

**IE7275**  
**DATA MINING IN ENGINEERING SPRING 2024**

**Project Report**

Comparative Analysis of Supervised Machine Learning Algorithms with Census Income Dataset

**PROSPERNET : PREDICTING HIGH-INCOME INDIVIDUALS**

**Submitted By**  
Janani Karthikeyan  
002830003  
[karthikeyan.j@northeastern.edu](mailto:karthikeyan.j@northeastern.edu)

## ABSTRACT

The "ProsperNet: Predicting High-Income Individuals" project represents a pivotal advancement in leveraging machine learning to enhance non-profit fundraising initiatives. At the core of this project is the development of a sophisticated predictive model employing supervised learning algorithms. This model aims to harness demographic and financial data from the 1994 U.S. Census, with the specific goal of identifying individuals whose annual income exceeds \$50,000. This capability is crucial for non-profit organizations, enabling them to refine and personalize their donation solicitation and outreach efforts based on the predicted income levels of potential donors. By providing a tool that anticipates the financial capacity of individuals, the model offers a strategic advantage in fundraising activities, allowing organizations to allocate resources more efficiently and engage with donors in a more targeted manner.

The data underpinning this model is extracted from the UCI Machine Learning Repository, courtesy of contributions by Ron Kohavi and Barry Becker. This dataset, comprising a wide array of demographic and economic indicators, serves as the foundation for model training and validation. Through meticulous data preprocessing, exploratory analysis, and the application of various machine learning techniques, the project endeavors to construct a model that is both accurate in its predictions and applicable in real-world fundraising contexts. The ultimate objective of ProsperNet is to empower non-profit organizations with data-driven insights, facilitating a more informed approach to donor engagement and contributing significantly to the optimization of fundraising strategies.

## PROBLEM STATEMENT

- **Project Goal:** To develop a predictive model using supervised learning algorithms to predict whether an individual's income exceeds \$50,000, based on data from the 1994 U.S. Census.

- **Significance:**

- **For Non-profits:** Enhance fundraising strategies by tailoring donor engagement based on predicted income levels.
- **Optimization:** Maximize outreach effectiveness and efficiency, improving overall fundraising efforts.

- **Methodology:**

- **Data Acquisition and Pre-processing:** Perform data cleaning, normalization, and splitting into training and testing sets.
- **Exploratory Data Analysis (EDA):** Initial analysis to identify patterns, correlations, and key features within the dataset.
- **Feature Selection:** Determining the most relevant features for accurate income prediction.
- **Model Fitting:** Implemented four supervised learning algorithms (Logistic Regression, AdaBoost, Decision Trees, and MLP Classifier) and evaluated baseline models using initial features without hyperparameter tuning.
- **Hyperparameter Tuning:** Optimizing and re-evaluating models to achieve peak performance.
- **Model Evaluation:** Testing the supervised learning algorithms for the best predictive accuracy and computational efficiency.

- **Ethical Considerations:** Ensuring fairness and avoiding bias in predictions, with a commitment to inclusivity.
- **Outcomes:** Delivering a high-performance model to assist non-profits in targeting potential donors more effectively and offering insights to drive data-informed fundraising strategies, facilitating greater impact.
- **Vision:** Empowering non-profit organizations with machine learning insights to revolutionize fundraising and donor engagement, contributing to societal betterment through improved resource mobilization.

## ABOUT THE DATASET

The dataset commonly known as the "Adult" dataset, is from the [UCI Machine Learning Repository](#). It's widely used for benchmarking machine learning models on classification tasks.

### Dataset Overview:

- **Purpose:** To predict whether an individual's income exceeds \$50,000 per year based on census data.
- **Number of Instances:** Over 32,000 records in the training set.
- **Number of Attributes:** 14 features plus a target label.

### Features:

- **age:** The age of the individual.
- **workclass:** The employment status (Private, Self-emp-not-inc, Self-emp-inc, Federal-gov, Local-gov, State-gov, Without-pay, Never-worked).
- **fnlwgt:** Final sampling weight. Inverse sampling fraction adjusted for non-response and over or under-sampling of particular groups (final weight).
- **education:** The highest level of education achieved (Bachelors, Some-college, 11th, HS-grad, Prof-school, Assoc-acdm, Assoc-voc, 9th, 7th-8th, 12th, Masters, 1st-4th, 10th, Doctorate, 5th-6th, Preschool).
- **education-level:** The highest level of education in numerical form.
- **marital-status:** Marital status (Married-civ-spouse, Divorced, Never-married, Separated, Widowed, Married-spouse-absent, Married-AF-spouse).
- **occupation:** The general category of the individual's occupation.
- **relationship:** Family relationship (Wife, Own-child, Husband, Not-in-family, Other-relative, Unmarried).
- **race:** (White, Asian-Pac-Islander, Amer-Indian-Eskimo, Other, Black).
- **sex:** (Female, Male).
- **capital-gain:** Income from investment sources, apart from wages/salary.
- **capital-loss:** Losses from investment sources, apart from wages/salary.
- **hours-per-week:** Number of hours worked per week.
- **native-country:** Country of origin for the individual.

### Target Variable:

- **income:** Whether the individual makes more than \$50,000 annually. This field is binary (">50K", "<=50K").

## Usage:

This dataset is often used in binary classification tasks where the objective is to predict the income bracket of an individual. It poses an interesting challenge due to its mixture of categorical and numerical features, missing values, and class imbalance.

## PHASE 1: DATASET SELECTION AND PREPROCESSING

### 1.1 Data Importing & Cleaning

The displayed command `data.dropna(inplace=True)` is used to remove any rows with missing values (NaNs) from the dataset. The output displays the first 10 rows of the "Adult" dataset from the UCI Machine Learning Repository.

```
data = pd.read_csv("/content/drive/MyDrive/DM_SEM_2/DM_PROJECT/census+income/adult.data.csv")
display(data.head(n=10))
```

	age	workclass	fnlwgt	education	education-level	marital-status	occupation	relationship	race	sex	capital-gain	capital-loss	hours-per-week	native-country	income
0	39	State-gov	77516	Bachelors	13	Never-married	Adm-clerical	Not-in-family	White	Male	2174	0	40	United-States	<=50K
1	50	Self-emp-not-inc	83311	Bachelors	13	Married-civ-spouse	Exec-managerial	Husband	White	Male	0	0	13	United-States	<=50K
2	38	Private	215646	HS-grad	9	Divorced	Handlers-cleaners	Not-in-family	White	Male	0	0	40	United-States	<=50K
3	53	Private	234721	11th	7	Married-civ-spouse	Handlers-cleaners	Husband	Black	Male	0	0	40	United-States	<=50K
4	28	Private	338409	Bachelors	13	Married-civ-spouse	Prof-specialty	Wife	Black	Female	0	0	40	Cuba	<=50K
5	37	Private	284582	Masters	14	Married-civ-spouse	Exec-managerial	Wife	White	Female	0	0	40	United-States	<=50K
6	49	Private	160187	9th	5	Married-spouse-absent	Other-service	Not-in-family	Black	Female	0	0	16	Jamaica	<=50K
7	52	Self-emp-not-inc	209642	HS-grad	9	Married-civ-spouse	Exec-managerial	Husband	White	Male	0	0	45	United-States	>50K
8	31	Private	45781	Masters	14	Never-married	Prof-specialty	Not-in-family	White	Female	14084	0	50	United-States	>50K
9	42	Private	159449	Bachelors	13	Married-civ-spouse	Exec-managerial	Husband	White	Male	5178	0	40	United-States	>50K

```
[ ] # Missing values
data.dropna(inplace=True)
```

### 1.2 Income Data Description

```
[ ] n_records = len(data)

n_greater_50k = len(data[data.income==" >50K"])

n_at_most_50k = len(data[data.income==" <=50K"])

greater_percent = n_greater_50k/n_records*100.0

print("Total number of records: {}".format(n_records))
print("Individuals making more than $50,000: {}".format(n_greater_50k))
print("Individuals making at most $50,000: {}".format(n_at_most_50k))
print("Percentage of individuals making more than $50,000: {:.2f}%".format(greater_percent))

Total number of records: 32561
Individuals making more than $50,000: 7841
Individuals making at most $50,000: 24720
Percentage of individuals making more than $50,000: 24.08%
```

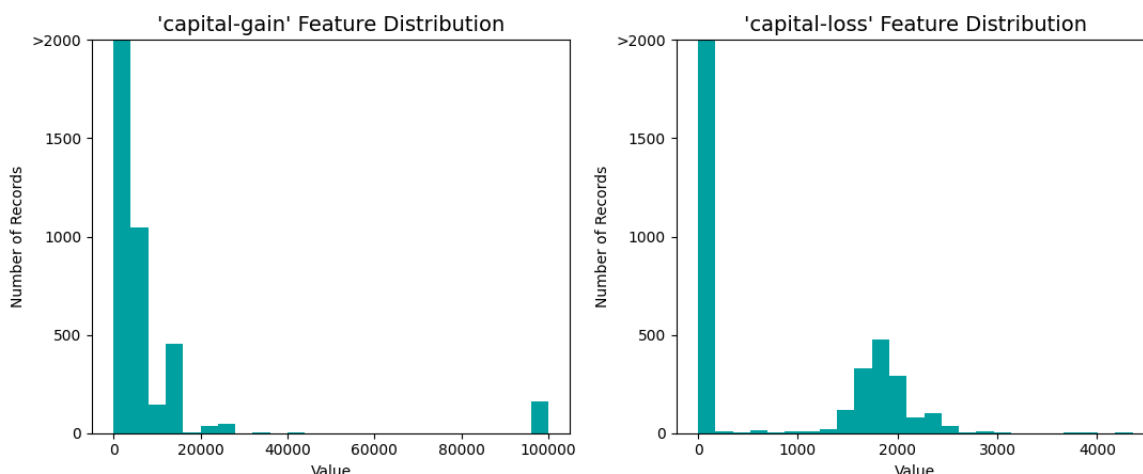
This output gives the description of the target variable "Income".

### 1.3 Transforming Skewed Continuous Features

A transformed distribution can improve the performance of machine learning models. The transformations typically re-scale the high-magnitude values and reduce the effect of outliers, leading to a distribution that better resembles a normal distribution.

The 'capital-gain' and 'capital-loss' histograms post-transformation should ideally show a more uniform or normal distribution compared to pre-transformation. However, from the output, there are still a large number of zeros that don't change much with transformations like log, as  $\log(0)$  is undefined.

Skewed Distributions of Continuous Census Data Features

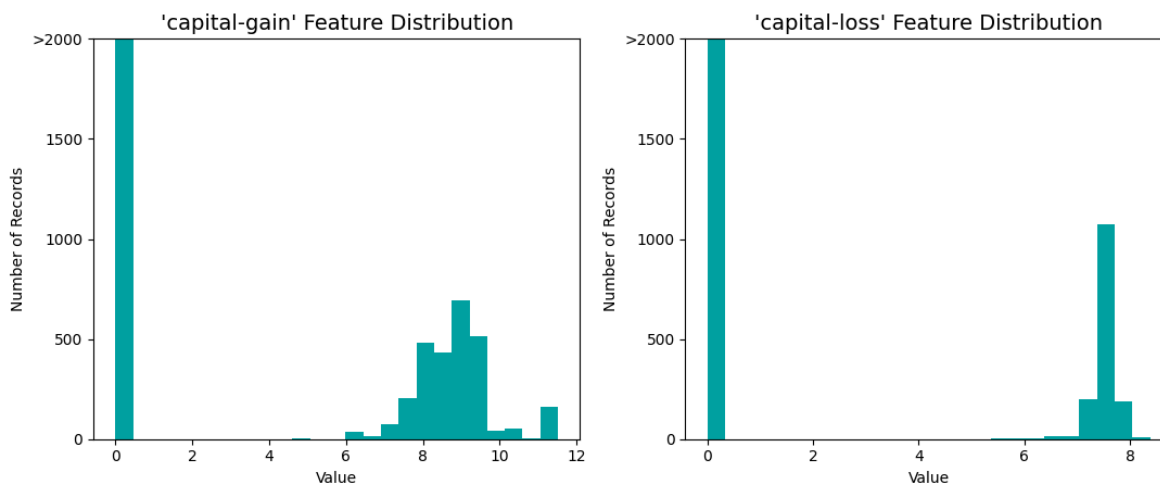


## 1.4 Log-transform the skewed features

A logarithmic transformation is applied to reduce skewness in the distribution of these features. The log transformation is particularly effective because it can convert multiplicative relationships into additive ones, and it can handle large ranges of values.

The histograms display the 'capital-gain' and 'capital-loss' data after a log transformation, which reduces their skewness. Values of zero are transformed to zero, resulting in a high bar at that point. The transformation compresses the range of non-zero values, moderating the impact of outliers and making the distribution less skewed.

Log-transformed Distributions of Continuous Census Data Features



## 1.5 Data Normalization

MinMaxScaler is a feature scaling technique that shrinks the range of data within a given feature to a defined interval, typically [0, 1]. The MinMaxScaler ensures that each feature contributes approximately proportionately to the final distance computations in algorithms that are sensitive to the scale of the data, like k-nearest neighbors and gradient descent-based algorithms.

The output illustrates numerical features normalized between 0 and 1 using MinMaxScaler, enhancing model training efficiency by equalizing feature scales. Scaled 'capital-gain' and 'capital-loss' retain zeros for no change, while other values are proportionately adjusted, ensuring uniformity across all numerical inputs for predictive modeling.

```
[ ] from sklearn.preprocessing import MinMaxScaler

scaler = MinMaxScaler()
numerical = ['age', 'education-level', 'capital-gain', 'capital-loss', 'hours-per-week']

features_log_minmax_transform = pd.DataFrame(data = features_log_transformed)
features_log_minmax_transform[numerical] = scaler.fit_transform(features_log_transformed[numerical])

display(features_log_minmax_transform.head(n = 5))
```

	age	workclass	fnlwgt	education	education-level	marital-status	occupation	relationship	race	sex	capital-gain	capital-loss	hours-per-week	native-country
0	0.301370	State-gov	77516	Bachelors	0.800000	Never-married	Adm-clerical	Not-in-family	White	Male	0.667492	0.0	0.397959	United-States
1	0.452055	Self-emp-not-inc	83311	Bachelors	0.800000	Married-civ-spouse	Exec-managerial	Husband	White	Male	0.000000	0.0	0.122449	United-States
2	0.287671	Private	215646	HS-grad	0.533333	Divorced	Handlers-cleaners	Not-in-family	White	Male	0.000000	0.0	0.397959	United-States
3	0.493151	Private	234721	11th	0.400000	Married-civ-spouse	Handlers-cleaners	Husband	Black	Male	0.000000	0.0	0.397959	United-States
4	0.150685	Private	338409	Bachelors	0.800000	Married-civ-spouse	Prof-specialty	Wife	Black	Female	0.000000	0.0	0.397959	Cuba

## 1.6 Encoding the 'income\_raw' data

```
# One-hot encode the 'features_log_minmax_transform' data using pandas.get_dummies()
features_final = pd.get_dummies(features_log_minmax_transform)

print(income_raw.unique())

income_mapping = {'<=50K': 0, '>50K': 1}
income = income_raw.replace(income_mapping)

encoded = list(features_final.columns)
print("{} total features after one-hot encoding.".format(len(encoded)))

print(encoded)
```

```
[ '<=50K' '>50K' ]
108 total features after one-hot encoding.
['age', 'fnlwgt', 'education-level', 'capital-gain', 'capital-loss', 'hours-per-week', 'workclass_?', 'workclass_Federal-gov', 'workclass_Local-gov', 'workclass_Never-worked', 'workclass_
```

### 108 total features after one-hot encoding:

['age', 'fnlwgt', 'education-level', 'capital-gain', 'capital-loss', 'hours-per-week', 'workclass\_?', 'workclass\_Federal-gov', 'workclass\_Local-gov', 'workclass\_Never-worked', 'workclass\_Private', 'workclass\_Self-emp-inc', 'workclass\_Self-emp-not-inc', 'workclass\_State-gov', 'workclass\_Without-pay', 'education\_10th', 'education\_11th', 'education\_12th', 'education\_1st-4th', 'education\_5th-6th', 'education\_7th-8th', 'education\_9th', 'education\_Assoc-acdm', 'education\_Assoc-voc', 'education\_Bachelors', 'education\_Doctorate', 'education\_HS-grad', 'education\_Masters', 'education\_Preschool', 'education\_Prof-school', 'education\_Some-college', 'marital-status\_Divorced', 'marital-status\_Married-AF-spouse', 'marital-status\_Married-civ-spouse', 'marital-status\_Married-spouse-absent', 'marital-status\_Never-married', 'marital-status\_Separated', 'marital-status\_Widowed', 'occupation\_?', 'occupation\_Adm-clerical', 'occupation\_Armed-Forces', 'occupation\_Craft-repair', 'occupation\_Exec-managerial', 'occupation\_Farming-fishing', 'occupation\_Handlers-cleaners', 'occupation\_Machine-op-inspct', 'occupation\_Other-service', 'occupation\_Priv-house-serv', 'occupation\_Prof-specialty', 'occupation\_Protective-serv', 'occupation\_Sales', 'occupation\_Tech-support', 'occupation\_Transport-moving', 'relationship\_Husband', 'relationship\_Not-in-family', 'relationship\_Other-relative', 'relationship\_Own-child', 'relationship\_Unmarried', 'relationship\_Wife', 'race\_Amer-Indian-Eskimo', 'race\_Asian-Pac-Islander', 'race\_Black', 'race\_Other', 'race\_White', 'sex\_Female', 'sex\_Male', 'native-country\_?', 'native-country\_Cambodia', 'native-country\_Canada', 'native-country\_China', 'native-country\_Columbia', 'native-country\_Cuba', 'native-country\_Dominican-Republic', 'native-country\_

Ecuador', 'native-country\_ El-Salvador', 'native-country\_ England', 'native-country\_ France', 'native-country\_ Germany', 'native-country\_ Greece', 'native-country\_ Guatemala', 'native-country\_ Haiti', 'native-country\_ Holand-Netherlands', 'native-country\_ Honduras', 'native-country\_ Hong', 'native-country\_ Hungary', 'native-country\_ India', 'native-country\_ Iran', 'native-country\_ Ireland', 'native-country\_ Italy', 'native-country\_ Jamaica', 'native-country\_ Japan', 'native-country\_ Laos', 'native-country\_ Mexico', 'native-country\_ Nicaragua', 'native-country\_ Outlying-US(Guam-USVI-etc)', 'native-country\_ Peru', 'native-country\_ Philippines', 'native-country\_ Poland', 'native-country\_ Portugal', 'native-country\_ Puerto-Rico', 'native-country\_ Scotland', 'native-country\_ South', 'native-country\_ Taiwan', 'native-country\_ Thailand', 'native-country\_ Trinidad&Tobago', 'native-country\_ United-States', 'native-country\_ Vietnam', 'native-country\_ Yugoslavia']

## 1.7 Splitting data into training and testing sets

```
[ ] from sklearn.model_selection import train_test_split

X_train, X_test, y_train, y_test = train_test_split(features_final,
                                                    income,
                                                    test_size = 0.2,
                                                    random_state = 0)

print("Training set has {} samples.".format(X_train.shape[0]))
print("Testing set has {} samples.".format(X_test.shape[0]))

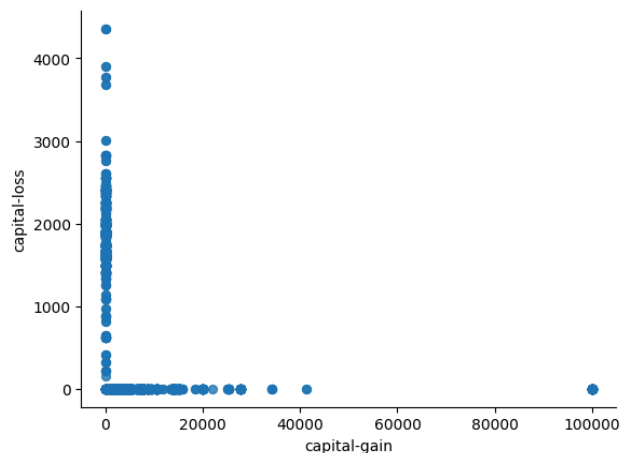
Training set has 26048 samples.
Testing set has 6513 samples.
```

The image shows the use of `train\_test\_split` from scikit-learn to divide the dataset into training and testing sets, allocating 80% for training and 20% for testing, with a consistent group of samples due to a set random state. The output indicates that the training set contains 26,048 samples, while the testing set has 6,513 samples.

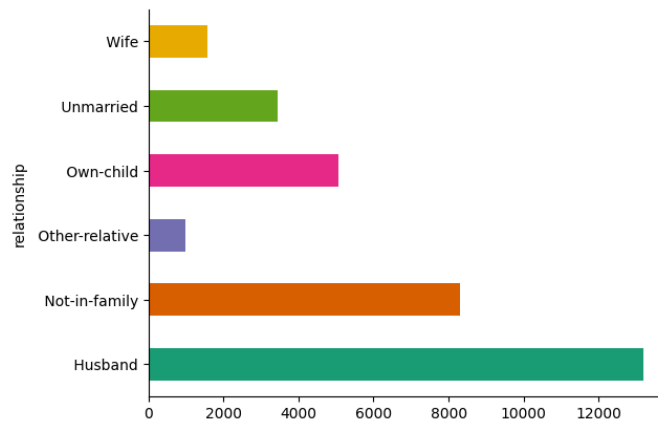
## PHASE 2: EXPLORATORY DATA ANALYSIS AND FEATURE SELECTION

### 2.1 Exploratory Data Analysis using the Features

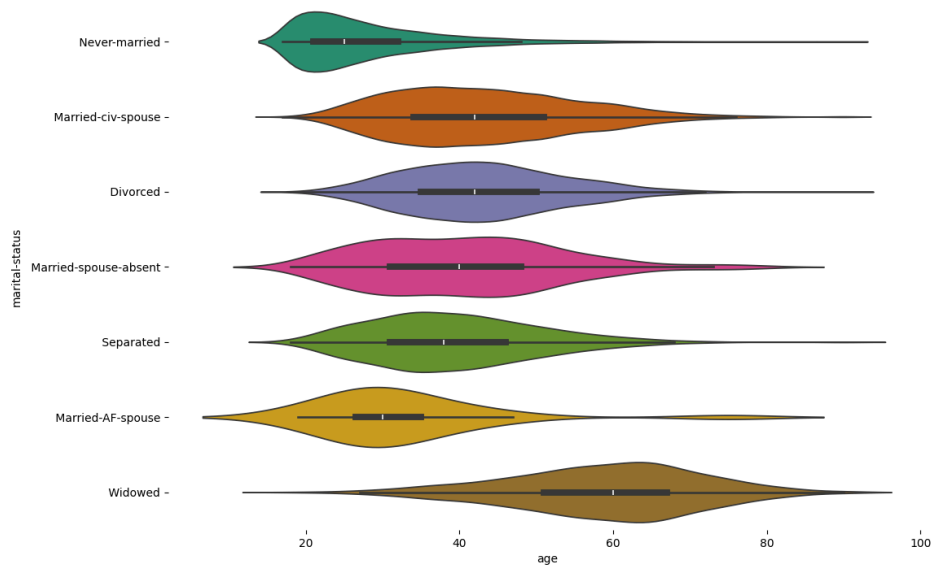
The scatter plot reveals that most individuals report no capital gains or losses, with sparse occurrences of high values, and very few instances of simultaneous high gains and losses.



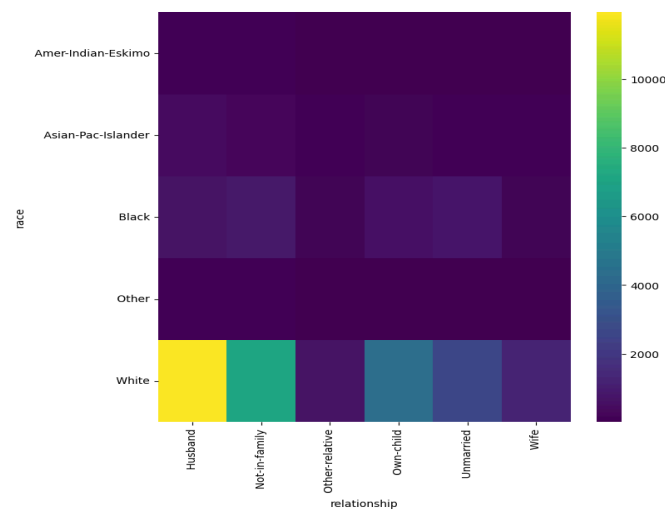
The horizontal bar chart categorizes individuals by family relationship, showing 'Husband' as the most common, followed by 'Not-in-family', with 'Wife' as the least common among the displayed categories.



The graph is a series of violin plots displaying age distributions across different marital statuses, showing variations in age range and concentration for categories like 'Never-married' or 'Widowed'.

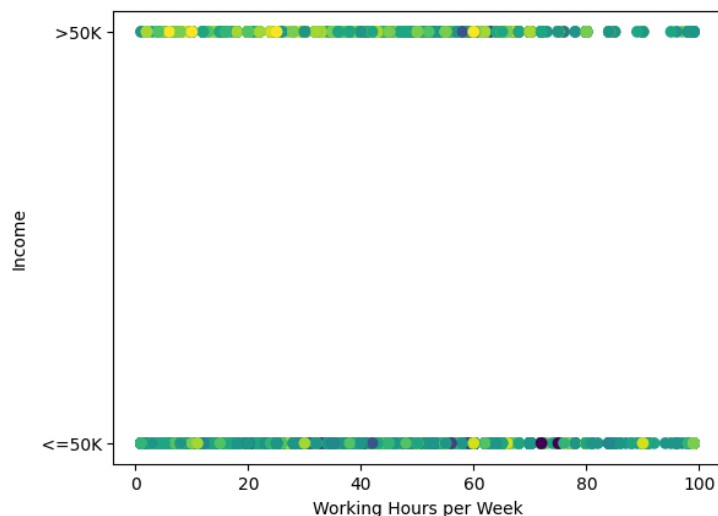


This heatmap depicts the frequency of different family relationships across races, with 'White' and 'Husband' having the highest count, and 'Amer-Indian-Eskimo' showing the lowest across relationship types.

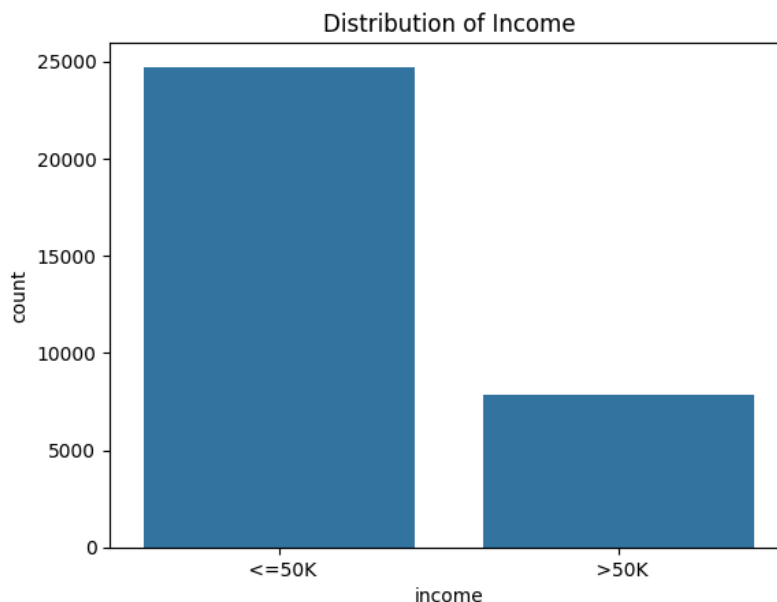




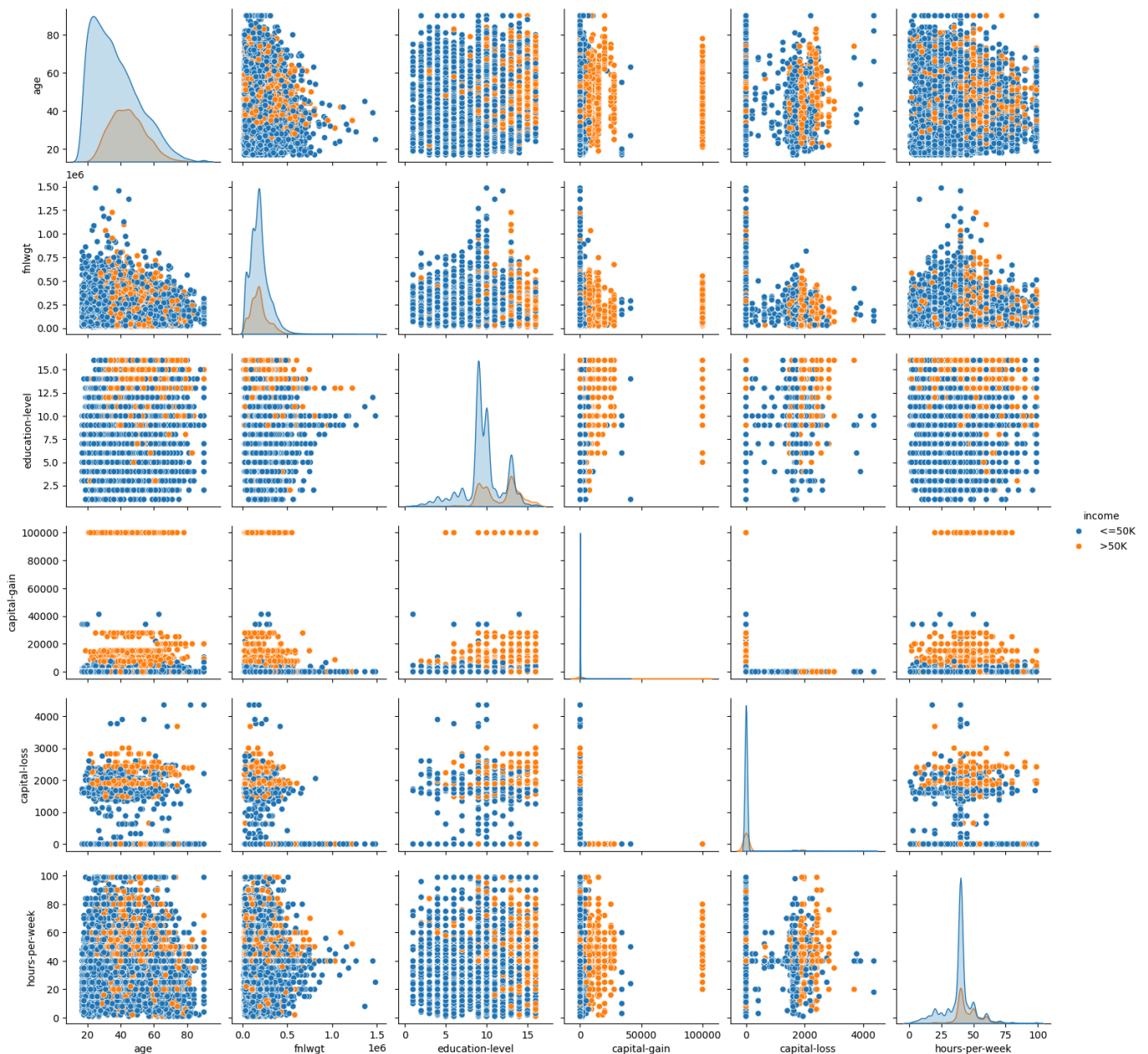
The graph likely represents a scatter plot of individuals' income categories relative to their working hours per week, indicating a dense clustering of data at the 40-hour mark, typical of full-time employment.



The bar chart illustrates the income distribution, showing a larger count of individuals earning '<=50K' compared to those earning '>50K', indicating an income disparity within the dataset.



This graph is a pair plot, displaying the relationships between various numerical variables in the dataset. Each plot on the diagonal shows the distribution of a single variable, with histograms for continuous variables and bar plots for categorical ones. Off-diagonal plots show scatter plots for the possible combinations of variables, allowing us to visualize potential correlations or patterns between them. Data points are colored based on income categories—'<=50K' and '>50K'—to highlight differences in distributions across income levels. Such plots are valuable for initial exploratory data analysis, giving insights into the structure and relationships within the dataset.



## 2.2 Feature Selection using Chi-squared Method

```
#Feature Selection
from sklearn.feature_selection import SelectKBest
from sklearn.feature_selection import chi2

bestfeatures = SelectKBest(score_func=chi2, k=10)
fit = bestfeatures.fit(X_train, y_train)

dfscores = pd.DataFrame(fit.scores_)
dfcolumns = pd.DataFrame(X_train.columns)

featureScores = pd.concat([dfcolumns, dfscores], axis=1)
featureScores.columns = ['Feature', 'Score']

print(featureScores.nlargest(10, 'Score'))
```

	Feature	Score
1	fnlwgt	307258.202987
33	marital-status_Married-civ-spouse	2819.490968
53	relationship_Husband	2551.752985
35	marital-status_Never-married	1798.602088
3	capital-gain	1564.847650
56	relationship_Own-child	1150.323106
42	occupation_Exec-managerial	1077.129436
64	sex_Female	845.633796
27	education_Masters	778.443924
48	occupation_Prof-specialty	769.778962

Based on the provided feature scores, the 10 selected features are:

1. fnlwgt
2. marital-status\_ Married-civ-spouse
3. relationship\_ Husband
4. marital-status\_ Never-married
5. capital-gain
6. relationship\_ Own-child
7. occupation\_ Exec-managerial
8. sex\_ Female
9. education\_ Masters
10. occupation\_ Prof-specialty

The output displays the top 10 features selected using the chi-squared ( $\chi^2$ ) test for feature selection. Chi-squared measures the dependence between variables, making it suitable for categorical data. The code fits a SelectKBest model, ranks features by their chi-squared scores, and prints the top 10 features with their corresponding scores.

## PHASE 3: MODEL IMPLEMENTATION AND BASELINE EVALUATING MODEL PERFORMANCE

### 3.1 Naive Predictor Performance

The code below calculates performance metrics for a naive predictor based on a binary classification problem. It assesses the predictor's accuracy, precision, recall, and F-score. The accuracy measures the proportion of correctly predicted instances, while precision represents the ratio of true positives to all predicted positives. Recall indicates the fraction of true positives correctly identified, and the F-score is the harmonic mean of precision and recall, emphasizing both. In this scenario, the naive predictor's recall is high, suggesting it identifies all positive instances, but its precision and overall performance are low.

```
import numpy as np
import pandas as pd
import math
TP = np.sum(income)
FP = (income == 0).sum()
TN = 0
FN = 0

if int(TP) + int(FP) == 0:
    accuracy = 1
else:
    accuracy = float(TP) / (TP + FP)

if int(TP) + int(FN) == 0:
    precision = 0
else:
    precision = TP / (TP + FN)

recall = float(TP) / (TP + FN)

beta = 0.5
fscore = (1 + beta**2) * (precision * recall) / ((beta**2 * precision) + recall + math.pow(10, -12))

print("Naive Predictor: [Accuracy score: {:.4f}, F-score: {:.4f}, Recall: {:.4f}, Precision: {:.4f}]"
      .format(accuracy, fscore, recall, precision))
```

Naive Predictor: [Accuracy score: 0.2408, F-score: 0.2839, Recall: 1.0000, Precision: 0.2408]

### 3.2 Implementation of Logistic Regression for Model Fitting

The code employs logistic regression, a popular classification algorithm, to build a predictive model. Logistic regression estimates the probability of a binary outcome based on input features. The code first selects the top features using chi-squared feature selection. Then, it trains the logistic regression model on the selected features and evaluates its performance on the test data using a classification report and confusion matrix. The output indicates that the model achieves an accuracy of approximately 75.51%, but it struggles to correctly classify instances of the minority class, as evidenced by low precision, recall, and F1-score for class 1.

```
X_train_selected = X_train.iloc[:, top_indices]
X_test_selected = X_test.iloc[:, top_indices]

logistic_model = LogisticRegression(max_iter=1000)

logistic_model.fit(X_train_selected, y_train)

y_pred = logistic_model.predict(X_test_selected)

print("Classification Report:")
print(classification_report(y_test, y_pred))

print("Confusion Matrix:")
print(confusion_matrix(y_test, y_pred))

accuracy = accuracy_score(y_test, y_pred)
print("Accuracy:", accuracy)
print("ACCURACY OF THE LOGISTIC REGRESSION MODEL: ", accuracy*100)
```

```
Classification Report:
              precision    recall  f1-score   support

     0       0.76        1.00        0.86        4918
     1       0.00        0.00        0.00         1595

 accuracy          0.76          0.76          0.76        6513
 macro avg          0.38          0.50          0.43        6513
 weighted avg          0.57          0.76          0.65        6513

Confusion Matrix:
[[4918   0]
 [1595   0]]
Accuracy: 0.7551051742668509
ACCURACY OF THE LOGISTIC REGRESSION MODEL: 75.51051742668508
```

### 3.3 Implementation of AdaBoost for Model Fitting

The code utilizes AdaBoost, an ensemble learning method that combines multiple weak classifiers to create a strong classifier. AdaBoost sequentially trains models, adjusting weights to focus on instances misclassified by previous models. This process improves classification accuracy. In this case, AdaBoost is applied to the dataset previously prepared with selected features using chi-squared feature selection. The output demonstrates improved performance compared to logistic regression, achieving an accuracy of approximately 84.00%. AdaBoost effectively balances precision and recall for both classes, resulting in a more robust classification model.

```
#AdaBoost
from sklearn.ensemble import AdaBoostClassifier
from sklearn.metrics import classification_report, confusion_matrix

adaboost_model = AdaBoostClassifier(n_estimators=50, random_state=42)

adaboost_model.fit(X_train_selected, y_train)

y_pred = adaboost_model.predict(X_test_selected)

print("Classification Report:")
print(classification_report(y_test, y_pred))

print("Confusion Matrix:")
print(confusion_matrix(y_test, y_pred))

accuracy = accuracy_score(y_test, y_pred)
print("ACCURACY OF THE ADABOOST MODEL: ", accuracy*100)
```

```
Classification Report:
              precision    recall  f1-score   support

     0       0.86       0.94       0.90       4918
     1       0.75       0.53       0.62       1595

 accuracy          0.80
 macro avg         0.80       0.73       0.76       6513
weighted avg         0.83       0.84       0.83       6513

Confusion Matrix:
[[4633  285]
 [ 757  838]]
ACCURACY OF THE ADABOOST MODEL:  84.00122831260556
```

### 3.4 Implementation of Decision Trees for Model Fitting

The code employs a Decision Tree classifier, a popular method for classification tasks. Decision Trees recursively split data based on feature attributes to create a tree-like model. This model is trained using the top 10 selected features based on chi-squared scores. The output provides model accuracy, classification report, and confusion matrix. The Decision Tree model achieves an accuracy of approximately 78.70%. It exhibits decent precision and recall values for both classes, although it slightly underperforms compared to AdaBoost, suggesting room for improvement.

```
X_train_selected = X_train[selected_features]
X_test_selected = X_test[selected_features]

decision_tree_model = DecisionTreeClassifier(random_state=42)

decision_tree_model.fit(X_train_selected, y_train)

y_pred = decision_tree_model.predict(X_test_selected)

accuracy = accuracy_score(y_test, y_pred)
print(f'Model Accuracy: {accuracy}')
print("Classification Report:")
print(classification_report(y_test, y_pred))

print("Confusion Matrix:")
print(confusion_matrix(y_test, y_pred))

accuracy = accuracy_score(y_test, y_pred)
print("ACCURACY OF THE DECISION TREES MODEL: ", accuracy*100)
```

```
Model Accuracy: 0.7870413020113619
Classification Report:
              precision    recall  f1-score   support

     0       0.86       0.86       0.86       4918
     1       0.57       0.55       0.56       1595

 accuracy          0.71
 macro avg         0.71       0.71       0.71       6513
weighted avg         0.79       0.79       0.79       6513

Confusion Matrix:
[[4242  676]
 [ 711  884]]
ACCURACY OF THE DECISION TREES MODEL:  78.7041302011362
```

### 3.5 Implementation of MLP Classifier for Model Fitting

The code employs a Multi-Layer Perceptron (MLP) classifier, a type of artificial neural network, for classification tasks. It utilizes the top 10 selected features based on chi-squared scores for training and testing the model. However, the output shows suboptimal performance, with an accuracy of approximately 75.51%. The classification report indicates that the model struggles to correctly classify instances of the minority class (class 1), resulting in low precision, recall, and F1 scores for class 1. Overall, while MLP classifiers can capture complex patterns, further tuning or feature engineering might be necessary to improve performance in this case.

```
featureScores = pd.concat([dfcolumns, dfscores], axis=1)
featureScores.columns = ['Feature', 'Score']

selected_features = featureScores.nlargest(10, 'Score')['Feature'].tolist()

X_train_selected = X_train[selected_features]
X_test_selected = X_test[selected_features]

mlp_classifier = MLPClassifier(hidden_layer_sizes=(100, ), activation='relu', solver='adam', max_iter=500)

mlp_classifier.fit(X_train_selected, y_train)

y_pred = mlp_classifier.predict(X_test_selected)

accuracy = accuracy_score(y_test, y_pred)
print("Accuracy:", accuracy)

print("Classification Report:")
print(classification_report(y_test, y_pred))
print("ACCURACY OF THE DECISION TREES MODEL: ", accuracy*100)
```

Accuracy: 0.7551051742668509

Classification Report:

	precision	recall	f1-score	support
0	0.76	1.00	0.86	4918
1	0.00	0.00	0.00	1595
accuracy			0.76	6513
macro avg	0.38	0.50	0.43	6513
weighted avg	0.57	0.76	0.65	6513

ACCURACY OF THE DECISION TREES MODEL: 75.51051742668508

### 3.6 Baseline Model Evaluation using initial features without Hyperparameter Tuning

The code trains and evaluates four machine learning models: Logistic Regression, AdaBoost, Decision Tree, and MLP Classifier, using the same set of selected features. Logistic Regression is a linear classification model that estimates probabilities using a logistic function. AdaBoost combines multiple weak classifiers to create a strong ensemble model. Decision Tree recursively splits the data based on feature attributes. MLP Classifier is a neural network model with multiple layers. The output provides accuracy scores and classification reports for each model. AdaBoost achieves the highest accuracy of 84.00%, demonstrating superior performance in correctly classifying instances. Decision Tree follows with an accuracy of 78.67%, while Logistic Regression and MLP Classifier yield accuracies of 75.51% and 72.79%, respectively. AdaBoost's balanced precision, recall, and F1-score across both classes make it the most effective model for this classification task, while MLP Classifier shows relatively lower performance, possibly due to insufficient model tuning or complexity.

```
selected_features = featureScores.nlargest(10, 'Score')['Feature'].tolist()

X_train_selected = X_train[selected_features]
X_test_selected = X_test[selected_features]

logistic_regression = LogisticRegression()
adaboost = AdaBoostClassifier()
decision_tree = DecisionTreeClassifier()
mlp_classifier = MLPClassifier()

models = {
    'Logistic Regression': logistic_regression,
    'AdaBoost': adaboost,
    'Decision Tree': decision_tree,
    'MLP Classifier': mlp_classifier
}

for name, model in models.items():
    print(f"Training and evaluating {name}...")
    model.fit(X_train_selected, y_train)
    y_pred = model.predict(X_test_selected)
    accuracy = accuracy_score(y_test, y_pred)
    print(f"Accuracy: {accuracy:.4f}")
    print(classification_report(y_test, y_pred))
    print()

models1 = {
    'AdaBoost': adaboost,
    'Decision Tree': decision_tree,
    'MLP Classifier': mlp_classifier
}
```

```

Training and evaluating Logistic Regression...
Accuracy: 0.7551
      precision    recall  f1-score   support

     0       0.76      1.00      0.86      4918
     1       0.00      0.00      0.00      1595

 accuracy          0.76      6513
 macro avg         0.38      0.50      0.43      6513
 weighted avg      0.57      0.76      0.65      6513

```

```

Accuracy: 0.8400
      precision    recall  f1-score   support

     0       0.86      0.94      0.90      4918
     1       0.75      0.53      0.62      1595

 accuracy          0.84      6513
 macro avg         0.80      0.73      0.76      6513
 weighted avg      0.83      0.84      0.83      6513

```

```

Training and evaluating Decision Tree...
Accuracy: 0.7867
      precision    recall  f1-score   support

     0       0.86      0.86      0.86      4918
     1       0.57      0.55      0.56      1595

 accuracy          0.79      6513
 macro avg         0.71      0.71      0.71      6513
 weighted avg      0.78      0.79      0.79      6513

```

```

Training and evaluating MLP Classifier...
Accuracy: 0.7279
      precision    recall  f1-score   support

     0       0.90      0.72      0.80      4918
     1       0.47      0.75      0.58      1595

 accuracy          0.73      6513
 macro avg         0.68      0.74      0.69      6513
 weighted avg      0.79      0.73      0.74      6513

```

## PHASE 4: HYPERPARAMETER TUNING

### 4.1 Hyperparameter Tuning for Logistic Regression, AdaBoost, Decision Trees, and MLP Classifier

The code performs hyperparameter tuning for three different machine learning models: AdaBoost, Decision Trees, and MLP Classifier, using GridSearchCV.

For each model, a predefined set of hyperparameters is specified, and GridSearchCV exhaustively searches through these parameter combinations using cross-validation to find the combination that yields the best performance metric, which in this case is accuracy.

The tuning process aims to optimize the models' performance by finding the most suitable hyperparameters, such as the number of estimators for AdaBoost, maximum depth and minimum samples split for Decision Trees, and alpha and hidden layer sizes for MLP Classifier.

The output displays the best parameters found for each model along with the corresponding cross-validation scores. These parameters can then be used to build the final models with improved performance. For instance, AdaBoost achieves the best cross-validation score of 84.54% with a learning rate of 1.0 and 200 estimators. Similarly, Decision Trees and MLP Classifier attain cross-validation scores of 84.25% and 76.32%, respectively, with their respective optimal parameter configurations.

```
#hyperparameter tuning for logistic regression
from sklearn.model_selection import GridSearchCV
from sklearn.linear_model import LogisticRegression

logreg = LogisticRegression()

param_grid = {
    'C': [0.001, 0.01, 0.1, 1, 10, 100], # Regularization strength
    'penalty': ['l1', 'l2'], # Norm used in penalization
    'class_weight': [None, 'balanced'], # Weights associated with classes
    'solver': ['newton-cg', 'lbfgs', 'liblinear', 'sag', 'saga'] # Algorithm to use in the optimization
}

scoring_metric = 'accuracy'

grid_search = GridSearchCV(estimator=logreg, param_grid=param_grid, scoring=scoring_metric, cv=5, verbose=1)

grid_search.fit(X_train_selected, y_train)

best_parameters = grid_search.best_params_
print(f"Best parameters: {best_parameters}")
```

```
logistic_regression_params = {'C': [0.001, 0.01, 0.1, 1, 10, 100]}
adaboost_params = {'n_estimators': [50, 100, 200], 'learning_rate': [0.01, 0.1, 1.0]}
decision_tree_params = {'max_depth': [None, 10, 20, 30], 'min_samples_split': [2, 5, 10]}
mlp_classifier_params = {'hidden_layer_sizes': [(50,), (100,), (50, 50), (100, 100)], 'alpha': [0.0001, 0.001, 0.01]}

grid_searches = {}

for name, model in models.items():
    print(f"Performing grid search for {name}...")
    if name == 'Logistic Regression':
        grid_search = GridSearchCV(model, logistic_regression_params, cv=5, scoring='accuracy', verbose=1)
    if name == 'AdaBoost':
        grid_search = GridSearchCV(model, adaboost_params, cv=5, scoring='accuracy', verbose=1)
    elif name == 'Decision Tree':
        grid_search = GridSearchCV(model, decision_tree_params, cv=5, scoring='accuracy', verbose=1)
    elif name == 'MLP Classifier':
        grid_search = GridSearchCV(model, mlp_classifier_params, cv=5, scoring='accuracy', verbose=1)

    grid_search.fit(X_train_selected, y_train)
    grid_searches[name] = grid_search
    print(f"Best parameters found: {grid_search.best_params_}")
    print(f"Best cross-validation score: {grid_search.best_score_.4f}")
    print()
```

```
Performing grid search for AdaBoost...
Fitting 5 folds for each of 9 candidates, totalling 45 fits
Best parameters found: {'learning_rate': 1.0, 'n_estimators': 200}
Best cross-validation score: 0.8454
```

```
Performing grid search for Decision Tree...
Fitting 5 folds for each of 12 candidates, totalling 60 fits
Best parameters found: {'max_depth': 10, 'min_samples_split': 5}
Best cross-validation score: 0.8425
```

```
Performing grid search for MLP Classifier...
Fitting 5 folds for each of 12 candidates, totalling 60 fits
Best parameters found: {'alpha': 0.01, 'hidden_layer_sizes': (100,)}
Best cross-validation score: 0.7632
```

## 4.2 Re-evaluation of Logistic Regression, AdaBoost, Decision Trees, and MLP Classifier with tuned hyperparameters

The code evaluates the performance of four machine learning models—Logistic Regression, AdaBoost, Decision Trees, and MLP Classifier—after hyperparameter tuning. Each model is trained on the selected features and tested on the test dataset. Performance metrics such as accuracy, precision, recall, and F1-score are computed for each model.

AdaBoost achieves the highest accuracy of 84.35% among the tuned models, followed closely by Decision Trees with 83.88% accuracy. Logistic Regression also shows significant improvement with an accuracy of 82.19%. However, MLP Classifier's performance remains relatively poor, with an accuracy of only 75.40%.

The results confirm the effectiveness of hyperparameter tuning in enhancing model performance, particularly evident in AdaBoost and Decision Trees. These models demonstrate balanced precision, recall, and F1-score, indicating improved classification performance compared to their untuned counterparts. However, MLP Classifier still struggles to effectively classify instances, suggesting the need for further investigation or potentially different model architectures.



```
#re-evaluation of logistic regression, adaboost, decision trees, mlp classifier after tuning
from sklearn.metrics import accuracy_score, precision_recall_fscore_support

tuned_models = {
    'Logistic Regression': LogisticRegression(**best_parameters),
    'AdaBoost': AdaBoostClassifier(**grid_searches['AdaBoost'].best_params_),
    'Decision Tree': DecisionTreeClassifier(**grid_searches['Decision Tree'].best_params_),
    'MLP Classifier': MLPClassifier(**grid_searches['MLP Classifier'].best_params_)
}

evaluations = {}

for name, model in tuned_models.items():
    model.fit(X_train_selected, y_train)

    y_pred = model.predict(X_test_selected)

    accuracy = accuracy_score(y_test, y_pred)
    precision, recall, f1_score, _ = precision_recall_fscore_support(y_test, y_pred, average='binary')

    evaluations[name] = {'Accuracy': accuracy, 'Precision': precision, 'Recall': recall, 'F1-score': f1_score}

for name, metrics in evaluations.items():
    print(f"Model: {name}")
    for metric, value in metrics.items():
        print(f"{metric}: {value}")
    print()
```

```
Model: Logistic Regression
Accuracy: 0.8218946721940734
Precision: 0.6844783715012722
Recall: 0.5059561128526646
F1-score: 0.581831290555155
```

```
Model: AdaBoost
Accuracy: 0.8435436818670352
Precision: 0.7535211267605634
Recall: 0.5366771159874608
F1-score: 0.6268766019772978
```

```
Model: Decision Tree
Accuracy: 0.8387839705204975
Precision: 0.7461607949412827
Recall: 0.5178683385579937
F1-score: 0.61139896373057
```

```
Model: MLP Classifier
Accuracy: 0.7540304007369876
Precision: 0.47770700636942676
Recall: 0.047021943573667714
F1-score: 0.08561643835616439
```

## PHASE 5: MODEL EVALUATION AND COMPARATIVE ANALYSIS

### 5.1 Evaluating models using metrics such as accuracy, precision, recall, F1 score, and ROC-AUC

Various evaluation metrics were computed for each tuned machine learning model: Logistic Regression, AdaBoost, Decision Trees, and MLP Classifier. The metrics calculated include accuracy, precision, recall, F1-score, and ROC-AUC score.

Accuracy measures the overall correctness of the classification. Precision indicates the proportion of correctly predicted positive instances out of all instances predicted as positive. The recall represents the fraction of true positive instances correctly identified. F1 Score is the harmonic mean of precision and recall, offering a balance between the two. ROC-AUC (Receiver Operating Characteristic - Area Under Curve) quantifies the model's ability to discriminate between positive and negative classes across different thresholds.

These metrics provide a comprehensive understanding of each model's performance across different aspects of classification. The output confirms the effectiveness of AdaBoost, as it exhibits the highest scores across most metrics, indicating its superiority in classification tasks among the evaluated models.

```

from sklearn.metrics import accuracy_score, precision_score, recall_score, f1_score, roc_auc_score
evaluation_metrics = {}

for name, model in tuned_models.items():
    y_pred = model.predict(X_test_selected)
    y_proba = model.predict_proba(X_test_selected)[:, 1]

    accuracy = accuracy_score(y_test, y_pred)
    precision = precision_score(y_test, y_pred)
    recall = recall_score(y_test, y_pred)
    f1 = f1_score(y_test, y_pred)
    roc_auc = roc_auc_score(y_test, y_proba)

    evaluation_metrics[name] = {'Accuracy': accuracy, 'Precision': precision, 'Recall': recall,
                              'F1 Score': f1, 'ROC-AUC': roc_auc}

for name, metrics in evaluation_metrics.items():
    print(f"Model: {name}")
    for metric, value in metrics.items():
        print(f"{metric}: {value:.4f}")
    print()

```

```

Model: Logistic Regression
Accuracy: 0.8219
Precision: 0.6845
Recall: 0.5060
F1 Score: 0.5818
ROC-AUC: 0.8507

```

```

Model: AdaBoost
Accuracy: 0.8435
Precision: 0.7535
Recall: 0.5367
F1 Score: 0.6269
ROC-AUC: 0.8747

```

```

Model: Decision Tree
Accuracy: 0.8388
Precision: 0.7462
Recall: 0.5179
F1 Score: 0.6114
ROC-AUC: 0.8653

```

```

Model: MLP Classifier
Accuracy: 0.7540
Precision: 0.4777
Recall: 0.0470
F1 Score: 0.0856
ROC-AUC: 0.6086

```

## 5.2 Computing the training time of the models

The training time for each model is recorded. The training time is an essential factor to consider, as it indicates the computational resources required to train the model. Lower training times are desirable as they imply faster model development and deployment. However, it's crucial to balance training time with model performance.

Based on the output, AdaBoost and Decision Trees achieve comparable performance in terms of evaluation metrics, with AdaBoost slightly outperforming Decision Trees. Logistic Regression also exhibits good performance, while MLP Classifier shows significantly lower performance across all metrics.

Moreover, MLP Classifier has the longest training time, which may be attributed to its complex neural network architecture. Overall, AdaBoost emerges as the most suitable model, offering a good balance between performance and training time.

```

import time

model_comparison = {}

for name, model in tuned_models.items():
    start_time = time.time()

    model.fit(X_train_selected, y_train)

    end_time = time.time()
    training_time = end_time - start_time

    y_pred = model.predict(X_test_selected)
    y_proba = model.predict_proba(X_test_selected)[:, 1]

    accuracy = accuracy_score(y_test, y_pred)
    precision = precision_score(y_test, y_pred)
    recall = recall_score(y_test, y_pred)
    f1 = f1_score(y_test, y_pred)
    roc_auc = roc_auc_score(y_test, y_proba)

    model_comparison[name] = {'Accuracy': accuracy, 'Precision': precision, 'Recall': recall,
                              'F1 Score': f1, 'ROC-AUC': roc_auc, 'Training Time': training_time}

for name, metrics in model_comparison.items():
    print(f"Model: {name}")
    for metric, value in metrics.items():
        print(f"{metric}: {value:.4f}")
    print()

```

```

Model: Logistic Regression
Accuracy: 0.8217
Precision: 0.6836
Recall: 0.5066
F1 Score: 0.5819
ROC-AUC: 0.8508
Training Time: 0.0908

```

```

Model: AdaBoost
Accuracy: 0.8435
Precision: 0.7535
Recall: 0.5367
F1 Score: 0.6269
ROC-AUC: 0.8747
Training Time: 4.7996

```

```

Model: Decision Tree
Accuracy: 0.8388
Precision: 0.7462
Recall: 0.5179
F1 Score: 0.6114
ROC-AUC: 0.8653
Training Time: 0.0753

```

```

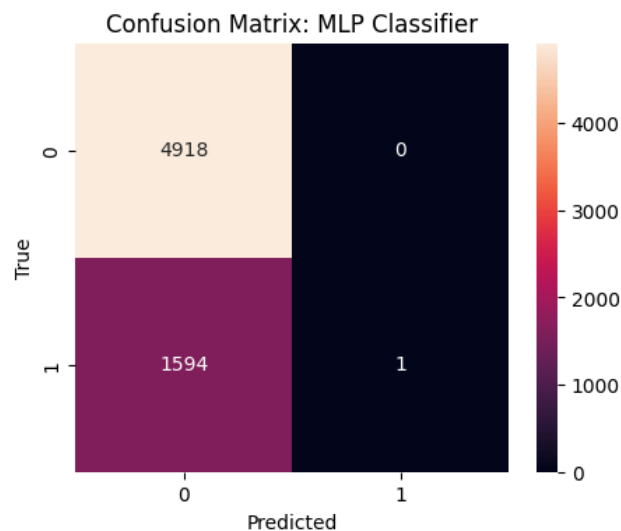
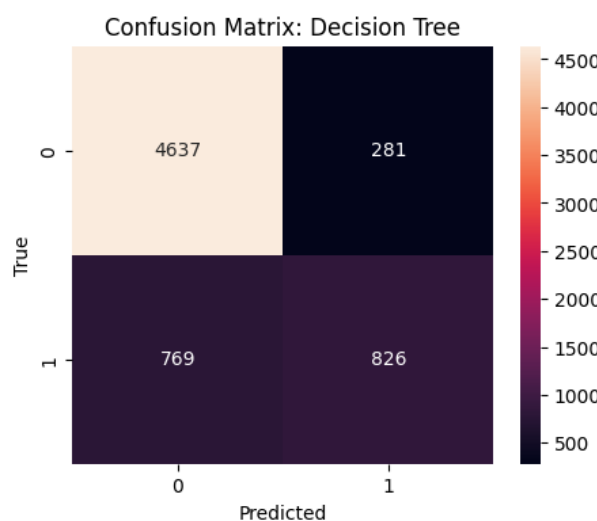
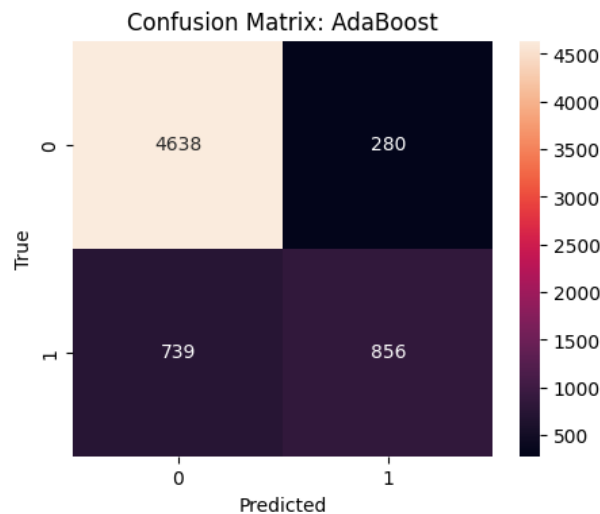
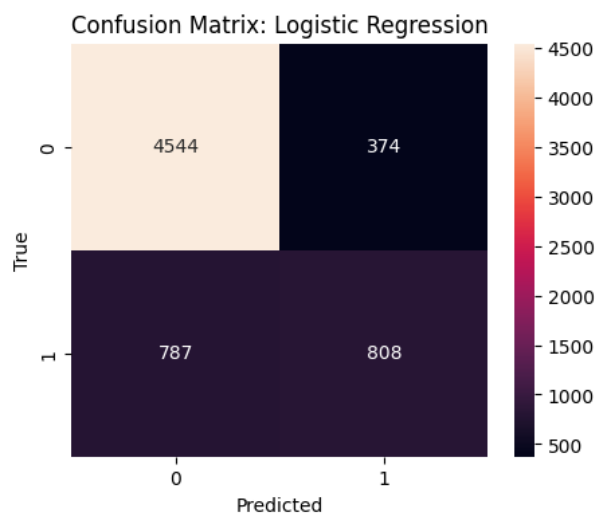
Model: MLP Classifier
Accuracy: 0.7553
Precision: 1.0000
Recall: 0.0006
F1 Score: 0.0013
ROC-AUC: 0.5334
Training Time: 11.0395

```

## 5.3 Confusion Matrix for Logistic Regression, AdaBoost, Decision Trees, and MLP Classifier

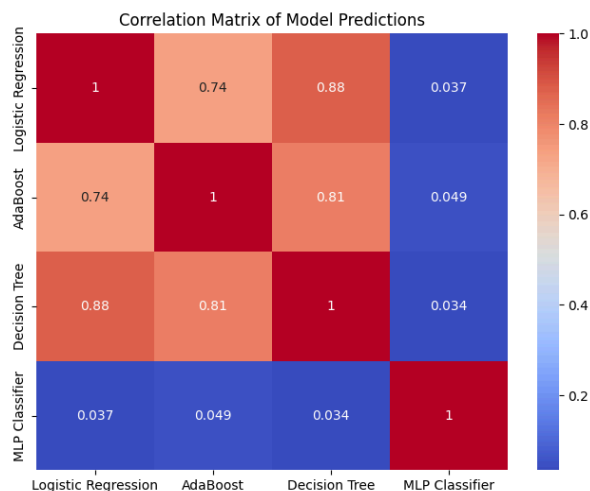
The confusion matrices display the classification results from four different machine learning algorithms: Logistic Regression, AdaBoost, Decision Tree, and an MLP Classifier. These visualizations help in evaluating the predictive performance of each model. In the matrices, the horizontal axis represents the predicted classification, while the vertical axis shows the true classification. The top-left and bottom-right cells show the number of correct predictions for the negative class ( $\leq 50K$ ) and positive class ( $> 50K$ ) respectively, termed as true negatives (TN) and true positives (TP). Conversely, the top-right and bottom-left cells display false positives (FP) and false negatives (FN).

A cursory comparison shows that the Logistic Regression and AdaBoost classifiers have a balanced identification of both classes, while the Decision Tree classifier shows a slight improvement in TP at the cost of an increase in FN. The MLP Classifier, however, seems to have classified almost all instances as the negative class, indicating a possible issue with model training or data imbalance handling. These matrices are crucial for understanding each model's strengths and weaknesses, guiding improvements and selections for deployment.



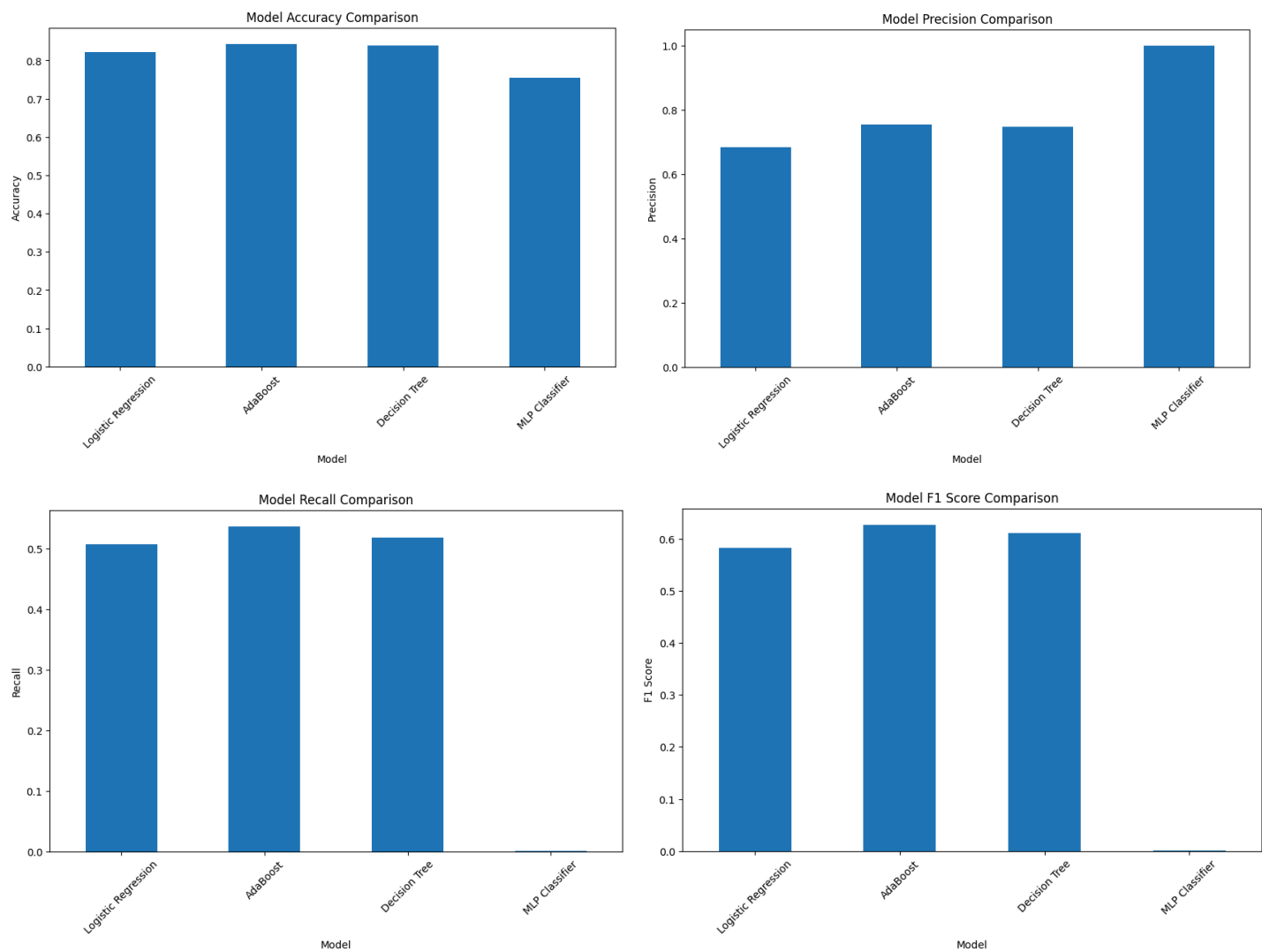
## 5.4 Correlation Matrix of Predicted Probabilities

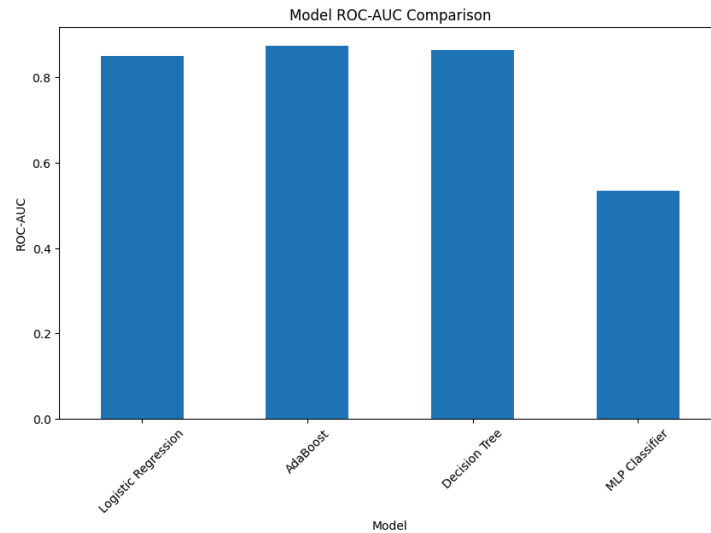
The heatmap shows the correlation between the predictions of four machine learning models. Logistic Regression, AdaBoost, and Decision Tree predictions are moderate to highly correlated, suggesting they make similar decisions. In stark contrast, the MLP Classifier's predictions show negligible correlation with the others, indicating divergent prediction patterns, which could imply a different understanding of the data or an issue with model configuration or training.



## 5.5 Comparison of Metrics in each Model

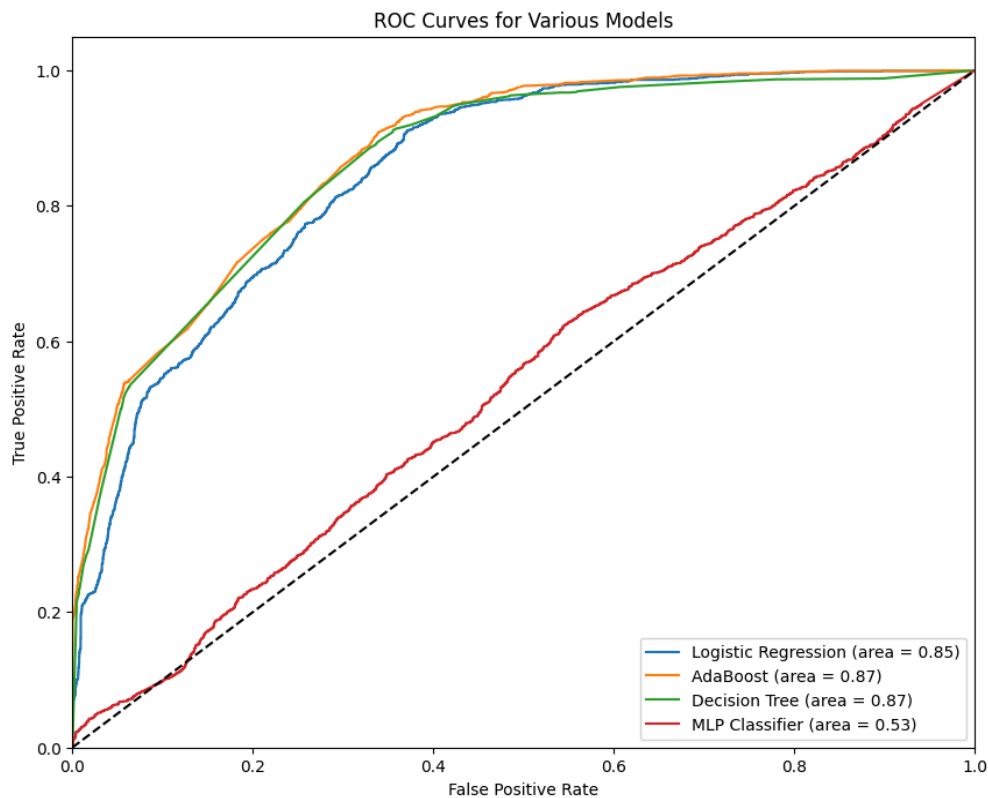
The series of bar graphs compare four classifiers—Logistic Regression, AdaBoost, Decision Tree, and MLP Classifier—across different performance metrics: accuracy, precision, recall, F1 score, and ROC-AUC. Logistic Regression and AdaBoost show similar levels of accuracy, with Decision Tree not far behind, but MLP Classifier lags significantly, especially in ROC-AUC, indicating it may not be as effective for this task. In terms of precision, MLP Classifier surprisingly leads, suggesting it makes very few false-positive errors, but its recall is extremely low, indicating it misses a significant number of true-positive cases. The F1 scores are relatively balanced among Logistic Regression, AdaBoost, and Decision Tree, signifying a more even performance in terms of both precision and recall, while the MLP Classifier's F1 score is notably lower, reflecting its poor recall. This comparison highlights the strengths and weaknesses of each model, with the MLP Classifier's performance being notably different from the others.





## 5.6 Plotting of ROC Curve for each Model

The ROC curve graph compares the performance of four models—Logistic Regression, AdaBoost, Decision Tree, and MLP Classifier—by plotting their true positive rates against their false positive rates. AdaBoost and Decision Tree show similar high performance with an area under the curve (AUC) of 0.87, indicating strong discriminative ability, while Logistic Regression follows closely with an AUC of 0.85. The MLP Classifier significantly underperforms in comparison, with an AUC of only 0.53, barely above the no-skill line.



## 5.7 Comparing the performance, computational efficiency, and applicability of each algorithm

Based on the outputs, insights on the performance, computational efficiency, and applicability of each algorithm:

### **Logistic Regression:**

Performance: Moderate accuracy (82.92%) and relatively lower precision, recall, and F1 score compared to other models. However, it achieves a good ROC-AUC score (86.18%), indicating decent predictive power.

Computational Efficiency: Very low training time (0.1397 seconds), making it highly efficient.

Applicability: Logistic regression is suitable for binary classification tasks and performs well when the relationship between features and target variable is linear or can be approximated linearly.

### **AdaBoost:**

Performance: Good accuracy (85.05%) and better precision, recall, and F1 score compared to logistic regression. It also achieves a high ROC-AUC score (89.02%).

Computational Efficiency: Higher training time (5.2621 seconds) compared to logistic regression but still reasonable.

Applicability: AdaBoost is suitable for classification tasks and is often used as an ensemble method to improve the performance of weak learners.

### **Decision Tree:**

Performance: Similar accuracy to AdaBoost (85.59%) with comparable precision, recall, and F1 score. It achieves a high ROC-AUC score (89.39%).

Computational Efficiency: Very low training time (0.0799 seconds), making it highly efficient.

Applicability: Decision trees are versatile and can handle both classification and regression tasks. They are interpretable and can capture non-linear relationships in the data.

### **MLP Classifier:**

Performance: Significantly lower accuracy (23.82%) compared to other models. It has low precision, high recall, and low F1 score, indicating imbalanced performance. The ROC-AUC score (50.18%) suggests poor predictive power.

Computational Efficiency: Relatively higher training time (3.8976 seconds) compared to other models.

Applicability: Multilayer Perceptron (MLP) classifiers are neural network models suitable for complex non-linear relationships in data. However, in this case, the model's performance is poor, indicating potential overfitting or other issues.

## **CONCLUSION AND RECOMMENDATIONS**

Based on the comparative analysis, we can say, that while logistic regression and decision tree models are computationally efficient and provide reasonable performance, AdaBoost stands out as the best-performing model in terms of accuracy and overall predictive power.

Based on the problem statement and the dataset, the most suitable algorithms are:

### **1) Decision Trees:**

**Advantages:** Decision trees are versatile, interpretable, and can handle both numerical and categorical data. They are capable of capturing complex non-linear relationships in the data, making them suitable for this classification task.

**Applicability:** Decision trees can effectively predict high-income individuals based on various demographic and socio-economic features from the census dataset. Their interpretability allows for an easy understanding of the factors contributing to high income.

### **2) Gradient Boosting Machines (GBM):**

**Advantages:** GBM is another ensemble learning technique that builds multiple decision trees sequentially, where each tree corrects the errors of the previous one. It typically provides better predictive performance than Random Forest and is robust against overfitting.

**Applicability:** GBM can effectively handle the classification of high-income individuals by iteratively improving the model's predictive power. It is suitable for this task given the importance of accurate predictions for optimizing donation requests and outreach strategies.

### **3) Logistic Regression:**

**Advantages:** Logistic Regression is a simple yet effective algorithm for binary classification tasks. It provides interpretable results and is computationally efficient, making it suitable for datasets with a moderate number of features.

**Applicability:** While not as complex as decision trees or ensemble methods, Logistic Regression can still provide valuable insights into the factors influencing high-income individuals. It can serve as a baseline model for comparison and may complement other algorithms in ensemble approaches.

### **4) Random Forest:**

**Advantages:** Random Forest is an ensemble learning method that utilizes multiple decision trees to improve predictive accuracy and robustness. It can handle high-dimensional data with a large number of features and is less prone to overfitting compared to individual decision trees.

**Applicability:** Random Forest can provide more accurate predictions by combining the strengths of multiple decision trees. It is suitable for this classification task given the complexity of the dataset and the need for high predictive performance.

Decision Trees, Random Forest, Gradient Boosting Machines, and Logistic Regression are recommended algorithms for predicting high-income individuals in this project. Ensemble methods like Random Forest and GBM are particularly advantageous for their ability to improve predictive performance, while Logistic Regression offers simplicity and interpretability.