marketing (customer segmentation).

# Types and Algorithms of Supervised Learning

## ■ Types:

- **Classification:** Predicting discrete labels (e.g., spam or not spam).
- **Regression:** Predicting continuous values (e.g., house prices).

## ■ Common Algorithms:

- Logistic Regression
- Decision Trees
- Support Vector Machines (SVM)
- Random Forests
- Neural Networks

## Example Scenario

### Scenario

Consider a bank that wants to forecast whether a loan applicant will default on a loan. The past data includes various applicants' features (income, credit history, etc.) and whether they defaulted or not (the label). A supervised learning algorithm can learn from this dataset to predict the likelihood of default for future applicants.

# Conclusion

Understanding supervised learning is fundamental to harnessing the potential of data mining for predictive analytics. It equips data scientists with methodologies to turn historical data into actionable insights.

## Code Snippet: Simple Example with Python

```
1 from sklearn.model_selection import train_test_split
2 from sklearn.linear_model import LogisticRegression
3 from sklearn.metrics import accuracy_score
4
5 # Sample Data
6 X = [[1, 2], [2, 3], [3, 4], [4, 5]] # Features
7 y = [0, 0, 1, 1] # Labels
8
9 # Split Data
10 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2)
11
12 # Model Training
13 model = LogisticRegression()
14 model.fit(X_train, y_train)
15
16 # Predictions
17 predictions = model.predict(X_test)
```



# What is Logistic Regression? - Definition

## Definition

**Logistic Regression** is a statistical method used for binary classification tasks, where the output or target variable can take one of two possible values or classes (e.g., 0 and 1, Yes and No, True and False).

# What is Logistic Regression? - Purpose and Key Concepts

## Purpose

The primary purpose of logistic regression is to model the probability that a given input point belongs to a particular category.

### ■ Binary Classification:

- Logistic regression is suitable for scenarios where the outcome is binary. Example: Predicting whether an email is “spam” (1) or “not spam” (0).

### ■ Odds and Probability:

- Logistic regression predicts outcomes using the concept of odds. The odds represent the ratio of the probability of an event occurring to the probability of it not occurring.
- Probability ( $P$ ) of an event occurring is transformed with the logistic function to ensure outputs are between 0 and 1.

# What is Logistic Regression? - Logistic Function and Example

## Logistic Function

The logistic function (also known as the sigmoid function) is used to model the relationship between one or more independent variables and the binary response variable:

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

Here,  $P(Y = 1|X)$  is the predicted probability of the outcome,  $X$  represents the independent variables, and  $\beta_0, \beta_1, \dots, \beta_n$  are coefficients that the model tries to learn.

## Example

Consider a healthcare scenario to predict if a patient has a disease (1) or not (0):

- **Input Variables:** Age, Blood Pressure
- **Outcome Variable:** Disease Status (1 = Disease, 0 = No Disease)

# Mathematical Foundation of Logistic Regression - Part 1

## Understanding the Logistic Function

The logistic function transforms a linear combination of input features into a probability value between 0 and 1, defined as:

$$f(z) = \frac{1}{1 + e^{-z}} \quad (2)$$

Where:

- $z$  is the linear combination:  $z = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_n X_n$
- $e$  is the base of the natural logarithm (approximately 2.71828).

## Mathematical Foundation of Logistic Regression - Part 2

### Key Properties of the Logistic Function:

- 1 **S-shaped Curve:** As  $z \rightarrow -\infty$ ,  $f(z) \rightarrow 0$  and as  $z \rightarrow +\infty$ ,  $f(z) \rightarrow 1$ .
- 2 **Threshold:** A common threshold at  $f(z) = 0.5$  classifies observations into class 0 or class 1.

### Example:

- For  $X = 0$ :  $f(0) = \frac{1}{1+1} = 0.5$ .
- As  $z$  increases (e.g.,  $z = 2$ ):  $f(2) \approx 0.88$ , indicating higher probability of class 1.

# Mathematical Foundation of Logistic Regression - Part 3

## Odds Ratio

The odds ratio indicates the odds of success versus failure:

$$\text{Odds} = \frac{P(Y = 1)}{P(Y = 0)} = \frac{f(z)}{1 - f(z)} \quad (3)$$

## Key Points:

- Odds ratio  $> 1$ : Positive association (greater odds of outcome).
- Odds ratio  $< 1$ : Negative association.

## Example:

- Odds ratio of 3 means the event is 3 times more likely to occur than not.

# Mathematical Foundation of Logistic Regression - Summary

## Summary of Key Concepts:

- **Logistic Regression:** Utilizes the logistic function for binary outcomes.
- **Logistic Function:** Maps real values into the  $(0, 1)$  interval.
- **Odds Ratio:** Measures effect size of predictors and assists in interpretation.

**Next Steps:** In the upcoming slide, we will discuss the **Assumptions of Logistic Regression** to ensure the validity of the model's outcomes.

# Assumptions of Logistic Regression

## Overview

Logistic Regression is a powerful statistical tool used for binary classification problems. To ensure the validity and reliability of the model's outputs, it is essential to adhere to certain key assumptions.



# Key Assumptions - 1

## 1 Binary Outcome Variable

- **Explanation:** The dependent variable must be binary (0/1, True/False, Yes/No).
- **Example:** Predicting whether a student passes (1) or fails (0) an exam based on study hours.

## 2 Independent Observations

- **Explanation:** Each observation in the dataset must be independent of others.
- **Example:** In a study to determine if diet impacts weight loss, the weight loss of one participant should not influence the weight loss of another.

## Key Assumptions - 2

### 3 Linearity of the Logit

- **Explanation:** There should be a linear relationship between the independent variables and the log-odds of the dependent variable.
- **Illustration:**

$$\log \left( \frac{p}{1-p} \right) = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_n X_n \quad (4)$$

### 4 No Multicollinearity

- **Explanation:** The independent variables should not be too highly correlated with each other.
- **Example:** Using both 'height' and 'weight' can be problematic if they are strongly correlated, making it difficult to determine the individual impact of each on the outcome.

## Key Assumptions - 3

### 5 Adequate Sample Size

- **Explanation:** Logistic regression requires a sufficient sample size to provide reliable estimates.
- **Key Point:** A rule of thumb is to have at least 10 events for each predictor variable to ensure the model is stable.

### Key Points to Emphasize

- Always check the assumptions before interpreting the results.
- Visualize data to understand relationships (scatter plots for linearity, VIF for multicollinearity).
- Consider transformations or variable selections if assumptions are violated.

# Conclusion

Understanding and validating the assumptions of logistic regression not only enhances the model's predictive power but also ensures more robust conclusions and insights drawn from the model. Prepare to apply these principles when implementing logistic regression in real-world scenarios!

# Implementing Logistic Regression - Overview

- Logistic regression is a statistical method for binary classification.
- It predicts the probability of a data point belonging to a certain category.
- Utilizes a logistic function to map predictions to probabilities.

# Implementing Logistic Regression - Steps

## 1 Import Libraries

```
1 import pandas as pd
2 import numpy as np
3 from sklearn.model_selection import train_test_split
4 from sklearn.linear_model import LogisticRegression
5 from sklearn.metrics import accuracy_score, confusion_matrix
```

## 2 Load and Prepare Dataset

```
1 data = pd.read_csv('data.csv') # Replace 'data.csv' with your dataset
```

## 3 Explore the Dataset

- Examine the first few rows with `data.head()`.
- Understand features and target variables.

# Implementing Logistic Regression - Continued Steps

## res Data Preprocessing

```
1 data['category'] = data['category'].map({'class_0': 0, 'class_1': 1})
```

## res Split the Dataset

```
1 X = data.drop('target', axis=1) # Replace 'target' with your dependent
  variable
2 y = data['target']
3 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2,
  random_state=42)
```

## res Initialize and Train the Model

```
1 model = LogisticRegression()
2 model.fit(X_train, y_train)
```

## res Make Predictions

# Implementing Logistic Regression - Evaluation

## ■ Evaluate the Model

```
1 accuracy = accuracy_score(y_test, y_pred)
2 cm = confusion_matrix(y_test, y_pred)
3 print("Accuracy:", accuracy)
4 print("Confusion Matrix:\n", cm)
```

## ■ Key Points

- Logistic function:  $\sigma(z) = \frac{1}{1+e^{-z}}$
- Evaluate your model using multiple metrics.
- Data quality is crucial for model performance.



# Implementing Logistic Regression - Conclusion

Implementing logistic regression with Python and Scikit-learn is straightforward. Key points of focus include:

- Importance of proper data preparation.
- Necessity of thorough model evaluation.
- Reliability of predictions hinges on input data quality.

# Data Preparation for Logistic Regression - Introduction

## Importance of Data Preprocessing

Data preprocessing is a critical step in building a logistic regression model. Well-prepared data enhances prediction accuracy and allows the model to learn effectively. Key components include:

- Feature Scaling
- Encoding Categorical Variables

# Data Preparation for Logistic Regression - Feature Scaling

## Feature Scaling

Feature scaling adjusts the range of the data. Logistic regression is sensitive to the scale of input features, thus scaling is essential.

- **Standardization:** Transforms data to have mean 0 and standard deviation 1.

$$z = \frac{x - \mu}{\sigma} \quad (5)$$

- **Normalization:** Rescales data to a range of [0, 1].

$$x' = \frac{x - \min(x)}{\max(x) - \min(x)} \quad (6)$$

## Example

A dataset may have features like age (0 to 100) and income (20,000 to 120,000). Without scaling, income could dominate due to its larger range.

# Data Preparation for Logistic Regression - Encoding Categorical Variables

## Encoding Categorical Variables

Logistic regression requires numerical inputs, necessitating conversion of categorical variables.

- **One-Hot Encoding:** Creates binary columns for each category. Useful for nominal variables.
  - Example: "Color" = ["Red", "Green", "Blue"] results in:
    - Color\_Red: [1, 0, 0]
    - Color\_Green: [0, 1, 0]
    - Color\_Blue: [0, 0, 1]
- **Label Encoding:** Assigns integers to categories. Appropriate for ordinal variables.
  - Example: "Size" = ["Small", "Medium", "Large"] maps to:
    - Small: 0
    - Medium: 1
    - Large: 2

# Data Preparation for Logistic Regression - Key Points and Code Snippet

## Key Points to Emphasize

- Careful data preparation enhances model performance.
- Scaling ensures features are on a similar scale, reducing bias.
- Correct encoding is essential for interpreting feature relationships.

## Code Snippet Example (Python)

```
1 from sklearn.preprocessing import StandardScaler, OneHotEncoder
2 from sklearn.compose import ColumnTransformer
3
4 # Data preparation pipeline
5 preprocessor = ColumnTransformer(
6     transformers=[
7         ('num', StandardScaler(), ['Age', 'Income']),
8         ('cat', OneHotEncoder(), ['Color'])
```

# Splitting the Dataset - Overview

## Best Practices for Dividing Data

In supervised learning, especially with logistic regression, it is vital to split your dataset into:

- **Training Set:** Used to fit the model.
- **Testing Set:** Used to evaluate model performance.

Proper dataset splitting is essential for model generalization to unseen data.

# Splitting the Dataset - Key Concepts

## 1 Training Set:

- Typically comprises **70-80%** of total data.
- Used for training the logistic regression model.

## 2 Testing Set:

- Contains the remaining **20-30%**.
- Validates the model's performance after training.

## 3 Validation Set (optional):

- A third set created from training data for tuning model parameters.
- Common split: 60% training, 20% validation, 20% testing.

# Splitting the Dataset - Best Practices

- **Random Sampling:**
  - Splits dataset randomly to ensure representation.
- **Stratified Sampling:**
  - Ensures proportional representation in both sets.
  - Important for binary classification tasks.
- **Avoiding Data Leakage:**
  - Prevents testing set information from corrupting the training set.

## Example

Consider a dataset with 1000 samples:

- **Total Samples:** 1000
- **Training Set:** 800 samples (80%)
- **Testing Set:** 200 samples (20%)



## Splitting the Dataset - Code Example

Here is a sample code snippet using scikit-learn for data splitting:

```
1 from sklearn.model_selection import train_test_split
2
3 # Assuming X is your feature set and y is the target variable
4 X_train, X_test, y_train, y_test = train_test_split(X, y,
5                                                    test_size=0.2,
6                                                    random_state=42,
7                                                    stratify=y)
```

### Key Points to Emphasize

- Importance of proper splitting influences model accuracy.
- Choose between Random and Stratified based on dataset needs.
- Consistency in random seed for reproducibility.

# Training the Model

## Understanding Logistic Regression Training

Logistic Regression is a statistical method for binary classification. It predicts the probability of a binary outcome based on one or more predictor variables (features). To train a logistic regression model, a dataset must be split into **training** and **testing** sets.

# Step-by-Step Process to Train the Model

## 1 Data Preparation:

- Ensure data is clean and preprocessed.
- Handle missing values, encode categorical variables, and normalize or scale numerical features when necessary.

## 2 Splitting the Dataset:

- Divide the dataset into **training (70-80%)** and **testing sets (20-30%)**, where the training set is used to train the model.

## 3 Choosing the Model: The logistic regression model can be expressed as:

$$P(Y = 1|X) = \frac{1}{1 + e^{-(\beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_n X_n)}} \quad (7)$$

# Fitting and Evaluating the Model

## res Fitting the Model:

- Use a training algorithm to find optimal parameters by minimizing the difference between observed outcomes and predicted probabilities. The common approach is **maximum likelihood estimation (MLE)**.
- Example in Python:

```
1 from sklearn.linear_model import LogisticRegression
2
3 model = LogisticRegression()
4 model.fit(X_train, y_train)
```

## res Evaluating the Model:

- Assess model performance using metrics such as **accuracy**, **precision**, **recall**, and **F1 score**.
- Use functions like `cross_val_score` or tools like `confusion_matrix` in scikit-learn.

## Making Predictions - Overview

In this slide, we will explore how to utilize a trained logistic regression model to make predictions on new, unseen data. This is the crucial step where our model moves from theory to practice, seeking to apply what it has learned.

# Making Predictions - Key Concepts

## 1 Model Prediction:

- After training the logistic regression model, we use learned coefficients to relate independent variables (features) to a dependent variable (target).
- New data is input into the logistic regression equation to generate probabilities for target classes.

## 2 Logistic Regression Equation:

$$P(y = 1|X) = \frac{1}{1 + e^{-(\beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_n X_n)}} \quad (8)$$

Here,  $P(y = 1|X)$  is the probability that the target variable  $y$  equals 1 given features  $X$ .

## 3 Thresholding:

- To convert predicted probabilities into class labels, apply a threshold (commonly 0.5):
  - If  $P(y = 1|X) > 0.5$ , classify as 1 (positive class).
  - If  $P(y = 1|X) \leq 0.5$ , classify as 0 (negative class).

## Making Predictions - Example

Assume we have a trained logistic regression model for predicting whether a student passes (1) or fails (0) an exam based on hours studied and previous grades.

- **New Data Input:**

- Hours studied = 5
- Previous grade = 80

- **Calculation:**

$$P(y = 1|X) = \frac{1}{1 + e^{-(\beta_0 + \beta_1 \cdot 5 + \beta_2 \cdot 80)}} \quad (9)$$

Let's say this results in a predicted probability of 0.7.

- **Classification:**

- Since  $0.7 > 0.5$ , we classify the student as passing the exam (1).

## Making Predictions - Key Points

- Predictions from a logistic regression model yield probabilities, not direct class predictions.
- Choosing the right threshold is crucial and can be adjusted depending on the application, such as medical diagnoses.
- Always back predictions with model evaluation metrics to gauge accuracy and reliability.



## Making Predictions - Code Snippet

```
1 import numpy as np
2 from sklearn.linear_model import LogisticRegression
3
4 # Assume X_train and y_train are your training data and target variable
5 model = LogisticRegression()
6 model.fit(X_train, y_train)
7
8 # New data for prediction
9 new_data = np.array([[5, 80]]) # New student data with hours studied and
    previous grade
10 probability = model.predict_proba(new_data)[:, 1] # Probability of passing
    (class 1)
11
12 # Thresholding to make class prediction
13 prediction = (probability > 0.5).astype(int)
14 print(f'Predicted Probability: {probability}, Class Prediction: {prediction}')
```

# Evaluating Model Performance - Overview

## Key Metrics for Logistic Regression

Evaluating the performance of a logistic regression model is crucial in understanding its effectiveness. Below, we detail the essential metrics used for this evaluation.

## Evaluating Model Performance - 1. Accuracy

- **Definition:** Measures the proportion of true results (both true positives and true negatives) among the total cases examined.

- **Formula:**

$$\text{Accuracy} = \frac{TP + TN}{TP + TN + FP + FN} \quad (10)$$

- **Where:**

- $TP$  = True Positives
- $TN$  = True Negatives
- $FP$  = False Positives
- $FN$  = False Negatives
- **Example:** If a model predicts 80 correct and 20 incorrect out of 100 instances, the accuracy is 80%.

## Evaluating Model Performance - 2. Precision

- **Definition:** Proportion of positive identifications that were actually correct.

- **Formula:**

$$\text{Precision} = \frac{TP}{TP + FP} \quad (11)$$

- **Example:** If out of 30 predicted positives, 20 are true positives and 10 are false positives, then precision is  $\frac{20}{30} = 66.67\%$ .

## Evaluating Model Performance - 3. Recall (Sensitivity)

- **Definition:** Proportion of actual positives that were correctly identified.

- **Formula:**

$$\text{Recall} = \frac{TP}{TP + FN} \quad (12)$$

- **Example:** If there are 25 actual positive cases and the model identifies 20 correctly, recall is  $\frac{20}{25} = 80\%$ .

## Evaluating Model Performance - 4. F1 Score

- **Definition:** Harmonic mean of precision and recall, balancing both metrics.

- **Formula:**

$$F1 = 2 \times \frac{\text{Precision} \times \text{Recall}}{\text{Precision} + \text{Recall}} \quad (13)$$

- **Example:** With precision of 66.67% and recall of 80%, the F1 score is calculated as:

$$F1 \approx 72.73\% \quad (14)$$

## Evaluating Model Performance - 5. ROC-AUC

- **Definition:** Plots true positive rate (Recall) against false positive rate (FPR) at various thresholds.
- **Interpretation:**
  - $AUC = 1$ : Perfect model
  - $AUC > 0.5$ : Better than random guessing
  - $AUC = 0.5$ : No discrimination ability
- **Example:** If a model has an AUC of 0.85, it indicates strong performance in distinguishing between classes.

## Evaluating Model Performance - Key Points

- Different metrics provide insights into different aspects of model performance.
- Accuracy can be misleading, especially in imbalanced datasets.
- F1 score is preferred over accuracy when dealing with uneven class distributions.
- ROC-AUC is valuable for understanding the trade-offs between true and false positive rates.



## Evaluating Model Performance - Conclusion

Understanding and utilizing these performance metrics allows data scientists to better gauge their logistic regression model's effectiveness and make informed decisions for improvement or deployment.

# Interpreting Model Coefficients

Understanding coefficients in logistic regression is essential for understanding the relationships between predictors and outcomes.

# Understanding Coefficients in Logistic Regression

## ■ Logistic Regression Overview:

- A statistical method for predicting binary classes (0 or 1).
- Models the probability of a given input belonging to a specific category.

## ■ Coefficients Explanation:

- Indicate strength and direction of association between predictors and log-odds of the outcome.
- Logistic regression equation:

$$P(Y = 1|X) = \frac{1}{1 + e^{-(\beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_n X_n)}} \quad (15)$$

where  $\beta_0$  is the intercept, and  $\beta_i$  are coefficients for predictors  $X_i$ .

# Interpreting Coefficients

- **Positive Coefficient:**

- Indicates an increase in odds of the outcome.
- Example:  $\beta_1 = 0.5$  implies a log-odds increase of 0.5 for a unit increase in  $X_1$ .

- **Negative Coefficient:**

- Indicates a decrease in odds of the outcome.
- Example:  $\beta_2 = -0.7$  implies a log-odds decrease of 0.7 for a unit increase in  $X_2$ .

# Odds Ratio and Example

## ■ Odds Ratio:

$$\text{Odds Ratio} = e^{\beta_i} \quad (16)$$

For example, if  $\beta_1 = 0.5$ , then  $e^{0.5} \approx 1.65$ .

## ■ Example of Logistic Regression Model:

- Predicting whether a student will pass an exam based on hours studied and attendance.
- Coefficients:  $\beta_1 = 0.4$  (hours studied),  $\beta_2 = 1.2$  (attendance).

## ■ Interpretation:

- Hours studied: Odds ratio  $e^{0.4} \approx 1.49$  (49% more likely to pass).
- Attendance: Odds ratio  $e^{1.2} \approx 3.32$ , indicates significantly higher likelihood of passing with increased attendance.

# Conclusion

- Understanding the direction and magnitude of coefficients is crucial for interpreting the model.
- Coefficients and their odds ratios provide insights into predictor contributions to the outcome probability.
- Effective interpretation supports informed decision-making and understanding of data trends.

# Common Issues and Solutions in Logistic Regression - Introduction

Logistic regression is a powerful tool for binary classification, but it can encounter several common issues:

- Multicollinearity
- Overfitting
- Underfitting

Understanding these issues and how to address them is crucial for robust model development.

# Common Issues and Solutions - Multicollinearity

## Multicollinearity

- **Definition:** Occurs when independent variables are highly correlated.
- **Impact:** Inflates standard errors, leading to misleading significance tests.
- **Detection:** Variance Inflation Factor (VIF):

$$VIF = \frac{1}{1 - R^2} \quad (17)$$

- **Solutions:**
  - Remove one of the correlated variables.
  - Combine correlated variables (e.g., PCA).
  - Use Lasso (L1) regression as regularization.



# Common Issues and Solutions - Overfitting and Underfitting

## Overfitting

- **Definition:** Model learns noise in addition to the pattern.
- **Impact:** Good training performance, poor validation/testing performance.
- **Solutions:**
  - Cross-Validation (e.g., k-fold).
  - Regularization (L1 or L2).
  - Simplification of the model.

## Underfitting

- **Definition:** Model too simple to capture trends.
- **Impact:** Poor performance on training and test data.
- **Solutions:**
  - Increase model complexity (add polynomial terms).

# Common Issues and Solutions - Conclusion and Example Code

## Key Takeaways

- Monitor multicollinearity using VIF.
- Balance model complexity to avoid overfitting/underfitting.
- Use appropriate techniques like cross-validation and regularization.

## Example Code Snippet (Python)

```
1 import pandas as pd
2 import statsmodels.api as sm
3
4 # Calculate VIF
5 from statsmodels.stats.outliers_influence import variance_inflation_factor
6
7 def calculate_vif(X):
8     vif_data = pd.DataFrame()
```

# Use Cases of Logistic Regression

## Understanding Logistic Regression

Logistic Regression is a statistical method used for binary classification problems, where the outcome variable is categorical and typically takes on one of two values. It predicts the probability that a given input point belongs to a particular category.

# Real-world Applications

## 1 Healthcare

- Disease Diagnosis: Predicts presence or absence of a disease based on patient data, e.g., age, weight, blood pressure.
- *Example*: Predicting whether a patient has diabetes (yes/no) based on glucose levels and other features.

## 2 Finance

- Credit Scoring: Evaluates whether an applicant will default on a loan based on historical data.
- *Example*: Identifying whether a credit applicant will default based solely on credit history and income.

## 3 Marketing

- Customer Retention: Analyzes whether customers are likely to churn based on their interaction with products.
- *Example*: Assessing if a customer will renew their subscription based on usage patterns.

## 4 E-commerce

- Purchase Prediction: Forecasts whether a user will make a purchase after engaging with an advertisement.

## Key Points and Conclusion

### Key Points to Emphasize

- **Binary Outcome:** Logistic regression is suitable for problems with two outcome categories.
- **Interpretability:** Coefficients provide insightful information about the relationship between predictors and the outcome.
- **Probabilistic Framework:** Instead of predicting binary outcomes directly, it predicts the probability that an instance will fall into a category.

### Logistic Regression Formula

The logistic function is defined mathematically as:

$$P(Y = 1|X) = \frac{1}{1 + e^{-(\beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_n X_n)}} \quad (18)$$

# Advanced Topics in Logistic Regression

## Introduction to Regularization

Regularization is a technique used in logistic regression to prevent overfitting, especially when we have a large number of features. It adds a penalty to the loss function based on the magnitude of the coefficients.

# Lasso Regularization (L1)

- **Definition:** Adds the absolute value of the coefficients as a penalty term to the loss function.
- **Loss Function:**

$$J(\theta) = -\frac{1}{m} \sum_{i=1}^m [y^{(i)} \log(h_{\theta}(x^{(i)})) + (1 - y^{(i)}) \log(1 - h_{\theta}(x^{(i)}))] + \lambda \sum_{j=1}^n |\theta_j| \quad (19)$$

- **Key Features:**
  - Encourages sparsity: Reduces some coefficients to exactly zero, performing feature selection.
  - Useful in high-dimensional datasets.
- **Example:** Helps identify key features in datasets with many features by setting irrelevant coefficients to zero.

## Ridge Regularization (L2)

- **Definition:** Adds the squared coefficients as a penalty term to the loss function.

- **Loss Function:**

$$J(\theta) = -\frac{1}{m} \sum_{i=1}^m [y^{(i)} \log(h_{\theta}(x^{(i)})) + (1 - y^{(i)}) \log(1 - h_{\theta}(x^{(i)}))] + \lambda \sum_{j=1}^n \theta_j^2 \quad (20)$$

- **Key Features:**

- Does not perform variable selection; shrinks coefficients but keeps them non-zero.
- Helps in dealing with multicollinearity.

- **Example:** Stabilizes estimates in datasets with correlated features, improving predictions through all predictors.



## Comparing Lasso and Ridge

Feature	Lasso (L1)	Ridge (L2)
Coefficient Shrinkage	Can reduce coefficients to zero	Shrinks coefficients but not to zero
Feature Selection	Yes	No
Best for	High-dimensional data	Correlated variables

## Key Points and Conclusion

- Regularization helps manage overfitting by adjusting model complexity.
- Choosing between Lasso and Ridge depends on dataset needs (feature selection vs. multicollinearity).
- Regularization parameters ( $\lambda$ ) should be tuned with techniques like cross-validation.

### Conclusion

Integrating Lasso and Ridge in logistic regression leads to more robust and generalizable models, essential for effective data analysis in complex datasets.

# Practical Exercise

## Objective:

In this hands-on activity, you will implement **Logistic Regression** on a sample dataset to classify binary outcomes. This exercise will help you understand how to train a model, interpret its results, and evaluate its performance.

## Key Concepts

- **Logistic Regression** is a statistical method for predicting binary classes.
- The outcome is modeled as a function of the independent variables using the logistic function.
- The formula for logistic regression is:

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

- $P(Y = 1|X)$ : Probability of the positive class (class 1)
- $\beta_0$ : Intercept
- $\beta_1, \beta_2, \dots, \beta_n$ : Coefficients for each independent variable  $X_1, X_2, \dots, X_n$

# Steps for Implementation

## 1. Load the Dataset:

Use a CSV file or a preloaded dataset available in your Python environment. For this exercise, we will use the **Iris dataset**.

```
1 import pandas as pd
2 from sklearn.datasets import load_iris
3
4 data = load_iris()
5 df = pd.DataFrame(data.data, columns=data.feature_names)
6 df['target'] = (data.target == 0).astype(int) # 1 for Setosa, 0 for
        others
```

## 2. Explore the Data:

Perform a preliminary analysis to understand the features and target variable.

```
1 print(df.head())
2 print(df.describe())
```

## Steps for Implementation (cont.)

### 3. Split the Data:

Divide the dataset into training and testing sets.

```
1 from sklearn.model_selection import train_test_split
2
3 X = df[data.feature_names]
4 y = df['target']
5 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2,
    random_state=42)
```

### 4. Create and Train the Model:

Utilize LogisticRegression from sklearn.

```
1 from sklearn.linear_model import LogisticRegression
2
3 model = LogisticRegression()
4 model.fit(X_train, y_train)
```

## Steps for Implementation (cont.)

### 5. Make Predictions:

Use the model to predict outcomes on the test set.

```
1 y_pred = model.predict(X_test)
```

### 6. Evaluate the Model:

Assess performance using metrics like accuracy, precision, and the confusion matrix.

```
1 from sklearn.metrics import accuracy_score, confusion_matrix
2
3 accuracy = accuracy_score(y_test, y_pred)
4 cm = confusion_matrix(y_test, y_pred)
5
6 print(f"Accuracy: {accuracy}")
7 print(f"Confusion Matrix:\n{cm}")
```

## Key Points to Emphasize

- Logistic regression is suited for binary classification problems.
- Properly splitting your data into training and testing sets is crucial for unbiased evaluation.
- Understanding evaluation metrics helps assess model performance beyond mere accuracy.



## Conclusion and Next Steps

### **Conclusion:**

By completing this exercise, you will acquire practical skills in applying logistic regression, understand its mechanics, and appreciate the importance of model evaluation.

### **Next Steps:**

Prepare for our concluding slide, encapsulating your key learnings and suggesting areas for further exploration in supervised learning concepts.

## Conclusion - Key Takeaways

As we wrap up our exploration of Logistic Regression, let's summarize the key points and encourage further exploration into supervised learning.

- 1 Logistic Regression Definition:** A statistical method for binary classification predicting the probability of a categorical dependent variable based on independent variables.
- 2 Sigmoid Function:** The logistic function transforms any real-valued number into a value between 0 and 1, suitable for estimating probabilities.
- 3 Interpretation of Coefficients:** Coefficients ( $\beta$ ) represent the log odds of the output being 1 (positive class) relative to input features  $X$ .
- 4 Loss Function:** Using Log Loss (Cross-Entropy Loss) to evaluate performance, minimizing the difference between predicted probabilities and actual outcomes.

## Conclusion - Practical Applications

Logistic regression has numerous practical applications across various fields:

- **Medicine:** Predicting disease presence.
- **Marketing:** Customer conversion analysis.
- **Social Sciences:** Understanding voting behavior.

Moreover, there are additional learning opportunities to further develop your understanding.

## Conclusion - Encouragement to Explore Further

This concludes our fundamental overview of Logistic Regression, a core technique in supervised learning. We encourage you to:

- Experiment with different datasets using the concepts learned today.
- Implement and analyze model performance using metrics such as accuracy, precision, and recall.
- Engage with advanced literature and online courses to expand your knowledge base and capabilities in machine learning.

Thank you for participating this week; your journey into supervised learning is just beginning!