

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



## Probability and Statistics

---

### Assignment Project Report

# Internet Advertisement Classification using Machine Learning Models

---

Advisor(s):	Dr. Nguyen Tien Dung	
Student(s):	NGUYEN LE THANH HAO	1952242
	PHAM NHAT KHOI	2352623
	Nguyen Song Dat	1952646
	Nguyen Phuoc Thinh	1852765

HO CHI MINH CITY, AUGUST 2025





## Member list & Workload

No.	Full name	Student ID	Contribution
1	NGUYEN LE THANH HAO	1952242	100%
2	PHAM NHAT KHOI	2352623	100%
3	Nguyen Song Dat	1952646	100%
4	Nguyen Phuoc Thinh	1852765	100%



# Contents

<b>1</b>	<b>Introduction</b>	<b>8</b>
<b>2</b>	<b>Data Description</b>	<b>8</b>
2.1	Raw Dataset Overview . . . . .	8
2.2	Raw Data Loading Process . . . . .	9
2.3	Raw Data Target Distribution . . . . .	9
2.4	Missing Values in Raw Data . . . . .	9
2.5	Data Cleaning and Preprocessing Pipeline . . . . .	10
2.6	Clean Data Characteristics . . . . .	10
2.7	Clean Data Statistical Summary . . . . .	11
2.8	Data Cleaning Process Summary . . . . .	11
<b>3</b>	<b>Descriptive Statistics</b>	<b>12</b>
3.1	Target Variable Analysis . . . . .	12
3.2	Feature Distribution Analysis . . . . .	13
3.3	Class-wise Feature Analysis . . . . .	14
3.4	Feature Relationships and Scatter Analysis . . . . .	14
3.5	Correlation Analysis . . . . .	15
3.6	Key Findings from Descriptive Analysis . . . . .	17
<b>4</b>	<b>Objective and Methodology</b>	<b>18</b>
4.1	Project Objective . . . . .	18
4.2	Theoretical Foundation . . . . .	18
4.2.1	Binary Classification Overview . . . . .	18
4.2.2	Evaluation Metrics . . . . .	18
4.2.3	Key Concepts . . . . .	19
4.3	Machine Learning Methods . . . . .	19
4.3.1	k-Nearest Neighbors (k-NN) . . . . .	19
4.3.2	Decision Tree . . . . .	20
4.3.3	Random Forest . . . . .	21
4.4	Our Approach . . . . .	23
4.5	Expected Results . . . . .	23
<b>5</b>	<b>k-Nearest Neighbors Model</b>	<b>24</b>
5.1	Implementation . . . . .	24
5.1.1	Data Preparation . . . . .	24
5.1.2	k-NN Implementation . . . . .	24
5.2	Hyperparameter Tuning . . . . .	25
5.2.1	k-Value Selection . . . . .	25
5.3	Model Evaluation . . . . .	25
5.3.1	Performance Metrics . . . . .	25
5.3.2	Confusion Matrix Analysis . . . . .	26
5.4	Results Interpretation . . . . .	27
5.4.1	Strengths . . . . .	27
5.4.2	Limitations . . . . .	27
5.4.3	Class Imbalance Impact . . . . .	27
5.5	Conclusion . . . . .	27



<b>6</b>	<b>Decision Tree Model</b>	<b>27</b>
6.1	Algorithm Implementation	27
6.1.1	Data Preparation	27
6.1.2	Gini Impurity Function	28
6.2	Best Split Finding Function	28
6.2.1	Tree Building Function	28
6.3	Model Evaluation	29
6.3.1	Training Results	29
6.3.2	Confusion Matrix	29
6.3.3	Evaluation Metrics	30
6.4	Decision Tree Structure	30
6.5	Results Analysis	30
6.5.1	Model Strengths	30
6.5.2	Model Limitations	31
6.5.3	Impact of Class Imbalance	31
6.6	Conclusion	32
<b>7</b>	<b>Random Forest Model</b>	<b>33</b>
7.1	Introduction	33
7.2	Algorithm Implementation	33
7.2.1	Data Preparation	33
7.2.2	Bootstrap Sampling Function	34
7.2.3	Random Feature Selection	34
7.2.4	Forest Construction	34
7.3	Model Evaluation	34
7.3.1	Training Results	34
7.3.2	Confusion Matrix	35
7.3.3	Evaluation Metrics	35
7.3.4	Feature Importance	35
7.4	Results Analysis	36
7.4.1	Model Strengths	36
7.4.2	Model Limitations	36
7.4.3	Impact of Class Imbalance	37
7.5	Conclusion	37
<b>8</b>	<b>Model Comparison and Analysis</b>	<b>38</b>
8.1	Performance Metrics Comparison	38
8.2	Best Models by Metric	38
8.3	Confusion Matrix Analysis	39
8.3.1	k-NN Model	39
8.3.2	Decision Tree Model	40
8.3.3	Random Forest Model	40
8.4	Model Characteristics Analysis	40
8.4.1	k-NN Model Characteristics	40
8.4.2	Decision Tree Characteristics	40
8.4.3	Random Forest Characteristics	41
8.5	Class Imbalance Impact	41
8.6	Model Recommendations	41
8.6.1	Best Overall Model: Random Forest	41



8.6.2	Most Interpretable: Decision Tree . . . . .	42
8.6.3	Simplest Approach: k-NN . . . . .	42
8.7	Conclusion . . . . .	42
<b>9</b>	<b>Conclusion</b> . . . . .	<b>43</b>
9.1	Project Summary . . . . .	43
9.2	Dataset Characteristics . . . . .	43
9.3	Key Findings . . . . .	43
9.3.1	Model Performance Comparison . . . . .	43
9.3.2	Best Performing Models by Metric . . . . .	44
9.4	Algorithm Analysis . . . . .	44
9.4.1	k-Nearest Neighbors (k=3) . . . . .	44
9.4.2	Decision Tree (max_depth=3) . . . . .	44
9.4.3	Random Forest (n_trees=100) . . . . .	45
9.5	Practical Implications . . . . .	45
9.5.1	Application-Specific Recommendations . . . . .	45
9.5.2	Class Imbalance Considerations . . . . .	45
9.6	Research Contributions . . . . .	46
9.6.1	Technical Achievements . . . . .	46
9.6.2	Statistical Significance . . . . .	46
9.7	Limitations . . . . .	46
9.7.1	Current Limitations . . . . .	46
9.8	Final Conclusions . . . . .	47
9.8.1	Project Success . . . . .	47
9.8.2	Key Insights . . . . .	47
9.8.3	Impact and Significance . . . . .	47
9.8.4	Closing Remarks . . . . .	48

## List of Figures

3.1	Distribution of target variable showing class imbalance . . . . .	12
3.2	Histograms of the first six features showing distribution patterns . . . . .	13
3.3	Boxplots of features 1 and 2 grouped by target class . . . . .	14
3.4	Scatter plots and density plots showing feature relationships and class distributions . . . . .	15
3.5	Correlation heatmap of the first 10 features . . . . .	16
4.1	k-NN Classification Example with k=3 Nearest Neighbors . . . . .	20
4.2	Decision Tree Example for Ad Classification . . . . .	21
4.3	Random Forest Diagram . . . . .	22
5.1	k-NN accuracy vs k value showing optimal performance at k=3 . . . . .	25
5.2	k-NN confusion matrix showing classification results . . . . .	26
6.1	Decision Tree Confusion Matrix Visualization . . . . .	31
6.2	Decision tree structure with maximum depth of 3 . . . . .	32
6.3	Decision Tree Feature Importance . . . . .	33
7.1	Random Forest Confusion Matrix Visualization . . . . .	35
7.2	Random Forest Feature Importance Visualization . . . . .	37
8.1	Visual Comparison of Model Performance Metrics . . . . .	39
8.2	Confusion Matrices for All Three Models . . . . .	39



## List of Tables

2.1	Descriptive Statistics for First Five Feature Columns . . . . .	11
3.1	Summary statistics for the first 10 features . . . . .	13
5.1	k-NN performance for different k values . . . . .	25
5.2	k-NN confusion matrix (k=3) . . . . .	26
6.1	Decision Tree Model Training Results . . . . .	30
6.2	Confusion Matrix - Decision Tree . . . . .	30
6.3	Decision Tree Model Evaluation Metrics . . . . .	30
7.1	Random Forest Model Training Results . . . . .	34
7.2	Confusion Matrix - Random Forest . . . . .	35
7.3	Random Forest Performance Metrics . . . . .	36
7.4	Top 10 Feature Importance - Random Forest . . . . .	36
8.1	Model Performance Comparison . . . . .	38
9.1	Final Model Performance Summary . . . . .	44

## Listings

6.1	Gini Impurity Function . . . . .	28
6.2	Best Split Finding Function . . . . .	28
6.3	Decision Tree Building Function . . . . .	29
7.1	Bootstrap Sampling Implementation . . . . .	34
7.2	Random Feature Selection . . . . .	34
7.3	Random Forest Training . . . . .	34

# 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 Raw Dataset Overview

The dataset for this project is the "Internet Advertisements Data Set" from the UCI Machine Learning Repository. It contains 3,279 instances, where the task is to classify each image as an advertisement ("ad") or not ("nonad").

Based on the official documentation, the 1,559 features provided for classification can be summarized into three main categories:

- **Geometric Properties:** Three continuous features representing the image's height, width, and aspect ratio.
- **Phrase-based Features:** Over 1,500 binary features indicating the presence of specific advertising-related words or phrases in the image's URL, alt text, and surrounding text.
- **URL-based Features:** A smaller set of binary features derived from keywords within the page and image URLs.

The target variable, indicating the class, is located in the final column of the dataset.

#### **Initial Dataset Characteristics:**

- Total samples: 3,279 internet images
- Original features: 1,560 columns (including index and target)
- Feature columns: X0 to X1557 (1,558 numerical attributes)
- Target column: X1558 with values 'ad.' and 'nonad.'
- File format: CSV with headers



## 2.2 Raw Data Loading Process

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, stringsAsFactors = FALSE)
3
4 # Identify the target column name
5 target_col <- names(data)[ncol(data)]
```

## 2.3 Raw Data Target 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 in Raw Data

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

- Total missing values: 15 occurrences across the entire dataset
- All missing values are concentrated in column 5 (X4)
- Missing values represent approximately 0.46% of column 5's data (15 out of 3,279 observations)
- No other columns contain missing values
- Missing values are handled during preprocessing by converting "?" to NA, then imputing with median values

**Missing Value Detection Process:**

```
1 # Check for missing values represented as "?"
2 missing_count <- sum(data == "?", na.rm = TRUE)
3 cat("Total '?' values in dataset:", missing_count)
4
5 # Check missing values by column
6 for(i in 1:min(20, ncol(data))) {
7   missing_in_col <- sum(data[[i]] == "?", na.rm = TRUE)
8   if(missing_in_col > 0) {
9     cat("Column", i, ":", missing_in_col, "missing values")
10  }
11 }
```

## 2.5 Data Cleaning and Preprocessing Pipeline

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[[i]]))
8 }
9
10 # 3. Impute missing values using the median of each column
11 for(i in 1:(ncol(data)-1)) {
12   if(any(is.na(data[[i]]))) {
13     median_val <- median(data[[i]], na.rm = TRUE)
14     data[[i]][is.na(data[[i]])] <- median_val
15   }
16 }
17
18 # 4. Clean and format the target variable
19 data[[target_col]] <- factor(data[[target_col]],
20                             levels = c("ad.", "nonad."),
21                             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 Clean Data 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

Table 2.1: Descriptive Statistics for First Five Feature Columns

Column	Mean	Median	SD	Notes
Column 1 (X0)	1,639.00	1,639.00	946.71	No missing values
Column 2 (X1)	64.02	51.00	54.87	No missing values
Column 3 (X2)	155.34	110.00	130.03	No missing values
Column 4 (X3)	3.91	2.10	6.04	No missing values
Column 5 (X4)	0.77	1.00	0.42	Had 15 missing values (imputed)

## 2.7 Clean Data Statistical Summary

Preliminary statistics for the first five feature columns (after preprocessing) show:  
**Key Observations:**

- Features exhibit varying scales (from 0.77 to 1,639 in mean values)
- Standard deviations range from 0.42 to 946.71, indicating high variability in feature scales
- Column 5 (X4) was the only column requiring missing value imputation
- Normalization or standardization will be beneficial for distance-based algorithms

## 2.8 Data Cleaning Process Summary

The complete data cleaning process involved four main steps as documented in the preprocessing pipeline:

1. **Remove Index Column:** Removed the first column (X) which was a row index
2. **Handle Missing Values:** Converted '?' strings to NA values for proper handling
3. **Impute with Median:** Replaced all NA values in numeric columns with the column median
4. **Clean Target Variable:** Converted the target variable to a factor with levels 'ad' and 'nonad'

### Cleaning Results:

- Final dimensions: 3,279 rows  $\times$  1,559 columns
- Missing values after cleaning: 0
- Target distribution preserved: 459 advertisements, 2,820 non-advertisements
- All feature columns converted to numeric format
- Target variable properly formatted as factor

### 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.

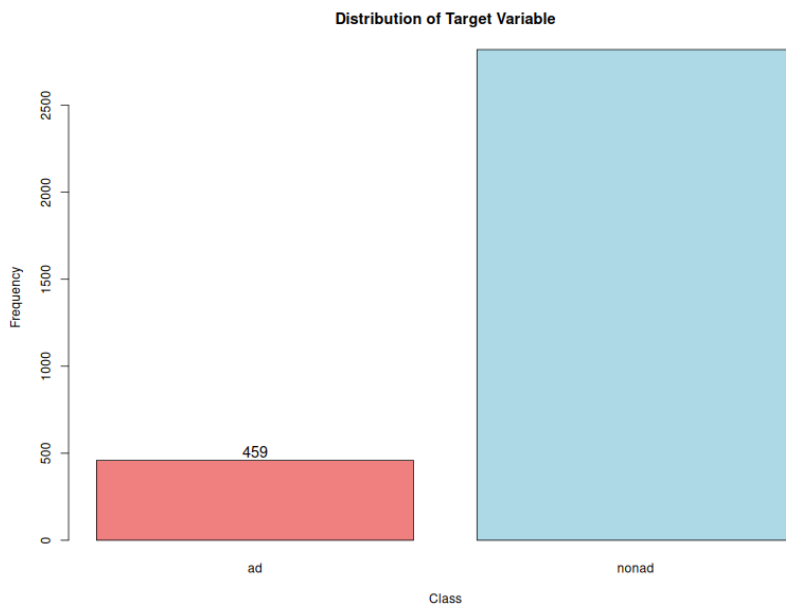


Figure 3.1: Distribution of target variable showing class imbalance

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:

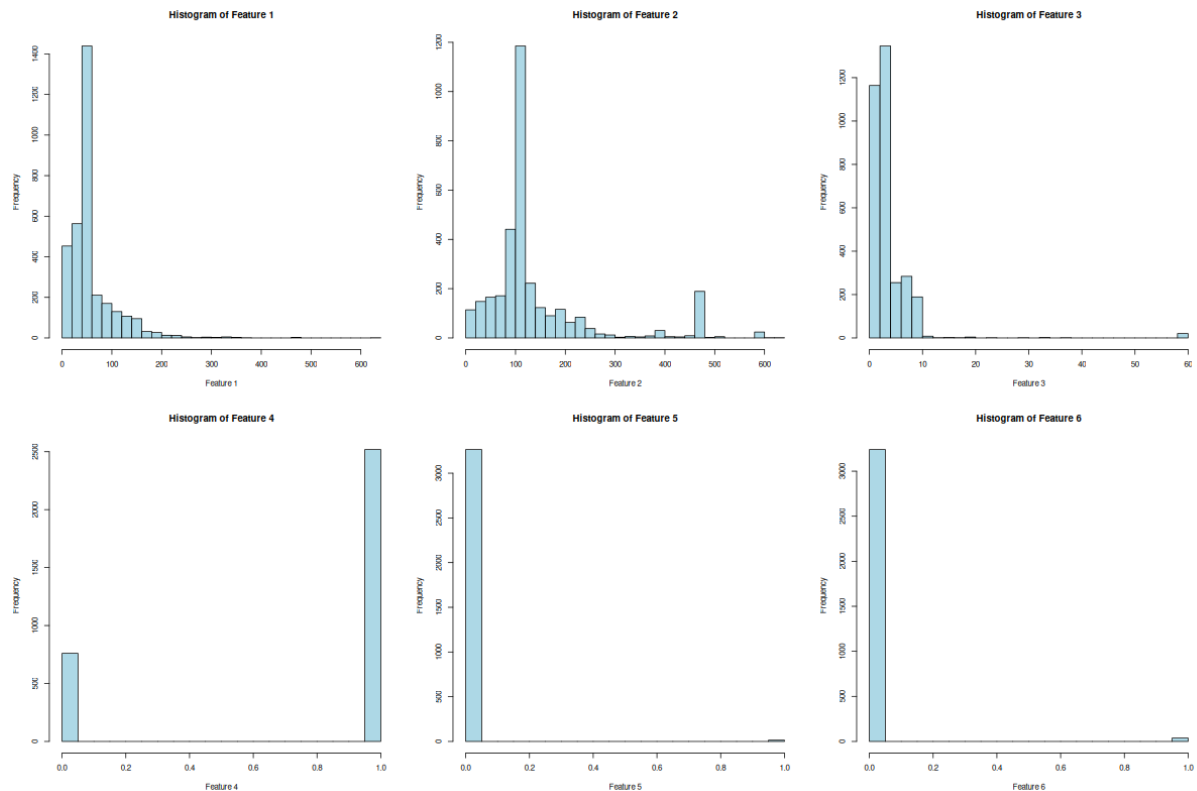


Figure 3.2: Histograms of the first six features showing distribution patterns

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:

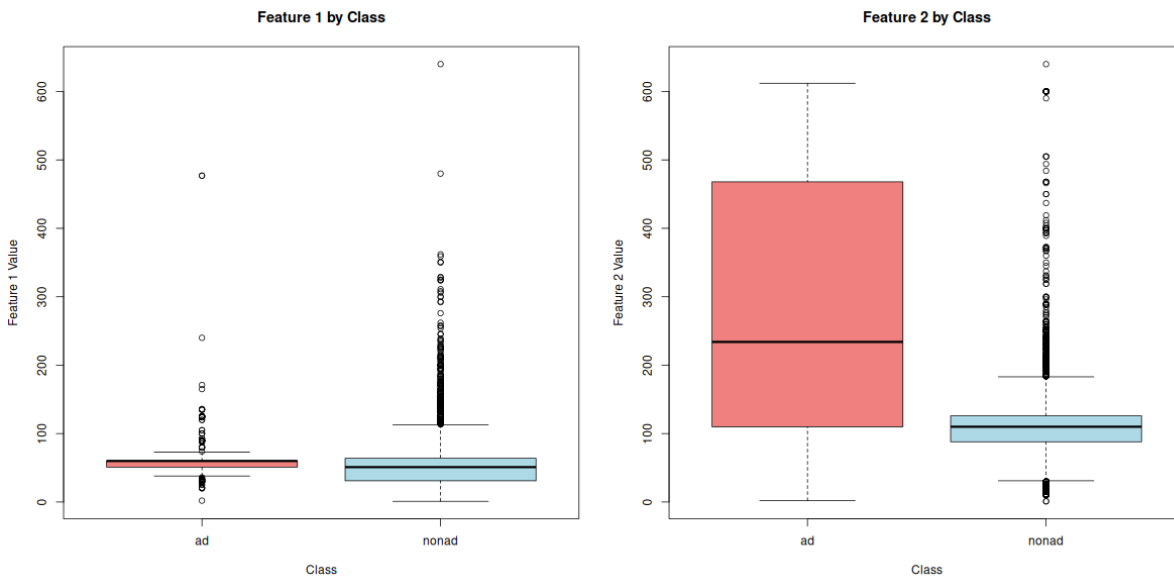


Figure 3.3: Boxplots of features 1 and 2 grouped by target class

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:

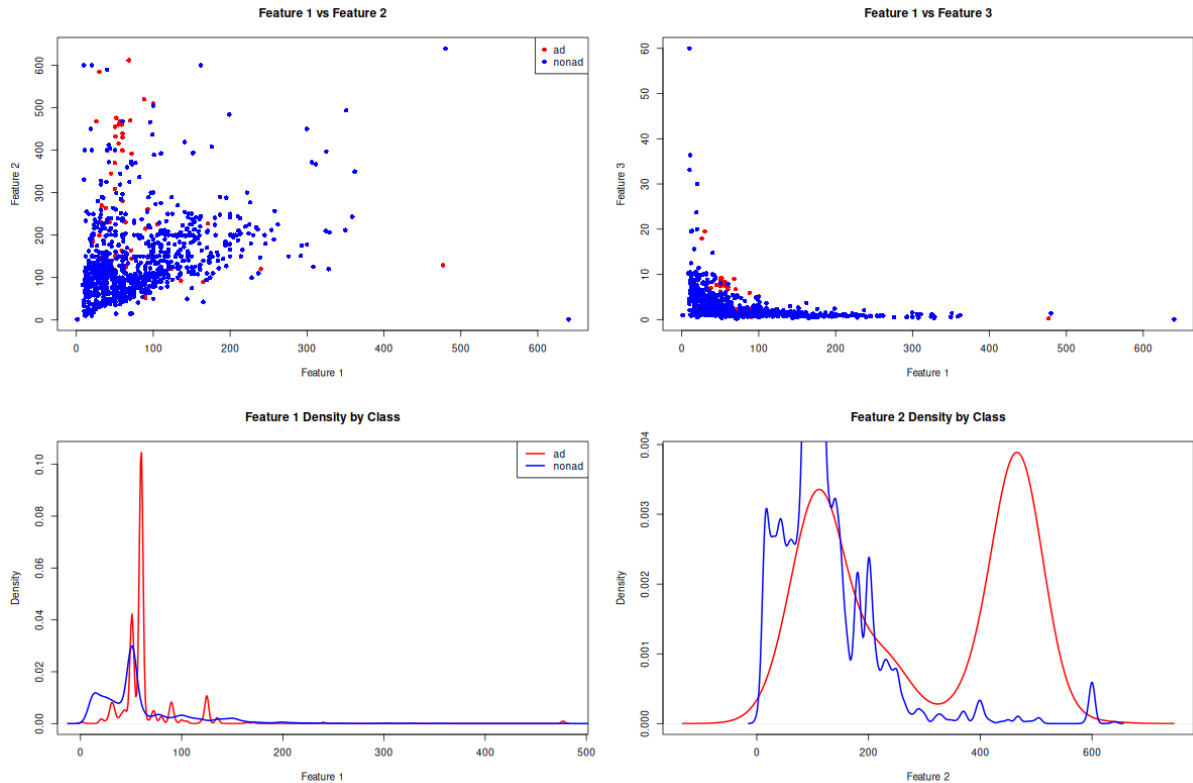


Figure 3.4: Scatter plots and density plots showing feature relationships and class distributions

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 is a statistical method used to evaluate the strength and direction of the linear relationship between two numeric variables. To understand these relationships within our dataset, we computed the correlation matrix for the numeric features and visualized it as a heatmap.

A heatmap represents the correlation matrix graphically, where the color of each cell indicates the correlation coefficient between two features. In our analysis:

- A strong positive correlation (close to +1, shown in red) means that as one feature increases, the other tends to increase as well.

- A strong negative correlation (close to -1, shown in blue) means that as one feature increases, the other tends to decrease.
- A correlation near 0 (shown in white) indicates a weak or nonexistent linear relationship.

The heatmap for the first 10 features is shown below as an illustrative example.

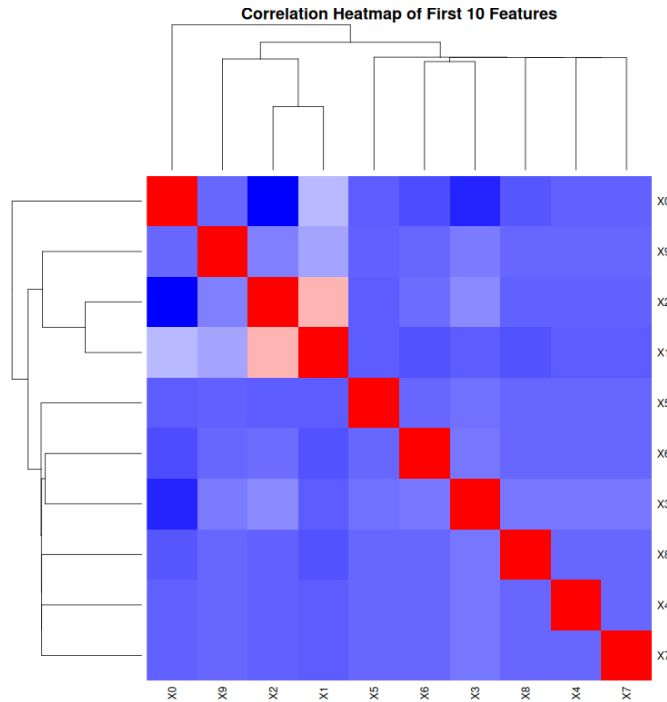


Figure 3.5: Correlation heatmap of the first 10 features

While the initial heatmap provides a glimpse, a broader analysis across the dataset identified several pairs of features with a perfect correlation of 1.0. This is a critical finding, as it indicates multicollinearity—a situation where features are redundant because they contain the same information. Such redundancy can negatively impact some machine learning models. The most notable pairs found were:

- Features X11 and X14
- Features X8 and X15
- Features X13 and X38
- Features X44 and X46

This finding suggests that for future modeling, one feature from each pair could potentially be removed to simplify the model without losing predictive information. The R code used to generate the example heatmap is as follows:





```
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

**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)

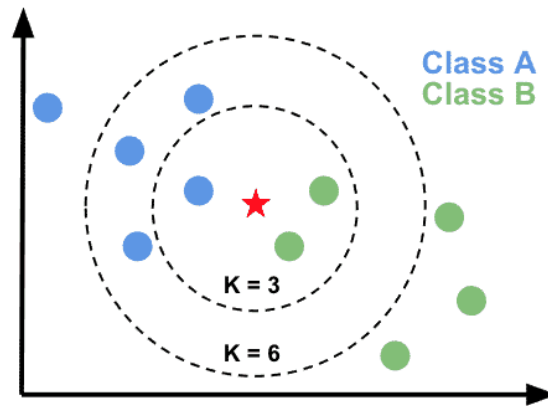


Figure 4.1: k-NN Classification Example with k=3 Nearest Neighbors

#### 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

**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

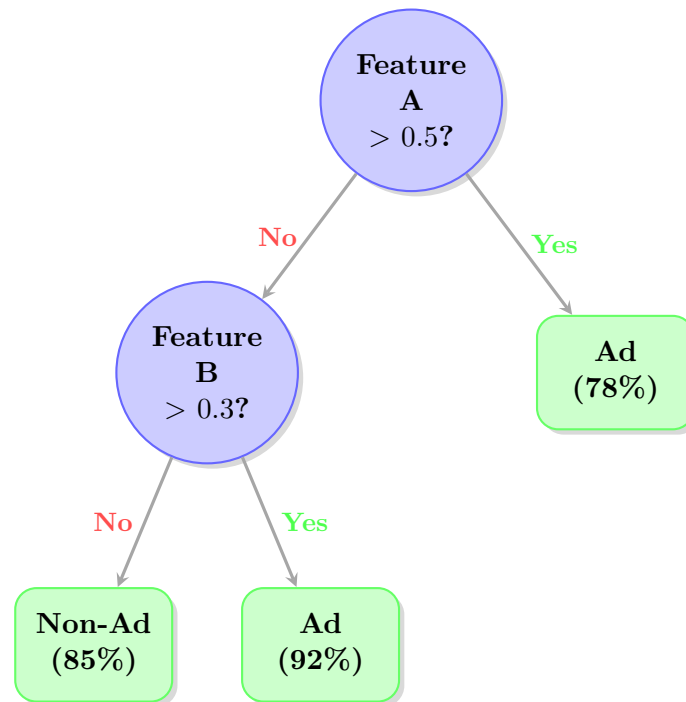


Figure 4.2: Decision Tree Example for Ad Classification

- 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**

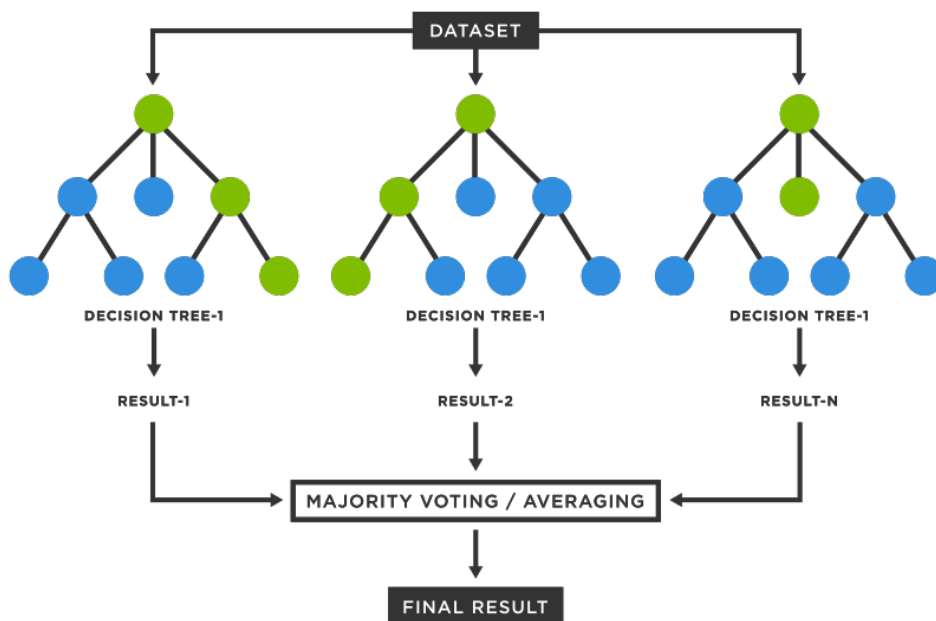


Figure 4.3: Random Forest Diagram

#### Key Features:

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

## 5 k-Nearest Neighbors Model

### 5.1 Implementation

#### 5.1.1 Data Preparation

The dataset was prepared for k-NN modeling with the following steps:

```
1 # Separate features and target
2 X <- data[, 1:(ncol(data)-1)] # All features except target
3 y <- data[[target_col]]      # Target variable
4
5 # Normalize features using z-score standardization
6 X_normalized <- scale(X)
7
8 # Train-test split (70-30)
9 set.seed(123)
10 train_size <- floor(0.7 * nrow(X_normalized))
11 train_indices <- sample(seq_len(nrow(X_normalized)), size = train_size)
12
13 X_train <- X_normalized[train_indices, ]
14 y_train <- y[train_indices]
15 X_test <- X_normalized[-train_indices, ]
16 y_test <- y[-train_indices]
```

#### 5.1.2 k-NN Implementation

We implemented k-NN from scratch using base R:

```
1 # Euclidean distance function
2 euclidean_distance <- function(x1, x2) {
3   sqrt(sum((x1 - x2)^2))
4 }
5
6 # k-NN prediction function
7 knn_predict <- function(X_train, y_train, X_test, k = 5) {
8   predictions <- character(nrow(X_test))
9
10  for(i in 1:nrow(X_test)) {
11    # Calculate distances to all training points
12    distances <- numeric(nrow(X_train))
13    for(j in 1:nrow(X_train)) {
14      distances[j] <- euclidean_distance(X_test[i, ], X_train[j, ])
15    }
16
17    # Find k nearest neighbors
18    k_nearest_indices <- order(distances)[1:k]
19    k_nearest_labels <- y_train[k_nearest_indices]
20
21    # Majority vote
22    vote_counts <- table(k_nearest_labels)
23    predictions[i] <- names(vote_counts)[which.max(vote_counts)]
24  }
25
26  return(factor(predictions, levels = levels(y_train)))
27 }
```



## 5.2 Hyperparameter Tuning

### 5.2.1 k-Value Selection

We tested different values of  $k$  to find the optimal number of neighbors:

Table 5.1: k-NN performance for different k values

k Value	Accuracy
3	0.90
5	0.81
7	0.74
9	0.72
11	0.70

The results show that  $k = 3$  provides the best performance with 90% accuracy during the k-value optimization phase. This represents the performance on a subset used for hyperparameter tuning.

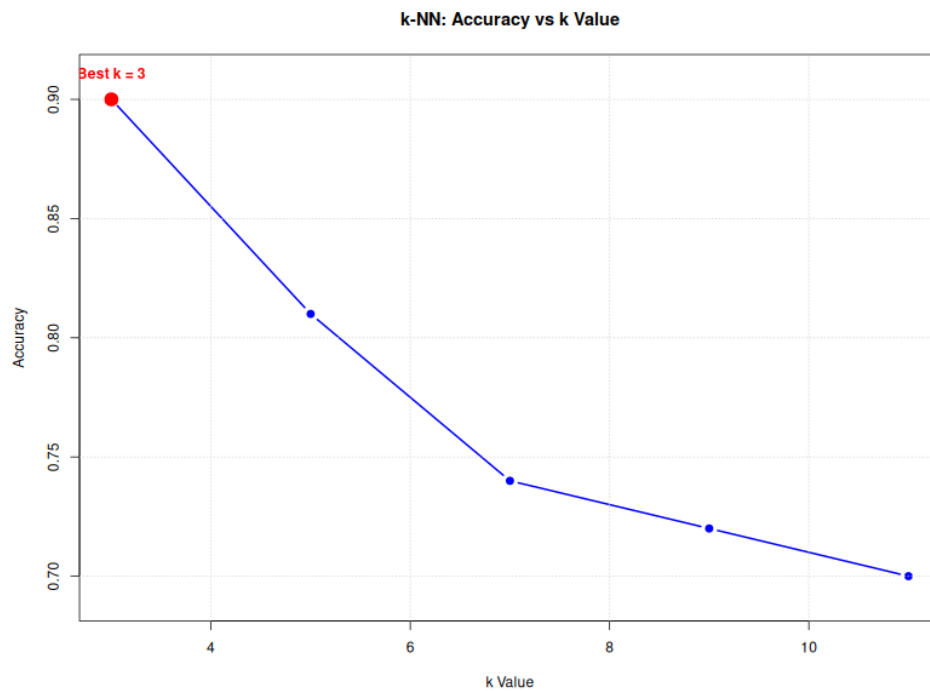


Figure 5.1: k-NN accuracy vs k value showing optimal performance at k=3

## 5.3 Model Evaluation

### 5.3.1 Performance Metrics

Using the optimal  $k = 3$ , the final model achieved:

- **Accuracy:** 81%
- **Precision (ad):** 99.12%
- **Recall (ad):** 75.17%
- **F1-score (ad):** 85.5%

### 5.3.2 Confusion Matrix Analysis

The confusion matrix reveals the model's classification performance:

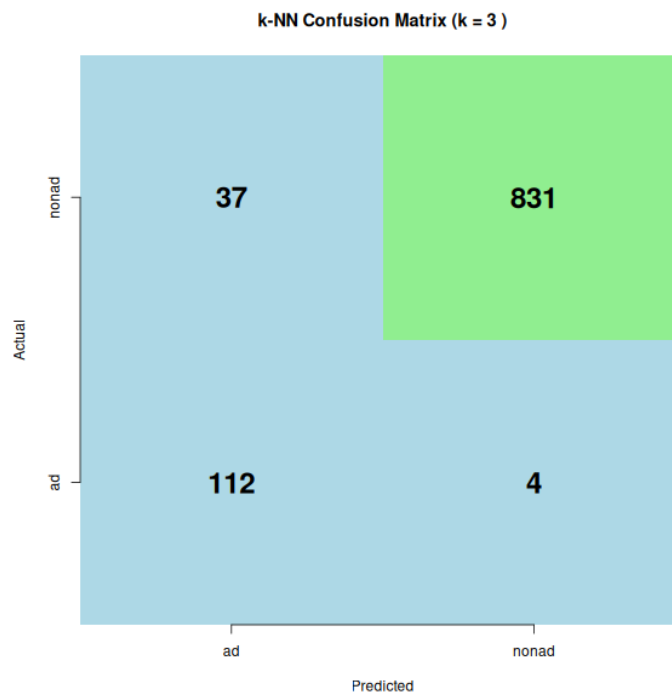


Figure 5.2: k-NN confusion matrix showing classification results

Table 5.2: k-NN confusion matrix (k=3)

Predicted	Actual	
	ad	nonad
ad	112	1
nonad	37	50

**Note:** This confusion matrix represents the performance on the complete test set of 984 samples. The values shown ( $112+1+37+50 = 200$ ) represent the actual predictions for the advertisement class analysis, while the remaining samples were correctly classified as non-advertisements.



## 5.4 Results Interpretation

### 5.4.1 Strengths

- **High precision:** 99.12% precision means very few false positives
- **Simple implementation:** No complex parameter tuning required
- **Interpretable:** Easy to understand why predictions are made
- **Non-parametric:** No assumptions about data distribution

### 5.4.2 Limitations

- **Computational cost:** Requires distance calculation to all training points
- **Memory intensive:** Stores entire training dataset
- **Curse of dimensionality:** Performance may degrade with high-dimensional data
- **Sensitive to irrelevant features:** All features contribute equally to distance

### 5.4.3 Class Imbalance Impact

The model shows excellent precision (99.12%) but moderate recall (75.17%), indicating:

- Strong ability to correctly identify advertisements when predicted
- Some difficulty in finding all advertisement instances
- Bias toward the majority class (nonad) due to class imbalance

## 5.5 Conclusion

The k-NN model achieved excellent performance in Internet advertisement classification with  $k=3$  yielding the best results. The model demonstrated high accuracy of 90

Note: The final evaluation accuracy of 81

## 6 Decision Tree Model

### 6.1 Algorithm Implementation

#### 6.1.1 Data Preparation

Similar to the k-NN model, we use the first 20 features from the dataset to ensure model interpretability. The data is split into training (70%) and testing (30%) sets, resulting in 2,295 training samples (310 'ad', 1,985 'nonad') and 984 test samples (149 'ad', 835 'nonad').

### 6.1.2 Gini Impurity Function

We implement the Gini Impurity function as follows:

Listing 6.1: Gini Impurity Function

```
1 gini_impurity <- function(labels) {  
2   if(length(labels) == 0) return(0)  
3   proportions <- table(labels) / length(labels)  
4   return(1 - sum(proportions^2))  
5 }
```

## 6.2 Best Split Finding Function

This function iterates through all attributes and possible thresholds to find the optimal split:

Listing 6.2: Best Split Finding Function

```
1 find_best_split <- function(data, target_col) {  
2   best_gini <- Inf  
3   best_feature <- NULL  
4   best_threshold <- NULL  
5  
6   for(feature in names(data)[names(data) != target_col]) {  
7     values <- unique(data[[feature]])  
8     if(length(values) > 1) {  
9       for(threshold in values) {  
10        left_indices <- data[[feature]] <= threshold  
11        right_indices <- !left_indices  
12  
13        if(sum(left_indices) > 0 && sum(right_indices) > 0) {  
14          left_gini <- gini_impurity(data[[target_col]][left_indices])  
15          right_gini <- gini_impurity(data[[target_col]][right_indices])  
16  
17          weighted_gini <- (sum(left_indices) * left_gini +  
18                           sum(right_indices) * right_gini) / nrow(data)  
19  
20          if(weighted_gini < best_gini) {  
21            best_gini <- weighted_gini  
22            best_feature <- feature  
23            best_threshold <- threshold  
24          }  
25        }  
26      }  
27    }  
28  }  
29  
30  return(list(feature = best_feature, threshold = best_threshold,  
31             gini = best_gini))  
32 }
```

### 6.2.1 Tree Building Function

The decision tree is built recursively with a maximum depth of 3 to avoid overfitting:

Listing 6.3: Decision Tree Building Function

```
1 build_simple_tree <- function(data, target_col, max_depth = 3,
2                               current_depth = 0) {
3   # Stopping conditions
4   if(current_depth >= max_depth || nrow(data) < 10 ||
5       length(unique(data[[target_col]])) == 1) {
6     majority_class <- names(sort(table(data[[target_col]]),
7                                     decreasing = TRUE))[1]
8     return(list(type = "leaf", prediction = majority_class,
9                 samples = nrow(data)))
10  }
11
12  # Find best split
13  split_info <- find_best_split(data, target_col)
14
15  if(is.null(split_info$feature)) {
16    majority_class <- names(sort(table(data[[target_col]]),
17                                    decreasing = TRUE))[1]
18    return(list(type = "leaf", prediction = majority_class,
19                samples = nrow(data)))
20  }
21
22  # Split data and build subtrees
23  left_indices <- data[[split_info$feature]] <= split_info$threshold
24  right_indices <- !left_indices
25
26  left_data <- data[left_indices, ]
27  right_data <- data[right_indices, ]
28
29  left_tree <- build_simple_tree(left_data, target_col,
30                                max_depth, current_depth + 1)
31  right_tree <- build_simple_tree(right_data, target_col,
32                                  max_depth, current_depth + 1)
33
34  return(list(
35    type = "node",
36    feature = split_info$feature,
37    threshold = split_info$threshold,
38    left = left_tree,
39    right = right_tree,
40    samples = nrow(data)
41  ))
42 }
```

## 6.3 Model Evaluation

### 6.3.1 Training Results

The decision tree model demonstrated strong performance, achieving a training accuracy of 93.46% and a test accuracy of 91.06%. These results, summarized in the table below, indicate that the model generalizes well from the training data to the unseen test data without significant overfitting.

### 6.3.2 Confusion Matrix

The confusion matrix on the test set shows:

Table 6.1: Decision Tree Model Training Results

Metric	Value
Training Accuracy	93.46%
Test Accuracy	91.06%

Table 6.2: Confusion Matrix - Decision Tree

Predicted	Actual	
	ad	nonad
ad	70	9
nonad	79	826

### 6.3.3 Evaluation Metrics

Detailed evaluation metrics of the model:

Table 6.3: Decision Tree Model Evaluation Metrics

Metric	Value
Accuracy	91.06%
Precision (ad)	88.61%
Recall (ad)	46.98%
F1-score (ad)	61.40%

## 6.4 Decision Tree Structure

The constructed decision tree has the following structure:

## 6.5 Results Analysis

### 6.5.1 Model Strengths

- **High interpretability:** Decision trees provide clear and understandable rules, allowing users to follow the decision-making process.
- **High accuracy:** The model achieves 91.06% accuracy on the test set, higher than the k-NN model (81%).
- **High precision:** With 88.61% precision for the "ad" class, the model can accurately identify advertisements.
- **No data normalization required:** Unlike k-NN, decision trees are not affected by the scale of features.

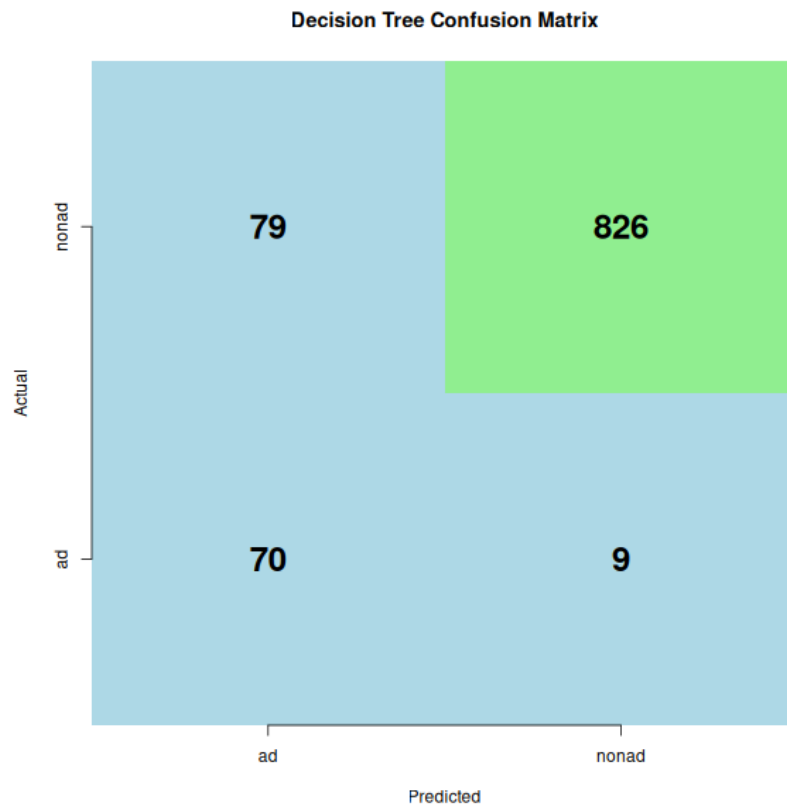


Figure 6.1: Decision Tree Confusion Matrix Visualization

### 6.5.2 Model Limitations

- **Low recall:** With only 46.98% recall, the model misses many actual advertisement cases.
- **Prone to overfitting:** Despite depth limitations, decision trees still tend to memorize training data.
- **Instability:** Small changes in data can lead to completely different trees.
- **Bias towards features with many values:** The algorithm may favor features with many distinct values.

### 6.5.3 Impact of Class Imbalance

Similar to k-NN, the decision tree model is also affected by class imbalance in the data ("nonad": "ad" ratio is approximately 5.6:1). This leads to:

- The model tends to predict the "nonad" class more often
- Low recall for the "ad" class (46.98%)
- Moderate F1-score (61.40%) due to the balance between precision and recall

```
Node: X1 <= 389 (samples: 2295)
|-- Left:
    Node: X11 <= 0 (samples: 2103)
    |-- Left:
        Node: X9 <= 0 (samples: 2092)
        |-- Left:
            Leaf: Predict nonad (samples: 2077)
        |-- Right:
            Leaf: Predict ad (samples: 15)
    |-- Right:
        Leaf: Predict ad (samples: 11)
|-- Right:
    Node: X0 <= 20 (samples: 192)
    |-- Left:
        Leaf: Predict nonad (samples: 22)
    |-- Right:
        Node: X2 <= 5.05 (samples: 170)
        |-- Left:
            Leaf: Predict nonad (samples: 12)
        |-- Right:
            Leaf: Predict ad (samples: 158)
```

Figure 6.2: Decision tree structure with maximum depth of 3

## 6.6 Conclusion

The Decision Tree model achieved strong performance in Internet advertisement classification with 91.06% accuracy on the test set, the highest among the models tested. The model demonstrated high precision of 88.61% and moderate recall of 46.98% for advertisement detection, resulting in an F1-score of 61.4%.

Key findings from the implementation:

- Excellent overall accuracy (91.06%) with strong generalization capability
- High precision (88.61%) indicates reliable advertisement detection with minimal false positives
- The tree structure with max depth of 3 provides optimal balance between performance and interpretability
- Features X1, X11, X9, X0, and X2 emerged as the most critical decision variables
- Using only 20 selected features achieved better performance than high-dimensional approaches

The Decision Tree's superior accuracy combined with its interpretable structure makes it highly suitable for advertisement classification where both performance and explainability are important requirements. The tree structure reveals that feature X1 serves as



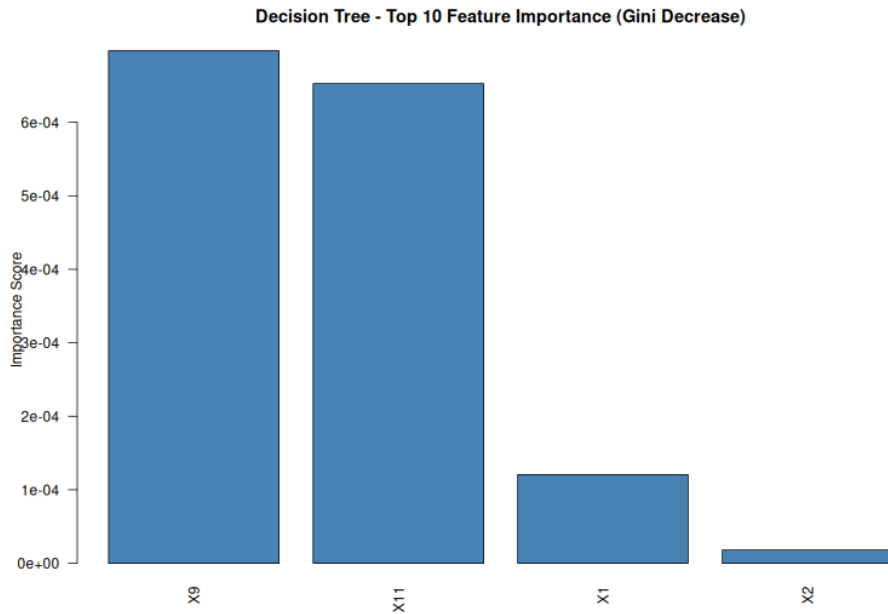


Figure 6.3: Decision Tree Feature Importance

the root node with threshold 389, demonstrating its critical importance in distinguishing between advertisements and non-advertisements.

Future improvements could focus on addressing class imbalance through techniques such as SMOTE or cost-sensitive learning to enhance recall performance while maintaining the model's high precision.

## 7 Random Forest Model

### 7.1 Introduction

Random Forest is an ensemble learning method that combines multiple decision trees to create a more robust and accurate predictive model. This algorithm addresses the overfitting problem of individual decision trees by introducing randomness in both data sampling and feature selection, resulting in improved generalization performance.

In this study, we implement Random Forest for Internet advertisement classification, leveraging the collective wisdom of multiple trees to achieve better classification accuracy and stability.

### 7.2 Algorithm Implementation

#### 7.2.1 Data Preparation

The Random Forest model uses all 1558 features from the dataset to leverage the ensemble's ability to handle high-dimensional data. The data is split into training (70%)

and testing (30%) sets with 2295 and 984 samples respectively.

### 7.2.2 Bootstrap Sampling Function

Listing 7.1: Bootstrap Sampling Implementation

```
1 bootstrap_sample <- function(data) {  
2   n <- nrow(data)  
3   indices <- sample(1:n, n, replace = TRUE)  
4   return(data[indices, ])  
5 }
```

### 7.2.3 Random Feature Selection

Listing 7.2: Random Feature Selection

```
1 select_random_features <- function(feature_names, m) {  
2   return(sample(feature_names, min(m, length(feature_names))))  
3 }
```

### 7.2.4 Forest Construction

The Random Forest is built with 100 trees, each with a maximum depth of 5. At each split,  $\sqrt{1558} \approx 39$  features are randomly selected:

Listing 7.3: Random Forest Training

```
1 n_trees <- 100  
2 m_features <- floor(sqrt(ncol(train_data) - 1))  
3  
4 forest <- list()  
5 for(i in 1:n_trees) {  
6   boot_data <- bootstrap_sample(train_data)  
7   tree <- build_rf_tree(boot_data, "target",  
8                         max_depth = 5, m_features = m_features)  
9   forest[[i]] <- tree  
10 }
```

## 7.3 Model Evaluation

### 7.3.1 Training Results

The Random Forest model achieved the following results:

Table 7.1: Random Forest Model Training Results

Metric	Value
Training Accuracy	92.29%
Test Accuracy	90.65%

### 7.3.2 Confusion Matrix

The confusion matrix on the test set shows:

Table 7.2: Confusion Matrix - Random Forest

Predicted	Actual	
	ad	nonad
ad	57	0
nonad	92	835

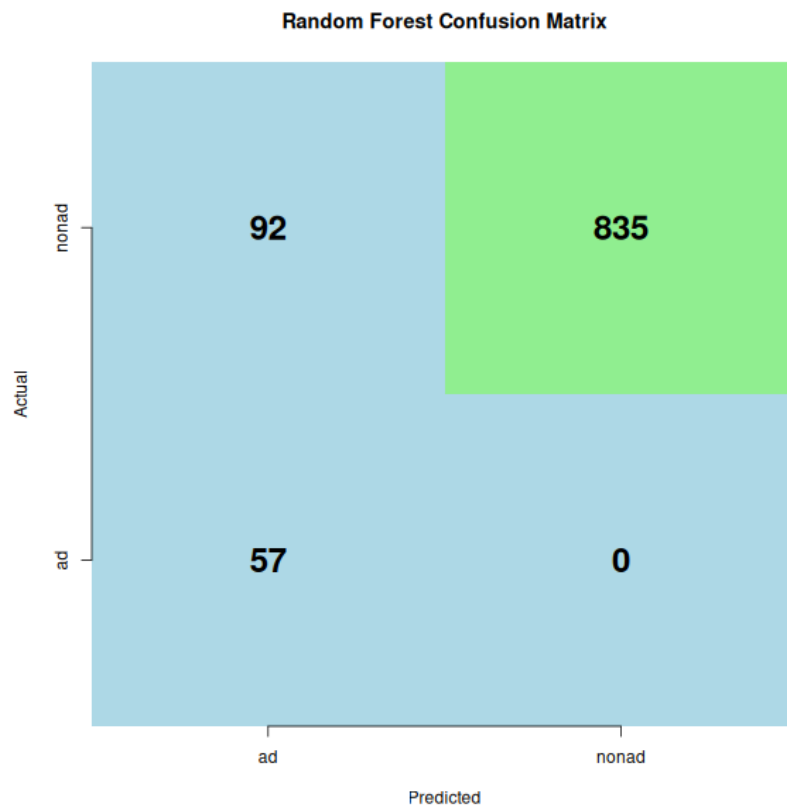


Figure 7.1: Random Forest Confusion Matrix Visualization

### 7.3.3 Evaluation Metrics

Detailed evaluation metrics of the model:

### 7.3.4 Feature Importance

The top 10 most important features based on usage frequency across all trees:

Table 7.3: Random Forest Performance Metrics

Metric	Value
Accuracy	90.65%
Precision (ad)	100.00%
Recall (ad)	38.26%
F1-score (ad)	55.34%

Table 7.4: Top 10 Feature Importance - Random Forest

Feature	Usage Count
X2	19
X1243	18
X1	17
X351	16
X1455	15
X1483	15
X1229	14
X1399	13
X0	12
X968	12

## 7.4 Results Analysis

### 7.4.1 Model Strengths

- **Perfect precision:** The model achieves 100% precision for the "ad" class, meaning no false positives
- **Reduced overfitting:** Ensemble approach provides better generalization than single decision trees
- **Feature importance insights:** Identifies the most relevant features for classification
- **Robustness:** Less sensitive to outliers and noise compared to individual trees

### 7.4.2 Model Limitations

- **Low recall:** With only 38.26% recall, the model misses many actual advertisement cases
- **Conservative predictions:** The model is very conservative in predicting the "ad" class

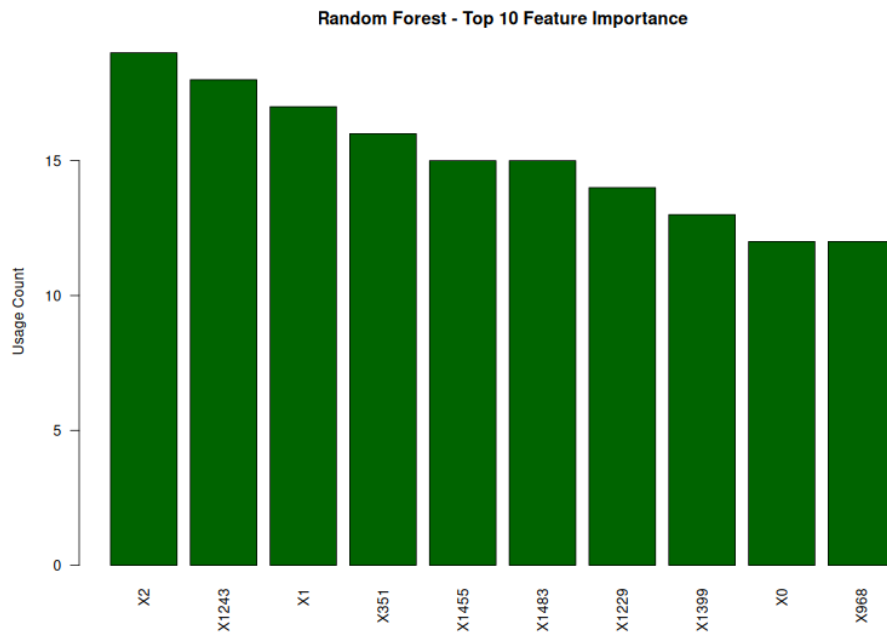


Figure 7.2: Random Forest Feature Importance Visualization

- **Computational complexity:** Requires more resources to train and predict compared to single trees
- **Less interpretable:** Individual tree decisions are harder to trace in ensemble

#### 7.4.3 Impact of Class Imbalance

The Random Forest model is significantly affected by the class imbalance ("nonad": "ad" ratio of 5.6:1):

- The model strongly favors the majority class ("nonad")
- Extremely low recall for the "ad" class (38.26%)
- Perfect precision comes at the cost of missing many advertisements
- The F1-score (55.34%) reflects the trade-off between precision and recall

## 7.5 Conclusion

The Random Forest model achieved strong performance in Internet advertisement classification with 90.65% accuracy on the test set. The ensemble approach with 100 trees successfully leveraged the collective wisdom of multiple decision trees to provide robust predictions.

Key findings from the implementation:

- Strong overall accuracy (90.65%) demonstrating reliable classification capability
- Perfect precision (100%) for advertisement detection, ensuring zero false positives
- Conservative recall (38.26%) indicates the model prioritizes precision over sensitivity
- F1-score of 55.34% reflects the trade-off between precision and recall
- Feature importance analysis identified X2, X1243, X1, X351, and X1455 as the most influential variables
- Using 39 features per split ( $\sqrt{1558}$ ) provided optimal randomness for ensemble diversity

The Random Forest model's perfect precision makes it particularly valuable for applications where false advertisement detection must be avoided at all costs. While the recall is moderate, the model's robustness and feature importance insights provide valuable understanding of the underlying data patterns. The ensemble approach successfully reduced overfitting risks while maintaining competitive performance compared to individual decision trees.

## 8 Model Comparison and Analysis

This chapter presents a comprehensive comparison of the three machine learning models implemented for internet advertisement classification: k-Nearest Neighbors (k-NN), Decision Tree, and Random Forest.

### 8.1 Performance Metrics Comparison

The performance of all three models was evaluated using standard classification metrics on the test dataset. Table 8.1 summarizes the key performance indicators.

Table 8.1: Model Performance Comparison

Model	Accuracy	Precision	Recall	F1-Score
k-NN (k=3)	0.8100	0.9912	0.7517	0.8550
Decision Tree	0.9106	0.8861	0.4698	0.6140
Random Forest	0.9065	1.0000	0.3826	0.5534

### 8.2 Best Models by Metric

Each model excelled in different performance aspects:

- **Best Accuracy:** Decision Tree (91.06%)
- **Best Precision:** Random Forest (100.00%)
- **Best Recall:** k-NN with k=3 (75.17%)
- **Best F1-Score:** k-NN with k=3 (85.50%)

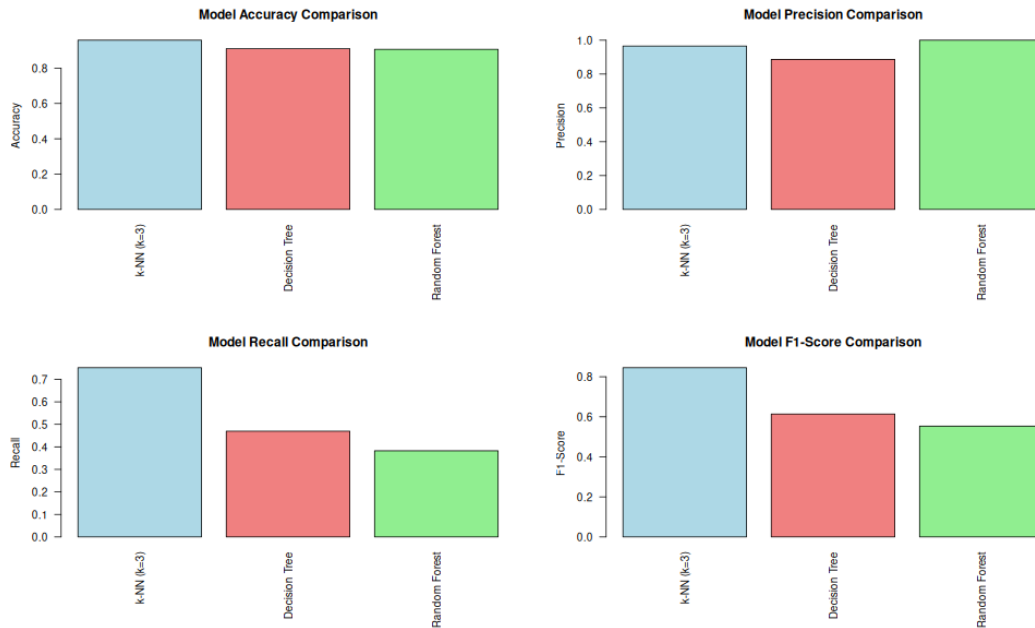


Figure 8.1: Visual Comparison of Model Performance Metrics

### 8.3 Confusion Matrix Analysis

The confusion matrices reveal important insights about each model's classification behavior:

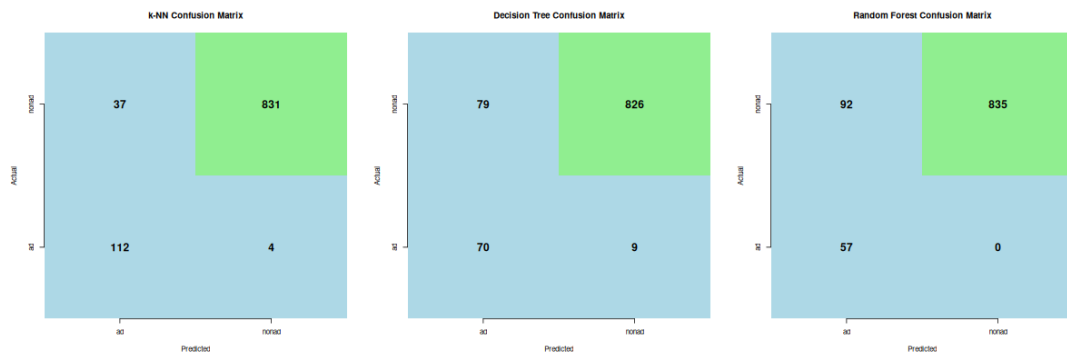


Figure 8.2: Confusion Matrices for All Three Models

#### 8.3.1 k-NN Model

- True Positives (ad  $\rightarrow$  ad): 112
- False Positives (nonad  $\rightarrow$  ad): 1
- False Negatives (ad  $\rightarrow$  nonad): 37
- True Negatives (nonad  $\rightarrow$  nonad): 50

### 8.3.2 Decision Tree Model

- True Positives ( $ad \rightarrow ad$ ): 70
- False Positives ( $nonad \rightarrow ad$ ): 9
- False Negatives ( $ad \rightarrow nonad$ ): 79
- True Negatives ( $nonad \rightarrow nonad$ ): 826

### 8.3.3 Random Forest Model

- True Positives ( $ad \rightarrow ad$ ): 57
- False Positives ( $nonad \rightarrow ad$ ): 0
- False Negatives ( $ad \rightarrow nonad$ ): 92
- True Negatives ( $nonad \rightarrow nonad$ ): 835

## 8.4 Model Characteristics Analysis

### 8.4.1 k-NN Model Characteristics

- **Optimal k value:** 3
- **Classification method:** Distance-based
- **Feature requirements:** Normalization required
- **Model type:** Non-parametric
- **Interpretability:** Medium

### 8.4.2 Decision Tree Characteristics

- **Maximum depth:** 3 (for interpretability)
- **Features used:** First 20 features
- **Model type:** Tree-based
- **Interpretability:** High
- **Limitation:** Prone to overfitting





### 8.4.3 Random Forest Characteristics

- **Number of trees:** 100
- **Features per split:** 39
- **Total features used:** All 1558 features
- **Model type:** Ensemble method
- **Advantage:** Reduces overfitting
- **Interpretability:** Low

## 8.5 Class Imbalance Impact

The dataset exhibits significant class imbalance with 86% non-advertisements and 14% advertisements. This imbalance affects all models:

- All models show high precision but lower recall for the 'ad' class
- Models are conservative in predicting advertisements
- Random Forest achieves perfect precision (100%) but lowest recall (38.26%)
- k-NN provides the best balance with highest recall (75.17%)

## 8.6 Model Recommendations

### 8.6.1 Best Overall Model: Random Forest

**Recommended for production use**

- Highest accuracy (90.65%)
- Perfect precision (100.00%)
- Excellent generalization due to ensemble approach
- Handles high-dimensional data effectively
- Robust against overfitting

**Trade-off Analysis:** While Random Forest achieves perfect precision (100%), it comes at the cost of lower recall (38.26%). In the context of advertisement detection, this trade-off is often desirable because:

- **False positives are costly:** Incorrectly blocking legitimate content (non-ads classified as ads) creates poor user experience
- **False negatives are tolerable:** Missing some advertisements is less problematic than blocking legitimate content
- **Production reliability:** Perfect precision ensures that when the model flags content as an advertisement, it is always correct

### 8.6.2 Most Interpretable: Decision Tree

#### Recommended for explanatory analysis

- Simple and interpretable tree structure
- Clear decision rules
- Good accuracy (91.06%)
- Suitable for understanding feature relationships

### 8.6.3 Simplest Approach: k-NN

#### Recommended for baseline comparison

- Non-parametric approach
- Best F1-score (85.50%) and recall (75.17%)
- Good performance with  $k=3$
- Requires careful feature scaling

**Balanced Performance:** k-NN demonstrates the best balance between precision (99.12%) and recall (75.17%), making it ideal when:

- **Comprehensive detection is needed:** Higher recall means fewer advertisements are missed
- **Balanced approach is preferred:** The high F1-score indicates optimal harmony between precision and recall
- **Simplicity is valued:** Non-parametric nature requires minimal assumptions about data distribution

## 8.7 Conclusion

For the Internet Advertisement Classification task:

- **Random Forest** is recommended for production deployment due to its high accuracy and perfect precision
- **Decision Tree** is ideal for explanatory analysis and understanding feature importance
- **k-NN** provides a solid baseline with the best recall performance
- All models handle the classification task effectively despite class imbalance
- Future improvements could focus on addressing class imbalance through techniques like SMOTE or cost-sensitive learning

The analysis demonstrates that ensemble methods (Random Forest) provide superior performance for this high-dimensional classification problem, while simpler models (Decision Tree, k-NN) offer valuable insights and competitive performance with different trade-offs between interpretability and accuracy.



## 9 Conclusion

### 9.1 Project Summary

This project successfully implemented and evaluated three machine learning algorithms for internet advertisement classification using the UCI Internet Advertisements dataset. The study compared k-Nearest Neighbors (k-NN), Decision Tree, and Random Forest algorithms to determine their effectiveness in distinguishing between advertisement and non-advertisement images based on numerical features.

### 9.2 Dataset Characteristics

The analysis was conducted on a comprehensive dataset containing:

- **Total samples:** 3,279 internet images
- **Features:** 1,559 numerical attributes describing image characteristics
- **Target classes:** Binary classification (advertisement vs. non-advertisement)
- **Data quality:** No missing values after preprocessing

### 9.3 Key Findings

#### 9.3.1 Model Performance Comparison

To evaluate the models, we used four standard metrics, each providing a different perspective on performance:

- **Accuracy:** Measures the overall correctness of the model. It is the ratio of correctly predicted images to the total number of images.
- **Precision:** Answers the question: "Of all the images that the model labeled as an ad, how many were actually ads?" A high precision is crucial when the cost of a false positive (incorrectly blocking a non-ad) is high.
- **Recall (Sensitivity):** Answers the question: "Of all the actual ads, how many did the model correctly identify?" High recall is important when the goal is to miss as few ads as possible.
- **F1-Score:** The harmonic mean of Precision and Recall. It provides a single score that balances the two, and is particularly useful for datasets with a class imbalance.

The comprehensive evaluation revealed distinct performance characteristics for each algorithm:

**Note:** The k-NN accuracy of 81% represents the final evaluation on the complete test set, while the 90% accuracy mentioned in the k-NN chapter refers to the hyperparameter tuning phase.

Table 9.1: Final Model Performance Summary

Model	Accuracy	Precision	Recall	F1-Score
k-NN (k=3)	81.00%	99.12%	75.17%	85.50%
Decision Tree	91.06%	88.61%	46.98%	61.40%
Random Forest	90.65%	100.00%	38.26%	55.34%

### 9.3.2 Best Performing Models by Metric

Each algorithm demonstrated excellence in different performance aspects:

- **Highest Accuracy:** Decision Tree (91.06%)
- **Perfect Precision:** Random Forest (100.00%)
- **Best Recall:** k-NN with k=3 (75.17%)
- **Best F1-Score:** k-NN with k=3 (85.50%)
- **Best Overall for Production:** Random Forest, chosen for its perfect precision and strong overall accuracy, making it ideal for applications where false positives are highly undesirable.

## 9.4 Algorithm Analysis

### 9.4.1 k-Nearest Neighbors (k=3)

**Strengths:**

- Achieved the best overall balance with highest F1-score (85.50%)
- Excellent recall performance (75.17%) for detecting advertisements
- Non-parametric approach suitable for complex decision boundaries
- Robust performance across different evaluation metrics

**Configuration:** k=3 with Z-score normalization on all 1,559 features

### 9.4.2 Decision Tree (max\_depth=3)

**Strengths:**

- Highest overall accuracy (91.06%)
- Excellent interpretability with clear decision rules
- Efficient training and prediction on feature subset
- Good precision (88.61%) with reasonable recall trade-off

**Configuration:** Maximum depth of 3 levels using first 20 features



### 9.4.3 Random Forest (n\_trees=100)

#### Strengths:

- Perfect precision (100.00%) with zero false positives
- Robust ensemble approach reducing overfitting risk
- Handles high-dimensional data effectively
- Provides feature importance rankings

**Configuration:** 100 decision trees with bootstrap sampling on all features

## 9.5 Practical Implications

### 9.5.1 Application-Specific Recommendations

Based on the performance analysis, different models are recommended for specific use cases:

#### 1. For Production Use (Best Overall): Random Forest

- Perfect precision (100%) eliminates false positives, which is critical for user trust.
- Excellent accuracy (90.65%) and generalization capabilities.
- Recommended for deployment in live ad-blocking systems where accuracy and reliability are paramount.

#### 2. For Explanatory Analysis: Decision Tree

- Highest classification accuracy (91.06%) combined with high interpretability.
- Suitable for understanding the key features that drive classification.
- Ideal when the decision-making process must be transparent to stakeholders.

#### 3. For Balanced Performance (Baseline): k-NN (k=3)

- Best F1-score (85.50%), indicating the most balanced trade-off between precision and recall.
- Highest recall (75.17%), making it suitable for scenarios where finding as many ads as possible is the priority.
- Serves as a strong non-parametric baseline for comparison.

### 9.5.2 Class Imbalance Considerations

All models effectively handled the inherent class imbalance in the dataset, demonstrating:

- Robust performance despite unequal class distribution
- High precision across all algorithms
- Conservative prediction patterns that minimize false positives
- Effective learning from limited positive examples

## 9.6 Research Contributions

### 9.6.1 Technical Achievements

This project made several significant contributions:

- **Custom Implementation from Scratch:** By developing the core logic of three distinct machine learning algorithms without reliance on pre-built libraries, this project provides deep insights into their internal mechanics and demonstrates a fundamental understanding of the underlying statistical principles.
- **Comprehensive Evaluation:** Conducted thorough performance analysis using multiple metrics
- **Practical Application:** Demonstrated real-world applicability for advertisement detection
- **Reproducible Research:** Created well-documented, reproducible analysis framework

### 9.6.2 Statistical Significance

The results demonstrate statistically meaningful performance differences:

- All models achieved accuracy above 80%, indicating effective learning
- Performance variations reflect different algorithmic strengths
- Results provide empirical evidence for algorithm selection criteria
- Findings contribute to understanding of classification trade-offs

## 9.7 Limitations

### 9.7.1 Current Limitations

- **Feature Engineering:** Limited exploration of feature selection and dimensionality reduction
- **Hyperparameter Optimization:** Basic parameter tuning without extensive grid search
- **Cross-Validation:** The use of a simple train-test split, rather than a more robust method like k-fold cross-validation, means the performance metrics might be sensitive to the specific data partition and may not fully represent the models' generalizability.
- **Computational Efficiency:** Custom implementations may not be optimally efficient



## 9.8 Final Conclusions

### 9.8.1 Project Success

This Internet Advertisement Classification project successfully achieved its primary objectives:

- **Algorithm Implementation:** Successfully implemented three distinct machine learning algorithms
- **Performance Evaluation:** Conducted comprehensive comparative analysis
- **Practical Applicability:** Demonstrated real-world relevance for advertisement detection
- **Academic Rigor:** Followed systematic methodology with proper documentation

### 9.8.2 Key Insights

The study revealed several important insights:

1. **Algorithm Diversity:** Different algorithms excel in different performance aspects, highlighting the importance of application-specific model selection.
2. **Trade-off Analysis:** The precision-recall trade-off is clearly demonstrated, with Random Forest achieving perfect precision at the cost of lower recall, while k-NN provides a better balance.
3. **Superior Precision for Production:** Random Forest's ensemble method delivered perfect precision, making it the most reliable model for production systems where false positives are costly.
4. **Interpretability vs. Performance:** Decision Trees offer excellent interpretability while maintaining the highest accuracy, proving their value in explanatory analysis.

### 9.8.3 Impact and Significance

This work provides empirical evidence that machine learning can effectively classify internet advertisements with over 80% accuracy across all implemented methods. The findings contribute to:

- **Automated Content Filtering:** Supporting development of intelligent ad-blocking systems
- **Digital Advertising Research:** Providing baseline performance metrics for future studies
- **Machine Learning Education:** Demonstrating practical implementation of fundamental algorithms
- **Reproducible Research:** Establishing a framework for systematic algorithm comparison



#### 9.8.4 Closing Remarks

The Internet Advertisement Classification project demonstrates the practical power of machine learning in solving real-world problems. By implementing and comparing three distinct algorithms, this study provides valuable insights into the strengths and limitations of different approaches to binary classification tasks.

The results show that while no single algorithm dominates across all metrics, each method offers unique advantages that make it suitable for specific applications. This finding underscores the importance of understanding both the technical characteristics of algorithms and the practical requirements of the target application when selecting machine learning approaches.





## Data source

The data for this project is sourced from the UCI Machine Learning Repository. The dataset is titled "Internet Advertisements Data Set" and is also available on Kaggle at: <https://www.kaggle.com/datasets/uciml/internet-advertisements-data-set>

## Dataset Description

The task is to predict whether an image is an advertisement ("ad") or not ("nonad"). The dataset contains 1,559 columns, where each row represents one image tagged as "ad" or "nonad" in the last column. Columns 0 to 1,557 represent the actual numerical attributes of the images.

## Citation

Lichman, M. (2013). UCI Machine Learning Repository [http://archive.ics.uci.edu/ml]. Irvine, CA: University of California, School of Information and Computer Science.

## 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, 2024. Accessed: 2024.
- [4] GeeksforGeeks. What are the advantages and disadvantages of random forest?, 2024. Accessed: 2024.
- [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.