

VIETNAM NATIONAL UNIVERSITY HO CHI MINH CITY  
HO CHI MINH CITY UNIVERSITY OF TECHNOLOGY  
FACULTY OF COMPUTER SCIENCE AND ENGINEERING



## Probability and Statistics

---

### Final Project Report

# Internet Advertisement Classification using Random Forest

---

Advisor(s): [Instructor Name]

Student(s):	NGUYỄN LÊ THANH HÀO	1952242
	PHẠM NHẬT KHÔI	2352623
	Nguyễn Song Đạt	1952646
	Nguyễn Phước Thịnh	1852765

HO CHI MINH CITY, AUGUST 2025





## Contents

<b>1</b>	<b>Introduction</b>	<b>5</b>
<b>2</b>	<b>Data Description</b>	<b>6</b>
2.1	Dataset Overview . . . . .	6
2.2	Data Loading . . . . .	6
2.3	Target Variable Distribution . . . . .	6
2.4	Missing Values . . . . .	6
2.5	Data Preprocessing . . . . .	7
2.6	Final Dataset Characteristics . . . . .	8
2.7	Basic Statistics . . . . .	8
<b>3</b>	<b>Descriptive Statistics</b>	<b>9</b>
3.1	Target Variable Analysis . . . . .	9
3.2	Feature Distribution Analysis . . . . .	10
3.3	Class-wise Feature Analysis . . . . .	12
3.4	Feature Relationships and Scatter Analysis . . . . .	12
3.5	Correlation Analysis . . . . .	14
3.6	Key Findings from Descriptive Analysis . . . . .	15
<b>4</b>	<b>Objective and Methodology</b>	<b>16</b>
4.1	Project Objective . . . . .	16
4.2	Theoretical Foundation . . . . .	16
4.2.1	Binary Classification Overview . . . . .	16
4.2.2	Evaluation Metrics . . . . .	16
4.2.3	Key Concepts . . . . .	17
4.3	Machine Learning Methods . . . . .	17
4.3.1	k-Nearest Neighbors (k-NN) . . . . .	17
4.3.2	Decision Tree . . . . .	19
4.3.3	Random Forest . . . . .	20
4.4	Our Approach . . . . .	21
4.5	Expected Results . . . . .	22

## List of Figures

3.1	Distribution of target variable showing class imbalance . . . . .	9
-----	---	---



3.2 Histograms of the first six features showing distribution patterns . . . . . 11

3.3 Boxplots of features 1 and 2 grouped by target class . . . . . 12

3.4 Scatter plots and density plots showing feature relationships and class distributions . . . . . 13

3.5 Correlation heatmap of the first 10 features . . . . . 14

4.1 k-NN Classification Example with k=3 Nearest Neighbors . . . . . 18

4.2 Decision Tree Example for Ad Classification . . . . . 19

4.3 Random Forest: Ensemble of Decision Trees for Ad Classification . . . . . 20

List of Tables

3.1 Summary statistics for the first 10 features . . . . . 10

Listings



# 1 Introduction

The proliferation of the internet has led to a massive increase in online advertising. While advertisements can be useful, they can also be intrusive and detract from the user experience. Therefore, the ability to automatically distinguish between advertising and non-advertising content is a significant challenge in modern data science. This project addresses this challenge by developing a model to classify internet images as either advertisements ('ad') or non-advertisements ('nonad').

This report details the process of analyzing the "Internet Advertisements Data Set" sourced from the UCI Machine Learning Repository<sup>[8]</sup>. We will explore the data, clean it for analysis, and apply several machine learning models to build an effective classifier. The primary goal is to construct a robust model that can accurately predict whether an image is an advertisement based on its various features.

## 2 Data Description

### 2.1 Dataset Overview

The dataset used in this project is the "Internet Advertisements Data Set" from the UCI Machine Learning Repository. This dataset contains information about internet images and their classification as advertisements or non-advertisements.

The original dataset consists of 3,279 observations with 1,560 columns. The first column serves as an index and is removed during preprocessing, leaving 1,559 features for analysis. The target variable, located in the last column (X1558), indicates whether an image is an advertisement ('ad.') or not ('nonad.').

### 2.2 Data Loading

The dataset is loaded from a CSV file using the following R code:

```
1 # Load the dataset from the CSV file
2 data <- read.csv("../add.csv", header = TRUE,
3                   stringsAsFactors = FALSE)
4 # Identify the target column name
5 target_col <- names(data)[ncol(data)]
```

### 2.3 Target Variable Distribution

The target variable shows a significant class imbalance:

- Advertisement images (ad): 459 observations (14%)
- Non-advertisement images (nonad): 2,820 observations (86%)

This imbalance is typical in advertisement detection problems and will need to be considered when building and evaluating classification models.

### 2.4 Missing Values

The dataset contains missing values represented as "?" strings. Initial analysis revealed:

- Total missing values: 15 occurrences



- All missing values are concentrated in column 5
- Missing values represent approximately 0.46% of column 5's data

## 2.5 Data Preprocessing

Several preprocessing steps were applied to prepare the data for analysis:

```
1 # 1. Remove the first column which is an unnecessary index
2 data <- data[, -1]
3
4 # 2. Handle missing values represented by "?"
5 # Convert "?" to NA for all feature columns
6 for(i in 1:(ncol(data)-1)) {
7   data[[i]] <- as.numeric(ifelse(data[[i]] == "?", NA, data[[
8     i]]))
9 }
10
11 # 3. Impute missing values using the median of each column
12 for(i in 1:(ncol(data)-1)) {
13   if(any(is.na(data[[i]]))) {
14     median_val <- median(data[[i]], na.rm = TRUE)
15     data[[i]][is.na(data[[i]])] <- median_val
16   }
17 }
18
19 # 4. Clean and format the target variable
20 data[[target_col]] <- factor(data[[target_col]],
21                               levels = c("ad.", "nonad."),
22                               labels = c("ad", "nonad"))
```

The preprocessing pipeline includes:

1. **Index removal:** The first column containing row indices was removed
2. **Missing value handling:** All "?" strings were converted to NA values
3. **Median imputation:** Missing values were replaced with the median of their respective columns

4. **Target variable formatting:** The target variable was converted to a factor with clear labels ("ad" and "nonad")

## 2.6 Final Dataset Characteristics

After preprocessing, the cleaned dataset has the following properties:

- Dimensions: 3,279 rows  $\times$  1,559 columns
- All missing values have been imputed
- Target variable is properly formatted as a factor
- All feature columns are numeric
- Class distribution remains: 459 advertisements, 2,820 non-advertisements

## 2.7 Basic Statistics

Preliminary statistics for the first five feature columns show:

- Column 1: Mean = 1,639, Median = 1,639, SD = 946.71
- Column 2: Mean = 64.02, Median = 51, SD = 54.87
- Column 3: Mean = 155.34, Median = 110, SD = 130.03
- Column 4: Mean = 3.91, Median = 2.1, SD = 6.04
- Column 5: Mean = 0.77, Median = 1, SD = 0.42

These statistics indicate varying scales across features, suggesting that normalization or standardization may be beneficial for certain machine learning algorithms.



## 3 Descriptive Statistics

This section presents a comprehensive exploratory data analysis (EDA) of the Internet Advertisements dataset to understand the data distribution, relationships between variables, and key characteristics that will inform our modeling approach.

### 3.1 Target Variable Analysis

The target variable distribution reveals a significant class imbalance in the dataset:

- Advertisement images (ad): 459 observations (14%)
- Non-advertisement images (nonad): 2,820 observations (86%)

This 1:6 ratio between advertisement and non-advertisement classes is typical in real-world advertisement detection scenarios and must be considered when evaluating model performance.

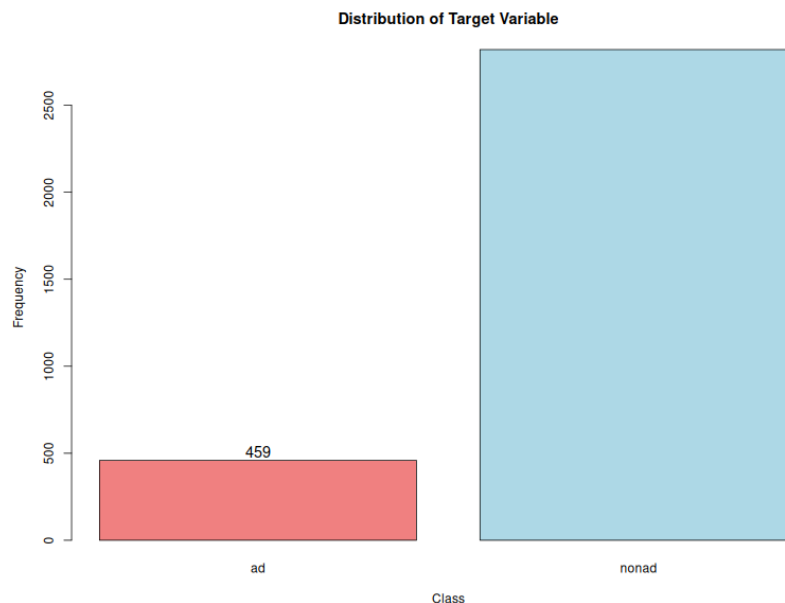


Figure 3.1: Distribution of target variable showing class imbalance

The target distribution visualization was generated using:

```
1 # Bar plot for target variable distribution
2 target_table <- table(data[[target_col]])
3 barplot(target_table,
```



```
4     main = "Distribution of Target Variable",  
5     xlab = "Class", ylab = "Frequency",  
6     col = c("lightcoral", "lightblue"),  
7     border = "black")
```

## 3.2 Feature Distribution Analysis

To understand the characteristics of individual features, we analyzed the distribution of the first 10 features in the dataset. The summary statistics reveal varying scales and distributions across features:

Table 3.1: Summary statistics for the first 10 features

Feature	Min	Q1	Median	Mean	Q3	Max	SD
X0	1.00	32.50	51.00	60.44	61.00	640	47.06
X1	1.00	90.00	110.00	142.89	144.00	640	112.56
X2	0.00	1.28	2.10	3.41	3.90	60	5.20
X3	0.00	1.00	1.00	0.77	1.00	1	0.42
X4	0.00	0.00	0.00	0.004	0.00	1	0.065

Histograms for the first six features show diverse distribution patterns:

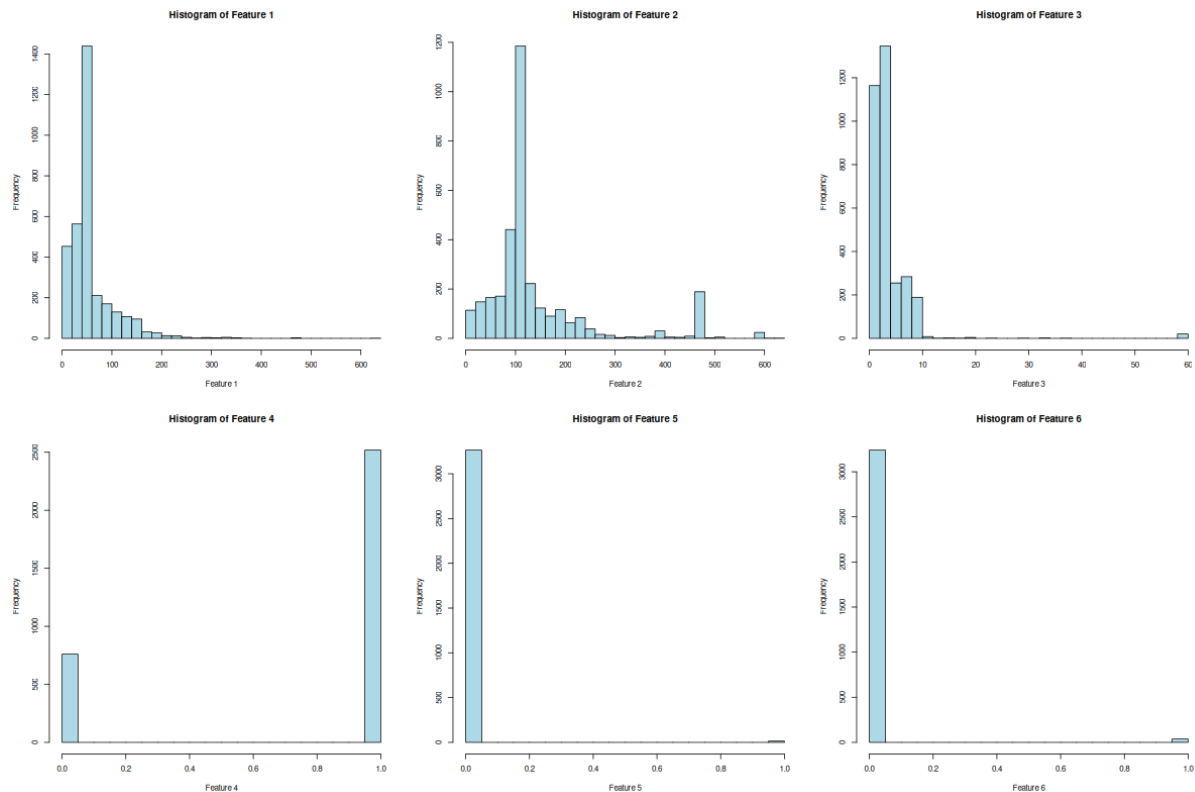


Figure 3.2: Histograms of the first six features showing distribution patterns

The histogram generation code:

```
1 # Histograms for the first 6 features
2 par(mfrow = c(2, 3))
3 for(i in 1:6) {
4   hist(data[[i]],
5         main = paste("Histogram of Feature", i),
6         xlab = paste("Feature", i),
7         ylab = "Frequency",
8         col = "lightblue",
9         border = "black",
10        breaks = 30)
11 }
```

### 3.3 Class-wise Feature Analysis

To understand how features differ between advertisement and non-advertisement classes, we created boxplots comparing the distribution of key features across classes:

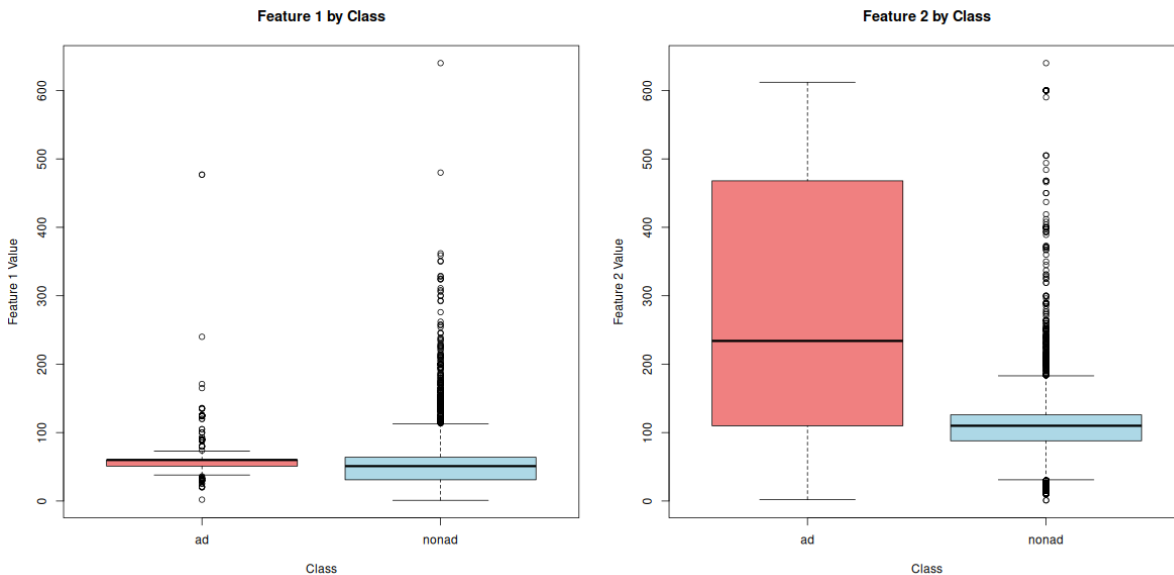


Figure 3.3: Boxplots of features 1 and 2 grouped by target class

The boxplot analysis was implemented as:

```
1 # Boxplots of features 1 and 2, grouped by target class
2 par(mfrow = c(1, 2))
3 boxplot(data[[1]] ~ data[[target_col]],
4         main = "Feature 1 by Class",
5         xlab = "Class", ylab = "Feature 1 Value",
6         col = c("lightcoral", "lightblue"))
7 boxplot(data[[2]] ~ data[[target_col]],
8         main = "Feature 2 by Class",
9         xlab = "Class", ylab = "Feature 2 Value",
10        col = c("lightcoral", "lightblue"))
```

### 3.4 Feature Relationships and Scatter Analysis

Scatter plots and density plots provide insights into feature relationships and class separability:

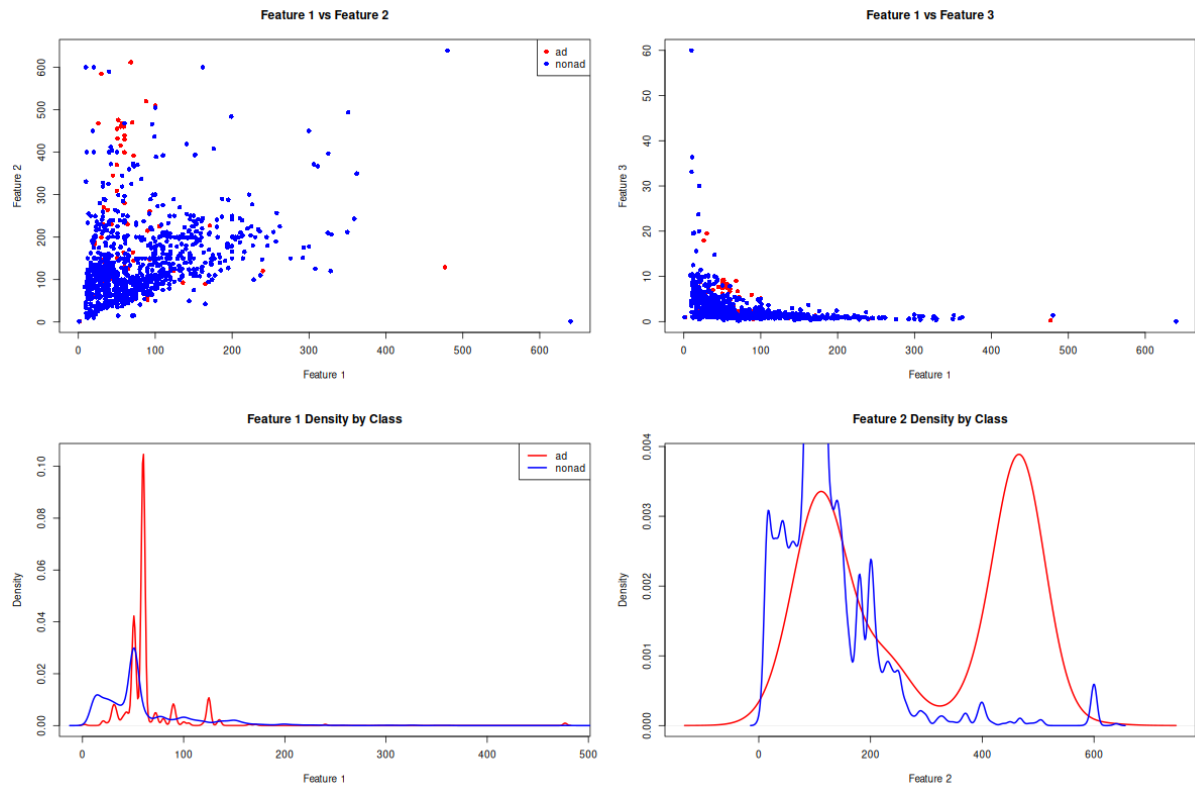


Figure 3.4: Scatter plots and density plots showing feature relationships and class distributions

The visualization combines scatter plots and density analysis:

```
1 # Scatter plots and density plots
2 par(mfrow = c(2, 2))
3 colors <- c("red", "blue")
4 class_colors <- colors[as.numeric(data[[target_col]])]
5 plot(data[[1]], data[[2]],
6       main = "Feature 1 vs Feature 2",
7       xlab = "Feature 1", ylab = "Feature 2",
8       col = class_colors, pch = 16)
9 legend("topright", legend = levels(data[[target_col]]),
10       col = colors, pch = 16)
```

### 3.5 Correlation Analysis

Correlation analysis reveals important relationships between features. The correlation heatmap for the first 10 features shows the strength of linear relationships:

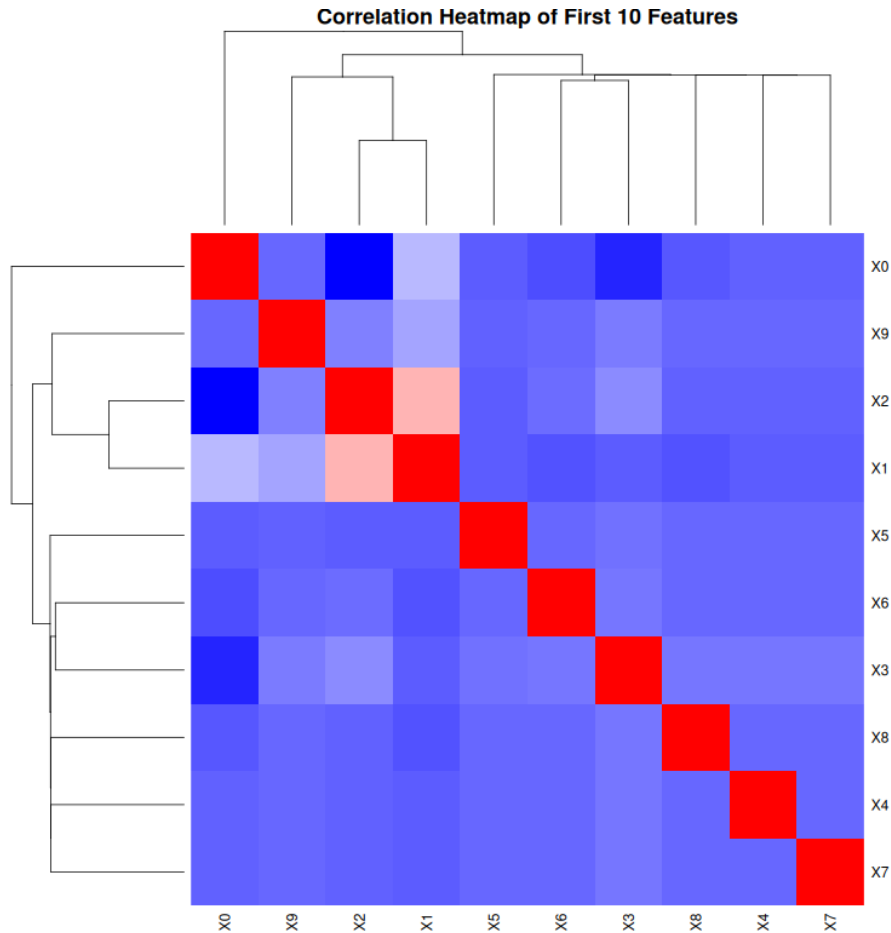


Figure 3.5: Correlation heatmap of the first 10 features

The correlation analysis identified several perfect correlations (correlation = 1.0) between different feature pairs, indicating potential redundancy in the dataset:

- Features X11 and X14: Perfect positive correlation
- Features X8 and X15: Perfect positive correlation
- Features X13 and X38: Perfect positive correlation
- Features X44 and X46: Perfect positive correlation

The correlation matrix was computed using:

```
1 # Correlation heatmap for the first 10 numeric features
2 numeric_data <- data[, sapply(data, is.numeric)]
3 cor_matrix <- cor(numeric_data[, 1:10], use = "complete.obs")
4 heatmap(cor_matrix,
5         symm = TRUE,
6         main = "Correlation Heatmap of First 10 Features",
7         col = colorRampPalette(c("blue", "white", "red"))
           (100))
```

### 3.6 Key Findings from Descriptive Analysis

The exploratory data analysis reveals several important characteristics:

1. **Class Imbalance:** The dataset has a significant imbalance with 86% non-advertisements
2. **Feature Diversity:** Features show diverse scales and distributions, suggesting the need for normalization
3. **Perfect Correlations:** Multiple feature pairs show perfect correlation, indicating potential redundancy
4. **Class Separability:** Some features show different distributions between advertisement and non-advertisement classes
5. **High Dimensionality:** With 1,559 features, dimensionality reduction techniques may be beneficial

These findings will inform our modeling approach, particularly regarding feature selection, data preprocessing, and evaluation metrics that account for class imbalance.

## 4 Objective and Methodology

### 4.1 Project Objective

This project develops machine learning models to automatically classify internet images as advertisements or non-advertisements. This binary classification addresses the challenge of distinguishing advertising content from regular content, which is essential for:

- **Content filtering:** Automatically detecting advertisement content
- **User experience:** Reducing intrusive advertising
- **Digital analysis:** Understanding advertisement patterns
- **Content moderation:** Supporting automated systems

With 1,559 features and class imbalance in our dataset, we aim to:

1. Build models that handle high-dimensional data effectively
2. Compare different machine learning approaches
3. Identify key features for advertisement classification
4. Provide practical implementation insights

### 4.2 Theoretical Foundation

#### 4.2.1 Binary Classification Overview

Our task is a **binary classification problem** where we predict one of two classes: advertisement (1) or non-advertisement (0). Given input features  $\mathbf{x}$  (1,559 dimensions), we want to learn a function  $f(\mathbf{x}) \rightarrow \{0, 1\}$  that makes accurate predictions.

**Goal:** Minimize classification errors by learning patterns from training data.

#### 4.2.2 Evaluation Metrics

We use multiple metrics to assess model performance, especially important for imbalanced datasets:



**Accuracy:** Overall correctness

$$Accuracy = \frac{\text{Correct Predictions}}{\text{Total Predictions}} = \frac{TP + TN}{TP + TN + FP + FN} \quad (4.1)$$

**Precision:** How many predicted ads are actually ads

$$Precision = \frac{TP}{TP + FP} \quad (4.2)$$

**Recall:** How many actual ads we correctly identified

$$Recall = \frac{TP}{TP + FN} \quad (4.3)$$

**F1-Score:** Balance between Precision and Recall

$$F1 = 2 \cdot \frac{Precision \times Recall}{Precision + Recall} \quad (4.4)$$

where TP=True Positives, TN=True Negatives, FP=False Positives, FN=False Negatives.

### 4.2.3 Key Concepts

**Bias-Variance Tradeoff:** Every model faces a balance between:

- **Bias:** Error from oversimplifying (underfitting)
- **Variance:** Error from being too sensitive to training data (overfitting)

**Cross-Validation:** We split data into k parts, train on k-1 parts, test on 1 part, repeat k times. This gives reliable performance estimates and helps select best parameters.

## 4.3 Machine Learning Methods

We employ three different approaches, each representing distinct learning paradigms:

### 4.3.1 k-Nearest Neighbors (k-NN)

**Core Idea:** Classify new data points based on the majority class of their k closest neighbors.

**How it works:**

1. Store all training data (**lazy learning**)

2. For a new point, find  $k$  nearest neighbors using distance
3. Assign the most common class among these  $k$  neighbors

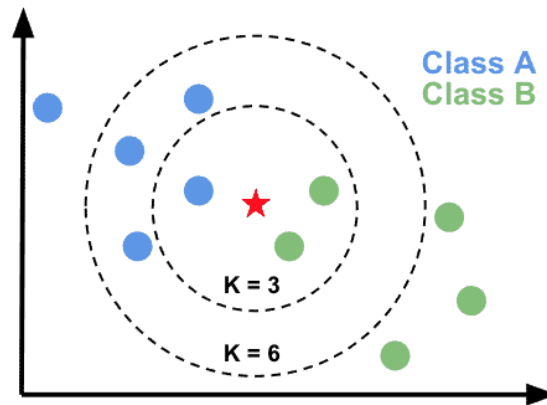


Figure 4.1:  $k$ -NN Classification Example with  $k=3$  Nearest Neighbors

#### Key Distance Formulas:

- **Euclidean:**  $d = \sqrt{\sum_{i=1}^p (x_i - y_i)^2}$  (straight-line distance)
- **Manhattan:**  $d = \sum_{i=1}^p |x_i - y_i|$  (city-block distance)

#### Advantages:

- Simple to understand and implement
- No assumptions about data distribution
- Works well with sufficient training data

#### Challenges:

- Slow prediction (must check all training points)
- Sensitive to irrelevant features
- Struggles with high dimensions (curse of dimensionality)

### 4.3.2 Decision Tree

**Core Idea:** Create a tree of yes/no questions to classify data, like a flowchart.

**How it works:**

1. Start with all training data at the root
2. Find the best feature and threshold to split data
3. Repeat for each branch until stopping criteria met
4. Make predictions by following the path from root to leaf

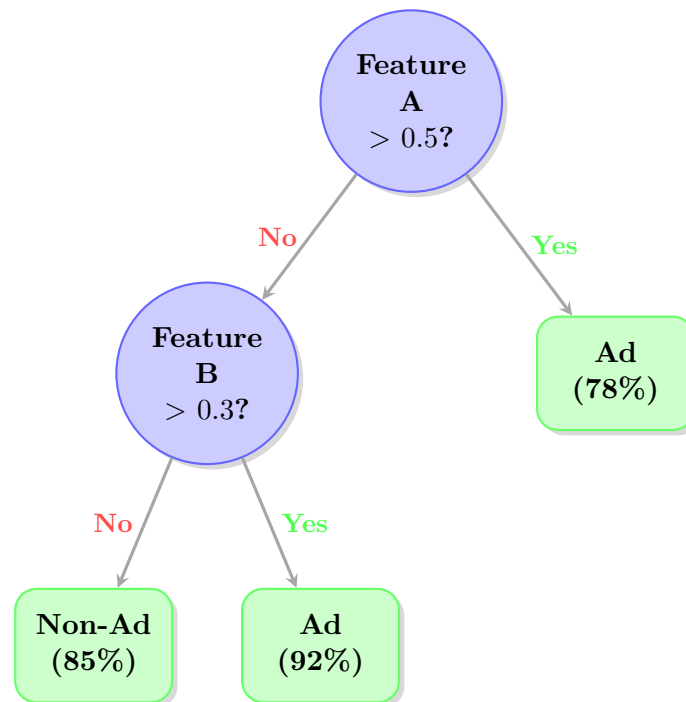


Figure 4.2: Decision Tree Example for Ad Classification

**Key Concepts:**

- **Gini Impurity:** Measures how "mixed" classes are in a node

$$Gini = 1 - \sum_{i=1}^c p_i^2 \quad (4.5)$$

- **Information Gain:** How much a split reduces uncertainty
- **Pruning:** Removing branches to prevent overfitting

**Advantages:**

- Easy to understand and interpret (white box)
- Handles both numerical and categorical features
- No need for feature scaling
- Automatically selects important features

**Challenges:**

- Prone to overfitting (memorizing training data)
- Unstable (small data changes = different trees)
- Can create overly complex rules

#### 4.3.3 Random Forest

**Core Idea:** Combine many decision trees to make better predictions (**ensemble method**).

**How it works:**

1. Create many different training datasets using **bootstrap sampling** (random sampling with replacement)
2. Train one decision tree on each dataset
3. For each tree, use only a random subset of features at each split
4. Combine all tree predictions by **majority voting**

**Key Features:**

- **Bootstrap Sampling:** Each tree sees different data
- **Feature Randomness:** Each split uses random feature subset
- **Majority Voting:** Final prediction = most common prediction

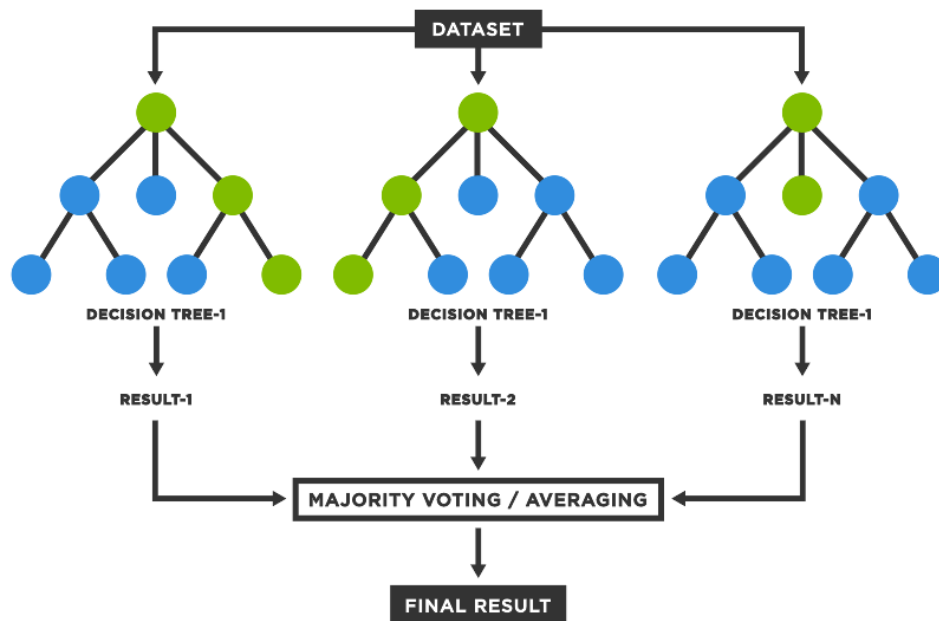


Figure 4.3: Random Forest: Ensemble of Decision Trees for Ad Classification

- **Out-of-Bag Error:** Built-in validation using unused samples

**Why it works better:**

- Individual trees may overfit, but averaging reduces this
- Different trees make different mistakes
- Combines strengths while canceling weaknesses

**Advantages:**

- Excellent performance with high-dimensional data
- Provides feature importance automatically
- Robust to overfitting
- Handles missing data well
- Less parameter tuning needed

### **Challenges:**

- Less interpretable than single trees (black box)
- Computationally expensive with many trees
- Memory intensive for large datasets

## **4.4 Our Approach**

We follow a systematic process:

### **1. Data Preparation:**

- Normalize features (important for k-NN)
- Split data: 70% training, 30% testing
- Handle class imbalance

### **2. Model Training:**

- k-NN: Find best k value using cross-validation
- Decision Tree: Use pruning to prevent overfitting
- Random Forest: Tune number of trees and features per split

### **3. Evaluation:**

- Compare using Accuracy, Precision, Recall, F1-score
- Analyze confusion matrices
- Examine feature importance (from Random Forest)

## **4.5 Expected Results**

This study will:

- Identify the best method for ad classification
- Discover which features are most important for detecting ads
- Provide practical recommendations for real-world systems
- Show how ensemble methods handle high-dimensional data

By comparing three different approaches, we ensure robust conclusions about the most effective techniques for internet advertisement classification.



## References

- [1] Leo Breiman. Random forests. *Machine learning*, 45(1):5–32, 2001.
- [2] Thomas Cover and Peter Hart. Nearest neighbor pattern classification. *IEEE transactions on information theory*, 13(1):21–27, 1967.
- [3] GeeksforGeeks. Random forest algorithm in machine learning, 2025. Accessed: 2025.
- [4] GeeksforGeeks. What are the advantages and disadvantages of random forest?, 2025. Accessed: 2025.
- [5] Trevor Hastie, Robert Tibshirani, and Jerome Friedman. *The elements of statistical learning: data mining, inference, and prediction*. Springer Science & Business Media, 2009.
- [6] IBM. What is random forest?, 2024. Accessed: 2024.
- [7] Built In. Random forest: A complete guide for machine learning, 2024. Accessed: 2024.
- [8] M. Lichman. UCI machine learning repository, 2013.
- [9] J Ross Quinlan. *C4. 5: programs for machine learning*. Morgan kaufmann, 1993.
- [10] SkillCamper. Random forest: Why ensemble learning outperforms individual models, 2024. Accessed: 2024.