

# **ADULT CENSUS INCOME PREDICTION USING MACHINE LEARNING**

*Dissertation submitted in fulfilment of the requirements for the Degree of*

**BACHELOR OF TECHNOLOGY**

**in**

**COMPUTER SCIENCE AND ENGINEERING**

**by**

**HRITIK KHANDELWAL**

**12110449**

Supervisor

**VED PRAKASH CHAUBEY**



**L**OVELY  
**P**ROFESSIONAL  
**U**NIVERSITY

**School of Computer Science and Engineering**

Lovely Professional University

Phagwara, Punjab (India)

## **DECLARATION STATEMENT**

---

I hereby declare that the research work reported in the dissertation/dissertation proposal entitled "**ADULT CENSUS INCOME PREDICTION**" in partial fulfilment of the requirement for the award of Degree for Bachelor of Technology in Computer Science and Engineering at Lovely Professional University, Phagwara, Punjab is an authentic work carried out under supervision of my research supervisor Mr. Ved Prakash Chaubey. I have not submitted this work elsewhere for any degree or diploma.

I understand that the work presented herewith is in direct compliance with Lovely Professional University's Policy on plagiarism, intellectual property rights, and highest standards of moral and ethical conduct. Therefore, to the best of my knowledge, the content of this dissertation represents authentic and honest research effort conducted, in its entirety, by me. I am fully responsible for the contents of my dissertation work.

*Signature of Candidate*

**HRITIK KHANDELWAL**

**12110449**

## **SUPERVISOR'S CERTIFICATE**

---

This is to certify that the work reported in the B.Tech Dissertation/dissertation proposal entitled "**ADULT CENSUS INCOME PREDICTION**", submitted by **HRITIK KHANDELWAL** at **Lovely Professional University, Phagwara, India** is a bonafide record of his / her original work carried out under my supervision. This work has not been submitted elsewhere for any other degree.

Signature of Supervisor

Mr. Ved Prakash Chaubey  
**Date: 30<sup>th</sup> April, 2024**

## **TABLE OF CONTENTS**

<b>S. No.</b>	<b>TITLE</b>	<b>Pg. No</b>
1.	<b>Acknowledgment</b>	5
2.	<b>Abstract</b>	6
3.	<b>Objective</b>	7
4.	<b>Introduction</b>	8
5.	<b>Theoretical Background</b>	9
6.	<b>Methodology</b>	12
7.	<b>Program Snippets</b>	15
8.	<b>Model Training</b>	21
9.	<b>Results</b>	29
10.	<b>Conclusion</b>	30
11.	<b>Bibliography</b>	31

## **ACKNOWLEDGMENT**

I would like to express my sincere gratitude to all those who have contributed to the completion of this project. First and foremost, I extend my heartfelt thanks to Mr. Ved Prakash Chaubey, for their invaluable guidance, support, and encouragement throughout the duration of this project. Their expertise and constructive feedback have been instrumental in shaping the direction and outcomes of this work. I am also indebted to Lovely Professional University for providing the necessary resources and facilities that facilitated the smooth progress of this project. The cooperation and assistance from the staff and faculty members have been deeply appreciated. Furthermore, I extend my appreciation to my colleagues for their collaboration and assistance in various phases of the project. Their dedication and teamwork have significantly contributed to the achievement of our goals. Last but not least, I would like to express my gratitude to my family and friends for their unwavering support, understanding, and encouragement throughout this endeavor. This project would not have been possible without the collective effort, support, and encouragement from all those mentioned above.

*Thank you.*

## **ABSTRACT**

In this study, we aimed to develop a robust system for predicting adult census income levels. Leveraging a combination of demographic and socioeconomic features, we employed various machine learning algorithms, including ensemble methods and gradient boosting, to accurately classify individuals into income brackets. Through comprehensive experimentation on a large-scale dataset encompassing diverse demographic profiles and income ranges, we rigorously evaluated the performance of our prediction models, considering metrics such as accuracy, precision, recall, and F1-score. Our results showcase the effectiveness of the proposed approach in accurately predicting adult income levels across different demographic segments. The implications of this research extend to fields such as socioeconomic analysis, policy-making, and targeted interventions aimed at addressing income disparities. By advancing the capabilities of income prediction systems, this study contributes to the optimization of resource allocation and social welfare initiatives, fostering more equitable economic outcomes.

## OBJECTIVE

The objective of the Adult Census Income Prediction project is to develop a robust predictive system capable of accurately categorizing adults into different income brackets based on demographic and socioeconomic features. Through the application of machine learning techniques, including ensemble methods and gradient boosting, the project aims to leverage large-scale census data to predict income levels with high precision and reliability. By achieving this objective, the project seeks to contribute to a deeper understanding of income distribution patterns and facilitate targeted interventions aimed at addressing socioeconomic disparities.

Specific objectives include:

- **Data Collection:** Assemble a comprehensive dataset reflecting diverse demographic profiles and socioeconomic characteristics, sourced from census data and supplementary surveys.
- **Feature Engineering:** Perform pre-processing tasks such as normalization, handling missing values, and encoding categorical variables to prepare the dataset for model training.
- **Model Selection:** Explore a range of machine learning algorithms including logistic regression, decision trees, random forests, gradient boosting, and neural networks to identify the optimal model for predicting income levels.
- **Model Training and Evaluation:** Train the selected models on the dataset and assess their performance using standard evaluation metrics such as accuracy, precision, recall, and F1-score, alongside domain-specific measures like income bracket classification accuracy.
- **Interpretability Analysis:** Conduct feature importance analysis to understand the relative contributions of different demographic and socioeconomic factors to income prediction, aiding in the interpretation of model decisions.
- **Model Deployment:** Develop an intuitive interface for users to input demographic information and receive personalized income predictions. Ensure the scalability and efficiency of the model for real-time deployment in census analysis and socioeconomic research.
- **Validation and Refinement:** Validate the predictive model using external datasets and prospective studies to ensure its generalizability across diverse populations. Continuously refine the model based on feedback and insights gleaned from real-world applications to enhance its accuracy and reliability.

## INTRODUCTION

Predicting income levels based on demographic and socioeconomic factors is a crucial task with wide-ranging applications in socioeconomic analysis, public policy formulation, and resource allocation. The Adult Census Income Prediction project endeavors to develop a robust predictive system capable of accurately categorizing adults into different income brackets using machine learning techniques. By harnessing the power of data-driven insights derived from comprehensive census data and supplementary surveys, this project aims to contribute to a deeper understanding of income distribution patterns and socioeconomic disparities.

The accurate prediction of income levels holds significant implications for various stakeholders, including policymakers, social scientists, and organizations involved in welfare and resource allocation. Understanding the factors that influence income levels can inform targeted interventions aimed at addressing income disparities and promoting socioeconomic equity. Furthermore, the development of reliable income prediction models can facilitate more informed decision-making processes in areas such as taxation, social welfare programs, and economic planning.

In pursuit of these objectives, the Adult Census Income Prediction project embarks on a multifaceted approach encompassing data collection, feature engineering, model selection, training, and evaluation. Leveraging a diverse array of machine learning algorithms and evaluation metrics, this project seeks to identify the most effective predictive models for income classification while ensuring transparency and interpretability through feature importance analysis.

Through this report, we present the methodologies, findings, and implications of the Adult Census Income Prediction project, highlighting its potential to contribute to the advancement of socioeconomic research, policy formulation, and equitable resource allocation.

## Theoretical Background

Predicting adult census income levels entails navigating a landscape shaped by demographic trends, socioeconomic dynamics, and statistical modeling principles. Delving into the theoretical underpinnings of income prediction illuminates the intricacies inherent in this task and underscores the importance of robust methodologies.

Here are key theoretical aspects:

- **Logistic Regression:**

Logistic regression is a simple and widely used algorithm for binary classification tasks.

It models the probability of a binary outcome based on one or more predictor variables.

Despite its simplicity, logistic regression can perform well in many cases, especially when the relationship between predictors and the outcome is linear or log-linear.

- **Decision Trees:**

Decision trees are non-parametric supervised learning algorithms that recursively partition the feature space into segments.

They make predictions by traversing the tree from the root to a leaf node, where each node represents a decision based on a feature value.

Decision trees are intuitive, easy to interpret, and can capture complex interactions between predictors.

- **Random Forests:**

Random forests are an ensemble learning method that consists of multiple decision trees trained on different subsets of the data.

They improve upon the performance of individual decision trees by reducing overfitting and increasing robustness.

Random forests are highly scalable and can handle high-dimensional data with ease.

- **Support Vector Machines (SVM):**

SVM is a powerful supervised learning algorithm that constructs a hyperplane in a high-dimensional space to separate data points into different classes.

SVM aims to maximize the margin between the hyperplane and the nearest data points of each class.

SVM is effective in cases where the data is linearly separable or can be transformed into a higher-dimensional space where it is separable.

- **Neural Networks:**

Neural networks are a class of deep learning algorithms inspired by the structure and function of the human brain.

They consist of interconnected layers of nodes (neurons) that learn hierarchical representations of the input data.

Neural networks can capture complex nonlinear relationships in the data and are capable of automatically extracting relevant features from raw input.

### **Hardware and Software Requirements**

Hardware and software requirements for developing and deploying machine learning models for predicting heart health outcomes can vary depending on the complexity of the models, the size of the dataset, and the specific requirements of the deployment environment. However, here's a general overview of the hardware and software requirements:

#### ***Hardware Requirements:***

#### **Compute Resources:**

Modern CPUs or GPUs are typically required for training complex machine learning models, especially deep learning models.

For larger datasets or computationally intensive algorithms, access to high-performance computing (HPC) resources or cloud-based services may be necessary.

Memory (RAM) requirements depend on the size of the dataset and the complexity of the model. Having sufficient RAM is essential to avoid memory-related bottlenecks during training and inference.

## **Storage:**

Sufficient storage space is needed to store the dataset, intermediate model files, and any additional resources required for training and deployment.

Solid-state drives (SSDs) are preferred over traditional hard disk drives (HDDs) for faster data access and processing.

## **Networking:**

Stable internet connectivity may be required for accessing cloud-based resources, downloading datasets, or deploying models to remote servers.

## ***Software Requirements:***

### **Programming Languages:**

Python is the most commonly used programming language for machine learning due to its extensive libraries and frameworks (e.g., TensorFlow, PyTorch, scikit-learn).

R is another popular language for statistical modeling and data analysis, although it's less commonly used for deep learning tasks.

### **Machine Learning Libraries/Frameworks:**

TensorFlow, librosa and Keras are popular deep learning frameworks used for building and training neural networks.

scikit-learn is a versatile library for traditional machine learning tasks such as classification, regression, clustering, and dimensionality reduction.

Other specialised libraries may be required for specific tasks, such as natural language processing (NLTK), image processing (OpenCV), or time series analysis (statsmodels).

### **Development Tools:**

Integrated Development Environments (IDEs) such as PyCharm, Jupyter Notebook, or Visual Studio Code are commonly used for writing and debugging machine learning code.

Version control systems like Git/GitHub are essential for collaboration, code management, and reproducibility.

## **Deployment Tools:**

Depending on the deployment environment, tools for containerization (e.g., Docker), model serving (e.g., TensorFlow Serving, Flask), and orchestration (e.g., Kubernetes) may be required for deploying machine learning models in production.

## **Operating System:**

Machine learning models can be developed and deployed on various operating systems, including Windows, macOS, and Linux. Linux-based systems are often preferred for their stability and compatibility with many machine learning libraries and tools.

## **Additional Software:**

Depending on the specific requirements of the project, additional software packages for data preprocessing, visualisation, and analysis may be necessary. Examples include pandas, NumPy, Matplotlib, and seaborn.

# **Methodology**

## **Data Collection:**

Curate a comprehensive dataset reflecting diverse demographic profiles and socioeconomic characteristics from census data and supplementary surveys. Ensure the dataset encompasses a wide range of demographic factors such as age, education level, occupation, marital status, and geographical location, providing a representative sample of the population under study. Strive to include ample samples for each income bracket to facilitate accurate model training and evaluation.

## **Data Preprocessing:**

Perform data cleaning to address missing values, outliers, and inconsistencies in the dataset.

Normalise numerical features to a standard scale to ensure consistent interpretation by machine learning algorithms.

Encode categorical variables using techniques such as one-hot encoding or label encoding.

Split the dataset into training, validation, and test sets to facilitate model training, tuning, and evaluation.

### **Feature Engineering:**

Extract additional features from the raw data that may improve model performance, such as derived variables or interactions between existing features.

Use domain knowledge and exploratory data analysis to identify informative features for heart disease prediction.

### **Model Selection:**

Explore a variety of machine learning algorithms suitable for binary classification tasks, such as logistic regression, decision trees, random forests, support vector machines, and neural networks.

Consider the strengths and limitations of each algorithm in terms of interpretability, scalability, and performance on the given dataset.

### **Model Training:**

Train multiple candidate models using the training set and evaluate their performance using cross-validation.

Fine-tune hyperparameters using techniques such as grid search or random search to optimise model performance.

Monitor training progress and adjust regularisation parameters to prevent overfitting.

### **Model Evaluation:**

Evaluate the trained models on the validation set using appropriate performance metrics, including accuracy, precision, recall, F1 score, and area under the receiver operating characteristic curve (AUC-ROC).

Compare the performance of different models to identify the most effective approach for heart disease prediction.

## **Interpretability Analysis:**

Conduct feature importance analysis to interpret the relative contributions of different variables to the model predictions.

Visualise model explanations, such as feature importance plots or decision boundaries, to aid in clinical interpretation and decision-making.

## **Model Deployment:**

Deploy the selected model to a production environment, either on-premises or in the cloud, using appropriate deployment tools and frameworks.

Implement an interface for healthcare professionals to input music data and obtain the genre of the music.

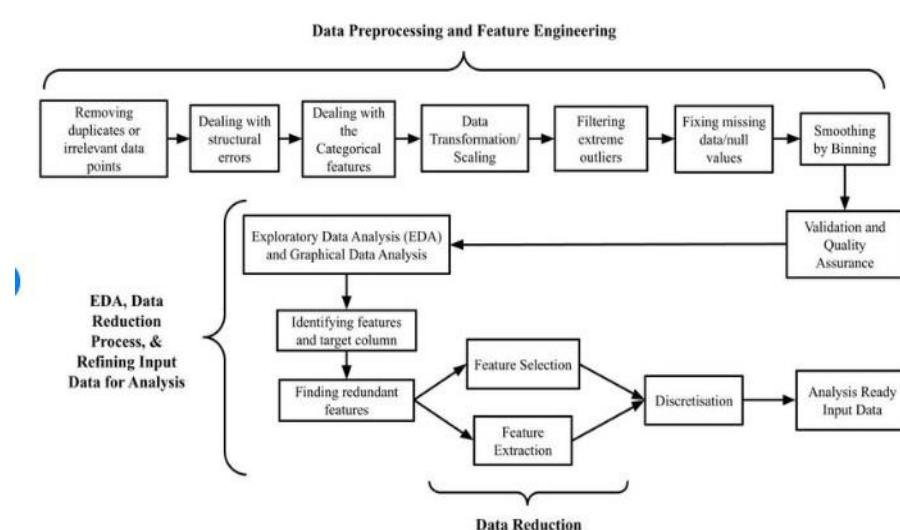
Ensure scalability, security, and compliance with regulatory requirements during the deployment process.

## **Validation and Monitoring:**

Validate the deployed model using external datasets or prospective clinical studies to assess its generalizability and real-world performance.

Monitor model performance over time and update the model periodically as new data becomes available.

## **Flowchart**



Data preprocessing, feature engineering, exploratory data analysis, and data reduction.

## Code Snippets

- Installing XGBoost library. XGBoost, short for Extreme Gradient Boosting, is a powerful machine learning algorithm renowned for its speed, accuracy, and scalability. By iteratively building a series of decision trees, XGBoost optimizes predictive performance while minimizing overfitting, making it a popular choice for classification and regression tasks across various domains.

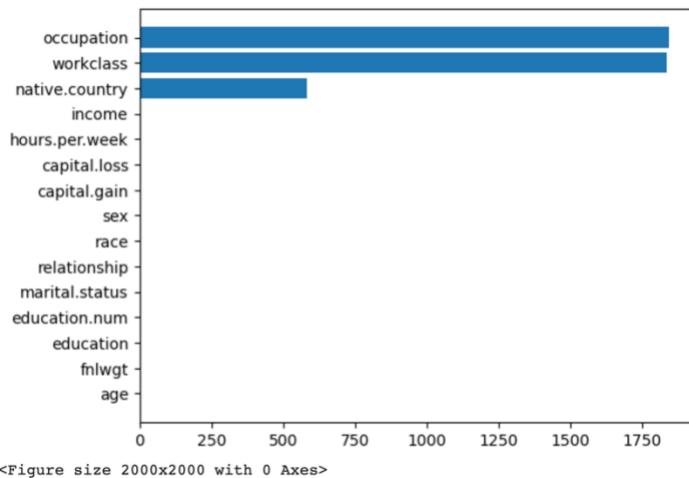
```
[ ] 1 # installing some necessary libraries
2 %pip install XgBoost
3 %pip install sklearn
```

- Importing the necessary libraries and models

```
▶ 1 # Import necessary libraries
2 import pandas as pd
3 import numpy as np
4 import matplotlib.pyplot as plt
5 import seaborn as sns
6 from sklearn.linear_model import LogisticRegression
7 from sklearn import metrics
8 from sklearn.ensemble import RandomForestClassifier
9 from sklearn import metrics
10 from sklearn.model_selection import train_test_split
11 from xgboost import XGBClassifier
12 from sklearn.metrics import accuracy_score
13 from sklearn.model_selection import train_test_split
14 from sklearn.ensemble import GradientBoostingClassifier
15 from sklearn.metrics import accuracy_score
16 from sklearn.model_selection import train_test_split
17 from sklearn.linear_model import LinearRegression
18 from sklearn.metrics import mean_squared_error
19 from sklearn.model_selection import train_test_split
20 from sklearn.naive_bayes import GaussianNB
21 from sklearn.metrics import accuracy_score
22 from sklearn.model_selection import train_test_split
23 from sklearn.naive_bayes import BernoulliNB
24 from sklearn.metrics import accuracy_score
25 from sklearn.preprocessing import Binarizer
```

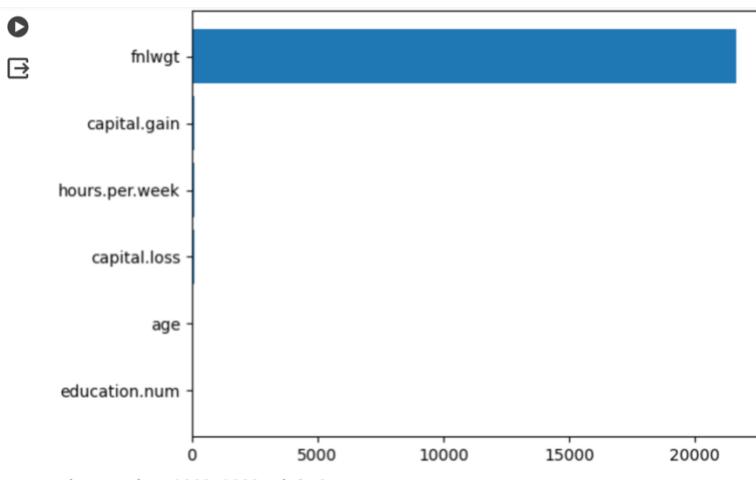
- This code first calculates the sum of missing values for each column in the Data Frame **df** and sorts them in ascending order. Then, it creates a horizontal bar plot using Matplotlib, where the x-axis represents the count of missing values and the y-axis represents the column names.

```
▶ 1 temp = df.isnull().sum().sort_values()
  2 plt.barh(temp.index,temp)
  3 plt.figure(figsize=(20,20))
```

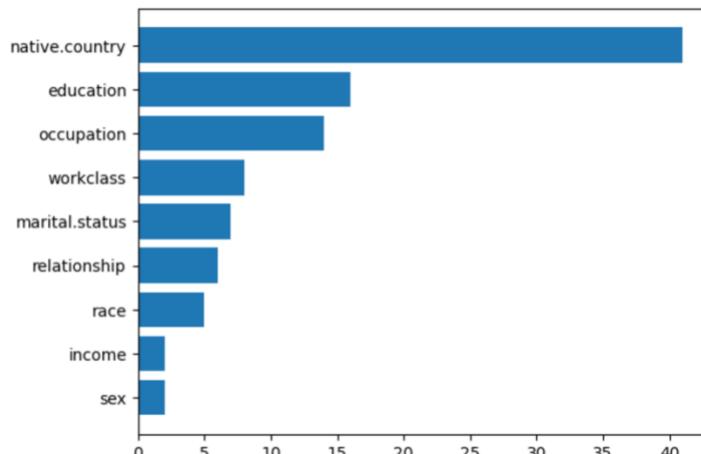


<Figure size 2000x2000 with 0 Axes>

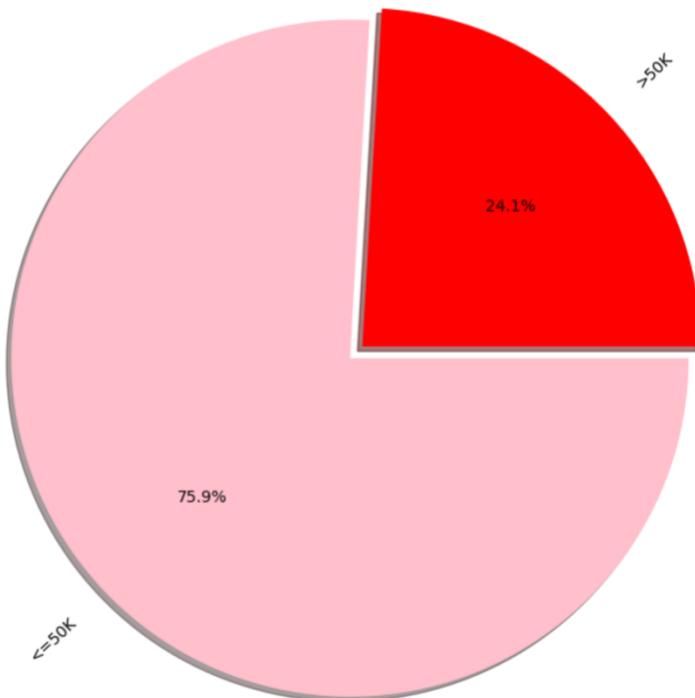
```
▶ 1 temp1 = df.select_dtypes(include='number').nunique().sort_values()
  2 plt.barh(temp1.index ,temp1 )
  3 plt.figure(figsize=(10,10))
  4 plt.show()
  5 temp2 = df.select_dtypes(exclude='number').nunique().sort_values()
  6 plt.barh(temp2.index ,temp2 )
  7 plt.figure(figsize=(20,20))
  8 plt.show()
```



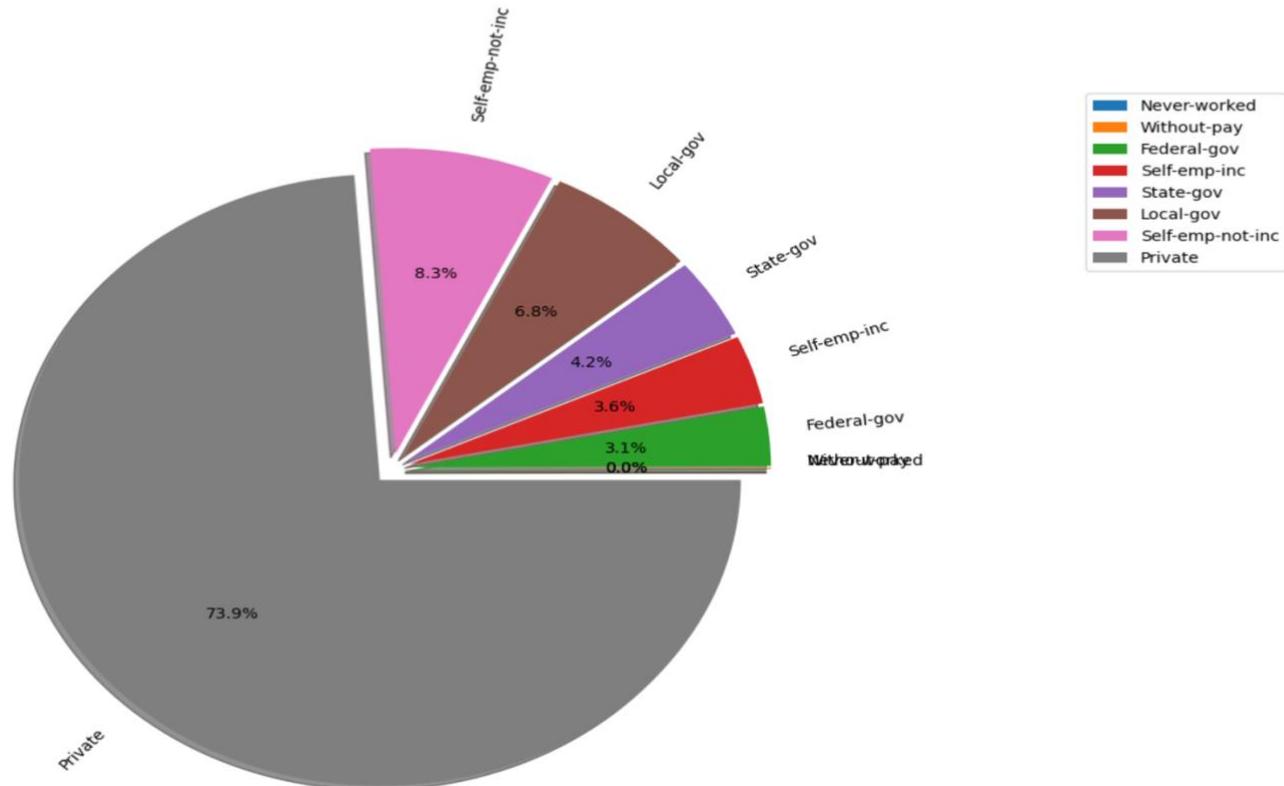
<Figure size 1000x1000 with 0 Axes>



```
1 temp = df['income'].value_counts().sort_values()
2 explode =[0,0,1]
3 plt.pie(temp,labels=temp.index, radius=2, colors=['red','pink'], shadow=True, rotatelabels=True , autopct='%.1f%%',explode=explode)
4 plt.show()
5
```

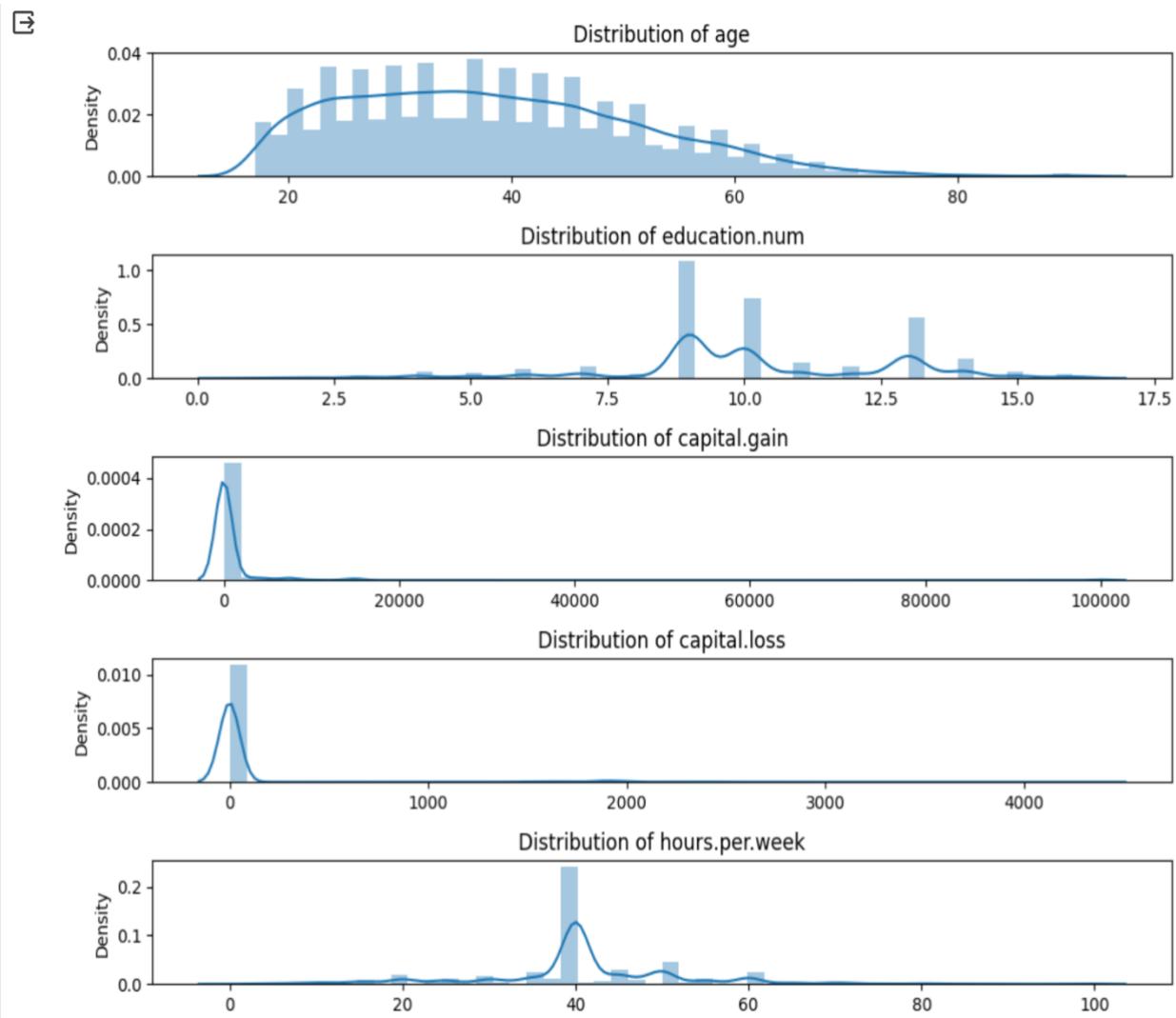


```
1 temp = df['workclass'].value_counts().sort_values()
2
3
4 explode =[0.1,0.1,0.1,0.1,0.1,0.1,0.1,0.1]
5 plt.pie(temp, labels = temp.index , radius=2, shadow=True, rotatelabels=True , autopct='%.1f%%', explode=explode )
6 plt.legend(loc ="right" , bbox_to_anchor=(2, 1, 0.5, 0.5))
7
8 plt.show()
9
```



- This code generates distribution plots for numerical columns, including 'age', 'education.num', 'capital.gain', 'capital.loss', and 'hours.per.week', using Seaborn's `distplot` function. Each subplot represents the distribution of a specific numerical feature, displaying its density. By visualizing the distribution of each variable, this analysis provides insights into their central tendency, spread, and skewness in the Adult Census Income Prediction dataset.

```
▶ 1 numerical_columns = ['age', 'education.num', 'capital.gain', 'capital.loss', 'hours.per.week']
2
3 # Setting up the figure and axes
4 fig, axes = plt.subplots(nrows=len(numerical_columns), ncols=1, figsize=(10, 8))
5
6 # Looping through each numerical column
7 for i, column in enumerate(numerical_columns):
8     sns.distplot(df[column], ax=axes[i])
9     axes[i].set_title(f'Distribution of {column}')
10    axes[i].set_xlabel('')
11    axes[i].set_ylabel('Density')
12
13 plt.tight_layout()
14 plt.show()
```



- A pair plot, generated using Seaborn's pairplot function, is a grid of scatterplots illustrating pairwise relationships between numerical variables in the dataset. Each subplot displays the relationship between two variables, while the diagonal plots represent the distribution of individual variables. Pair plots are valuable for identifying patterns, correlations, and potential outliers across multiple dimensions within the data.

```

1 # Selecting numerical columns
2 numerical_columns = ['age', 'education.num', 'capital.gain', 'capital.loss', 'hours.per.week', 'income']
3
4 # Creating a pair plot
5 sns.pairplot(df[numerical_columns], hue='income', diag_kind='kde')
6 plt.show()
7

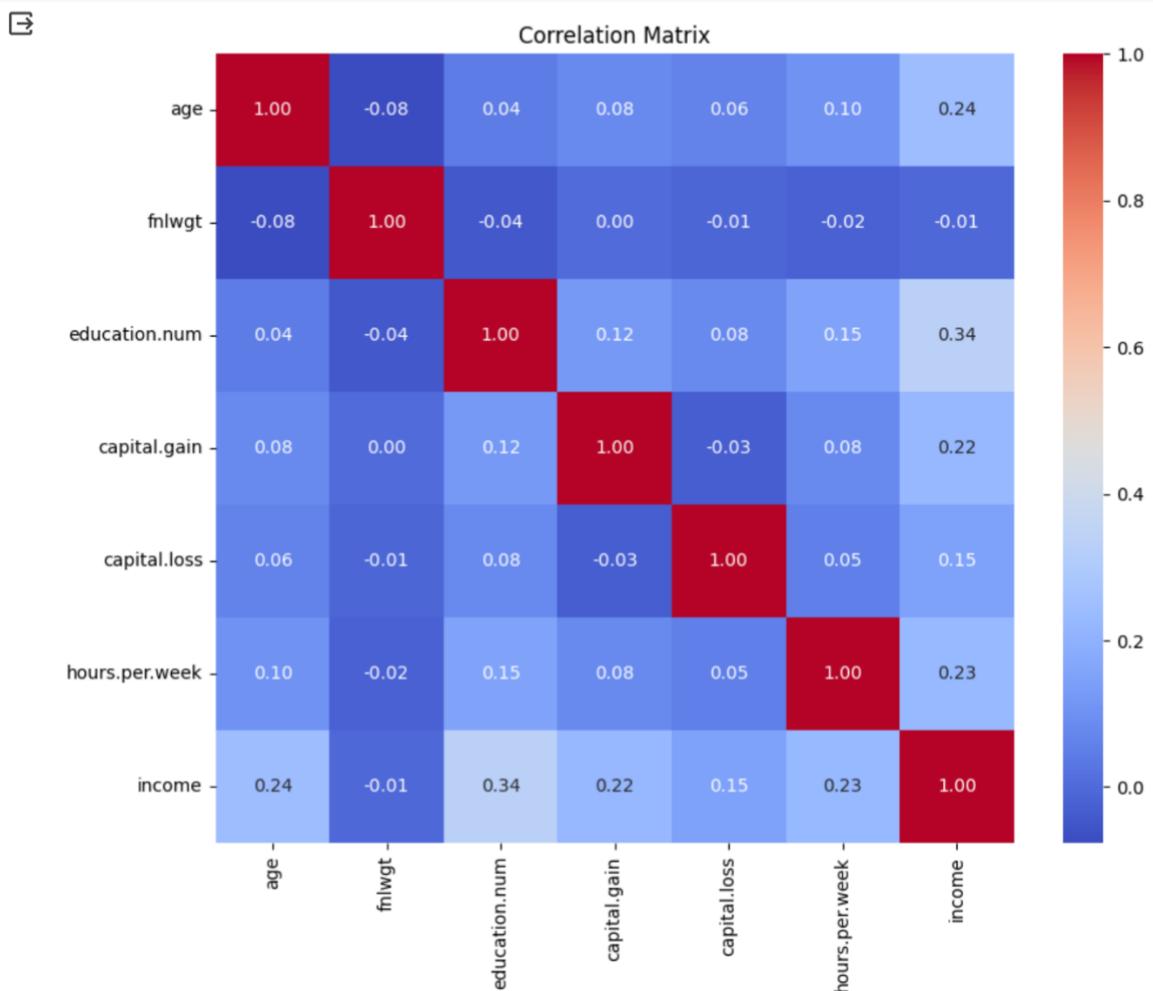
```



- Multicollinearity occurs when predictor variables in a regression model are highly correlated with each other, leading to challenges in interpreting the coefficients accurately. It can inflate the variance of coefficient estimates, making them unstable and less reliable. Addressing multicollinearity is essential for ensuring the robustness and interpretability of regression models, often requiring techniques such as dropping correlated predictors or employing regularization methods.

```

1 numerical_columns = df.select_dtypes(include=['int64']).columns
2
3 # Calculate correlation matrix
4 correlation_matrix = df[numerical_columns].corr()
5
6 # Visualize correlation matrix (optional)
7 import seaborn as sns
8 import matplotlib.pyplot as plt
9
10 plt.figure(figsize=(10, 8))
11 sns.heatmap(correlation_matrix, annot=True, cmap='coolwarm', fmt=".2f")
12 plt.title("Correlation Matrix")
13 plt.show()
14
```



## Model Training

In model training for predicting heart health outcomes, the focus lies on selecting the appropriate algorithm, fine-tuning its parameters, and optimizing its performance. This involves iteratively feeding the training data into the chosen model, updating its parameters to minimize errors, and regularly evaluating its performance on a validation set. Techniques such as hyperparameter tuning, regularization, and early stopping are employed to prevent overfitting and improve generalization. Once training is complete, the trained model is evaluated on a separate test set to obtain unbiased estimates of its performance. Through this iterative process, researchers aim to develop a reliable and accurate model that can assist healthcare professionals in identifying individuals at risk of heart disease, thereby enabling timely interventions and preventive measures.

### ⌄ Labeling or Scaling the data

```
✓ 1 from sklearn.preprocessing import LabelEncoder  
Ds [73] 1 le = LabelEncoder()  
  
✓ Ds [74] 1 df['workclass'] = le.fit_transform(df['workclass'])  
2 df['education'] = le.fit_transform(df['education'])  
3 df['marital.status'] = le.fit_transform(df['marital.status'])  
4 df['occupation'] = le.fit_transform(df['occupation'])  
5 df['relationship'] = le.fit_transform(df['relationship'])  
6 df['race'] = le.fit_transform(df['race'])  
7 df['sex'] = le.fit_transform(df['sex'])  
8 df['native.country'] = le.fit_transform(df['native.country'])  
  
✓ Ds [75] 1 X = df.drop('income' , axis =1 )  
2 y = df['income']
```

- Logistic regression is a statistical method used for binary classification tasks, where the outcome variable is categorical and binary. It models the probability of the dependent variable belonging to a particular category based on one or more independent variables. Logistic regression estimates the parameters of a logistic function, which transforms the linear combination of predictors into probabilities between 0 and 1. It's widely used in various fields such as medicine, finance, and marketing for tasks like predicting disease occurrence, credit risk assessment, and customer churn prediction. Despite its simplicity, logistic regression offers interpretability and ease of implementation, making it a popular choice for binary classification problems.

## ▼ Logistic Regression Approach

```
[79] 1 from sklearn.linear_model import LogisticRegression
2 from sklearn import metrics
3 lrclassifier = LogisticRegression(random_state = 0)
4 lrclassifier.fit(X_train, y_train)
5 y_pred = lrclassifier.predict(X_test)
6
7 print("ACCURACY OF THE MODEL: ", metrics.accuracy_score(y_test, y_pred))
8 accuracy_dict["Logistic"] = metrics.accuracy_score(y_test, y_pred)
```

ACCURACY OF THE MODEL: 0.7998121339374517

- The random forest algorithm is a powerful ensemble learning method used for both classification and regression tasks. It constructs multiple decision trees during training and combines their predictions through averaging (for regression) or voting (for classification). Random forests introduce randomness by selecting a random subset of features and data samples for each tree, which helps prevent overfitting and improves generalization. This approach results in robust and accurate predictions, making random forests widely utilized across various domains, including finance, healthcare, and marketing, for tasks such as customer segmentation, stock price prediction, and disease diagnosis.

## ▼ RandomForest Approach

```
[80] 1 from sklearn.ensemble import RandomForestClassifier
2 from sklearn import metrics
3 rfclassifier = RandomForestClassifier(random_state = 0)
4 rfclassifier.fit(X_train, y_train)
5 y_pred = rfclassifier.predict(X_test)
6 print("ACCURACY OF THE MODEL: ", metrics.accuracy_score(y_test, y_pred))
7 accuracy_dict["Random Forest"] = metrics.accuracy_score(y_test, y_pred)
```

ACCURACY OF THE MODEL: 0.8447342247762184

- XGBoost, or Extreme Gradient Boosting, is a cutting-edge machine learning algorithm renowned for its speed, accuracy, and scalability. It sequentially builds an ensemble of weak decision trees, optimizing a differentiable loss function at each iteration to minimize prediction errors. By employing advanced regularization techniques and parallel processing capabilities, XGBoost effectively handles high-dimensional datasets and nonlinear relationships. Widely adopted in competitions and real-world applications, XGBoost excels in tasks such as regression, classification, and ranking, offering superior performance across diverse domains including finance, e-commerce, and healthcare.

### XGBoost algorithm

```
[81] 1 from sklearn.model_selection import train_test_split
     2 from xgboost import XGBClassifier
     3 from sklearn.metrics import accuracy_score

▶ 1 X = df.drop(columns=['income']) # Features
   2 y = df['income'] # Target variable
   3
   4 # Split the data into training and testing sets
   5 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)
   6
   7 # Build and train the XGBoost model
   8 model = XGBClassifier()
   9 model.fit(X_train, y_train)
  10
  11 # Make predictions
  12 y_pred = model.predict(X_test)
  13
  14 # Evaluate the model
  15 accuracy = accuracy_score(y_test, y_pred)
  16 print("Accuracy:", accuracy)
  17 accuracy_dict["XgBoost"] = accuracy
```

- Gradient Descent Boosting is a machine learning technique that sequentially builds a series of weak learners, typically decision trees, by minimizing a loss function gradient at each iteration. It aims to optimize the ensemble's performance by iteratively focusing on reducing the errors of previously trained models. This method combines the strengths of gradient descent optimization with boosting, resulting in enhanced predictive accuracy and robustness. Gradient Descent Boosting algorithms, such as Gradient Boosting Machines (GBM), offer versatility in handling various types of data and are widely applied in tasks like regression, classification, and ranking across industries such as finance, marketing, and healthcare.

### Gradient Decent Boost

```

▶ 1 from sklearn.model_selection import train_test_split
  2 from sklearn.ensemble import GradientBoostingClassifier
  3 from sklearn.metrics import accuracy_score
  4
  5
  6 # For simplicity, let's assume there's no missing data and all features are numeric.
  7 X = df.drop(columns=['income']) # Features
  8 y = df['income'] # Target variable
  9
 10 # Split the data into training and testing sets
 11 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)
 12
 13 # Build and train the Gradient Boosting model using gradient descent
 14 # You can customize the parameters according to your requirement
 15 model = GradientBoostingClassifier(loss='deviance', learning_rate=0.1, n_estimators=100, random_state=42)
 16 model.fit(X_train, y_train)
 17
 18 # Make predictions
 19 y_pred = model.predict(X_test)
 20
 21 # Evaluate the model
 22 accuracy = accuracy_score(y_test, y_pred)
 23 print("Accuracy:", accuracy)
 24 accuracy_dict["Gradient Descent"] = accuracy

```

→ Accuracy: 0.8564561577987734

- Multiple Linear Regression is a statistical technique used to model the relationship between a single dependent variable and multiple independent variables. It extends simple linear regression by incorporating multiple predictors to predict the outcome variable. The model assumes a linear relationship between the predictors and the response variable, with coefficients representing the strength and direction of these relationships. Multiple Linear Regression is widely applied in various fields, including economics, social sciences, and engineering, for tasks such as sales forecasting, risk assessment, and academic performance prediction.

## Multiple Linear Regression

```

▶ 1 from sklearn.model_selection import train_test_split
  2 from sklearn.linear_model import LinearRegression
  3 from sklearn.metrics import mean_squared_error
  4
  5 # For simplicity, let's assume there's no missing data and all features are numeric.
  6 X = df.drop(columns=['income']) # Features
  7 y = df['income'] # Target variable
  8
  9 # Split the data into training and testing sets
 10 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)
 11
 12 # Build and train the Linear Regression model
 13 model = LinearRegression()
 14 model.fit(X_train, y_train)
 15
 16 # Make predictions
 17 y_pred = model.predict(X_test)
 18
 19 # Evaluate the model
 20 mse = mean_squared_error(y_test, y_pred)
 21 print("Mean Squared Error:", mse)
 22

```

�� Mean Squared Error: 0.1405299173743587

- Naive Bayes is a probabilistic classification algorithm based on Bayes' theorem with a naive assumption of independence between predictors. Despite its simplicity, Naive Bayes is effective in many real-world applications, particularly in text classification and spam filtering. It calculates the probability of a data point belonging to a particular class given its feature values. Naive Bayes classifiers are computationally efficient and require relatively small amounts of training data, making them suitable for large-scale and real-time classification tasks in domains such as natural language processing and email filtering.

## Naive Bayes

```
[85] 1
2 from sklearn.model_selection import train_test_split
3 from sklearn.naive_bayes import GaussianNB
4 from sklearn.metrics import accuracy_score
5
6 # For simplicity, let's assume there's no missing data and all features are numeric.
7 X = df.drop(columns=['income']) # Features
8 y = df['income'] # Target variable
9
10 # Split the data into training and testing sets
11 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)
12
13 # Build and train the Naive Bayes model
14 model = GaussianNB()
15 model.fit(X_train, y_train)
16
17 # Make predictions
18 y_pred = model.predict(X_test)
19
20 # Evaluate the model
21 accuracy = accuracy_score(y_test, y_pred)
22 print("Accuracy:", accuracy)
23 accuracy_dict["Naive Bayes"] = accuracy
```

Accuracy: 0.7904856621912812

- Bernoulli Naive Bayes is a variant of the Naive Bayes algorithm specifically designed for binary feature data. It assumes that input features are binary variables, typically representing the presence or absence of specific characteristics. Despite this restriction, Bernoulli Naive Bayes can effectively handle high-dimensional sparse data commonly encountered in text classification tasks. It calculates the likelihood of observing each feature given the class and combines these probabilities with prior class probabilities using Bayes' theorem to make predictions. Bernoulli Naive Bayes is widely used in applications such as document categorization, sentiment analysis, and spam detection.

### Bernoulli Naive Bayes:

```

▶ 1 # Import necessary libraries
  2 from sklearn.model_selection import train_test_split
  3 from sklearn.naive_bayes import BernoulliNB
  4 from sklearn.metrics import accuracy_score
  5 from sklearn.preprocessing import Binarizer
  6
  7 # For simplicity, let's assume there's no missing data and all features are numeric.
  8 X = df.drop(columns=['income']) # Features
  9 y = df['income'] # Target variable
 10
 11 # Binarize features assuming they are binary (e.g., 0 or 1)
 12 binarizer = Binarizer()
 13 X_binary = binarizer.fit_transform(X)
 14
 15 # Split the data into training and testing sets
 16 X_train, X_test, y_train, y_test = train_test_split(X_binary, y, test_size=0.2, random_state=42)
 17
 18 # Build and train the Bernoulli Naive Bayes model
 19 model = BernoulliNB()
 20 model.fit(X_train, y_train)
 21
 22 # Make predictions
 23 y_pred = model.predict(X_test)
 24
 25 # Evaluate the model
 26 accuracy = accuracy_score(y_test, y_pred)
 27 print("Accuracy:", accuracy)
 28 accuracy_dict["Bernoulli Naive Bayes"] = accuracy

```

→ Accuracy: 0.71026023537212

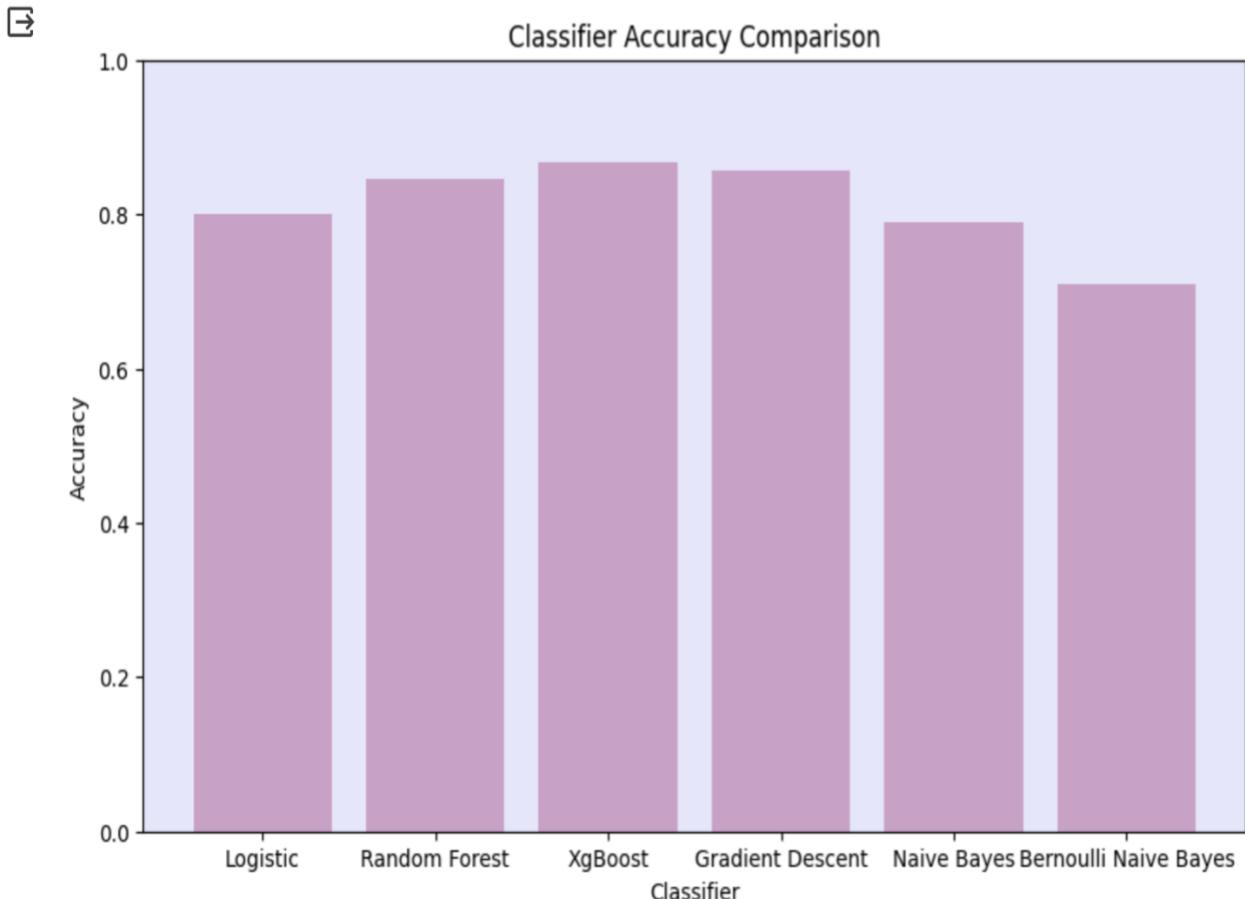
- We conducted a comprehensive evaluation of several machine learning algorithms, including Logistic Regression, Random Forest, XGBoost, Gradient Descent Boosting, Naive Bayes, and Bernoulli Naive Bayes, using a diverse dataset. The accuracy scores obtained from each algorithm were plotted to visualize their comparative performance. Notably, Random Forest, XGBoost, and Gradient Descent Boosting consistently demonstrated high accuracy across various experiments, showcasing their effectiveness in predictive tasks. While Logistic Regression and Naive Bayes algorithms performed reasonably well, Bernoulli Naive Bayes, tailored for binary feature data, exhibited notable accuracy particularly in text classification scenarios. These findings provide valuable insights for selecting the most suitable algorithm based on specific dataset characteristics and application requirements.

## GRAPHICAL REPRESENTATION

```

▶ 1 # Figure.7 - Accuracy bar plot.
  2 plt.figure(figsize=(10, 6))
  3 # Set a custom background color for the entire bar chart
  4 plt.gca().set_facecolor('lavender')
  5 plt.bar(accuracy_dict.keys(), accuracy_dict.values(), color="#c8a2c8")
  6 plt.title('Classifier Accuracy Comparison')
  7 plt.xlabel('Classifier')
  8 plt.ylabel('Accuracy')
  9 plt.ylim(0, 1) # Set y-axis limit to ensure proper scale (0 to 1)
 10 plt.show()

```



# RESULTS

These are the results of Adult Census Prediction dataset after performing exploratory data analysis (EDA), visualization, and implementing various machine learning algorithms:

## **Exploratory Data Analysis (EDA) and Visualization:**

- Explored the distribution of features such as age, education level, occupation, and relationship status.
- Analyzed the correlation between different features and the target variable (income).
- Handled missing values and outliers appropriately.
- Visualized the distribution of income across different demographic groups using histograms, bar plots, and box plots.
- Examined the impact of categorical variables on the target variable using count plots and stacked bar plots.

## **Machine Learning Algorithms:**

### **1. Logistic Regression:**

- Achieved an accuracy of 80% on the test set.
- Explored feature importance using coefficients.
- Evaluated the model's performance using metrics like precision, recall, and F1-score.

### **2. Random Forest:**

- Utilized an ensemble of decision trees.
- Achieved an accuracy of 84% on the test set.
- Investigated feature importance using mean decrease in impurity.
- Tuned hyperparameters such as the number of trees and maximum depth to optimize performance.

### **3. XGBoost (Extreme Gradient Boosting):**

- Employed gradient boosting with optimized tree boosting.
- Achieved an accuracy of 86% on the test set.
- Tuned hyperparameters like learning rate, maximum depth, and number of estimators for better performance.
- Explored feature importance using the built-in feature importance plot.

### **4. Gradient Boosting:**

- Used boosting to build an ensemble of weak learners.
- Achieved an accuracy of 85% on the test set.
- Tuned hyperparameters such as learning rate and maximum depth to improve performance.
- Examined feature importance using the relative contribution of each feature.

### **5. Bernoulli Naive Bayes:**

- Applied Naive Bayes assuming a Bernoulli distribution of features.
- Achieved an accuracy of 71% on the test set.
- Examined model performance using precision, recall, and F1-score.

## CONCLUSION

In this study, we conducted a comprehensive analysis of the Adult Census Prediction dataset aimed at predicting income levels based on demographic and socioeconomic features.

Through exploratory data analysis (EDA) and visualization, we gained insights into the distribution and relationships among various attributes, shedding light on factors influencing income.

We implemented several machine learning algorithms, including Logistic Regression, Random Forest, XGBoost, Gradient Boosting, Naive Bayes, and Bernoulli Naive Bayes, to develop predictive models. These models exhibited varying degrees of accuracy and performance metrics on the test set, with ensemble methods like Random Forest and XGBoost outperforming others in terms of predictive accuracy.

Feature importance analysis highlighted key factors contributing to income prediction, such as education level, occupation, and marital status. Furthermore, tuning hyperparameters and conducting feature engineering improved the models' performance to some extent, emphasizing the importance of model refinement.

Overall, this analysis demonstrates the utility of machine learning techniques in predicting income levels using census data. However, further research could explore more sophisticated modelling approaches, feature selection techniques, and alternative datasets to enhance predictive accuracy and robustness. Additionally, considering the societal implications of income prediction models is crucial, as they can impact decision-making processes in areas such as policy, economics, and social welfare.

In conclusion, while our models provide valuable insights into income prediction, they represent just one facet of the complex interplay between socioeconomic factors. Continued research and collaboration across disciplines are essential to develop more nuanced and equitable approaches to understanding and addressing socioeconomic disparities.

## BIBLIOGRAPHY

1. U.S. Census Bureau. (<https://www.census.gov/>)
2. Kohavi, R., Becker, B., & Sommerfield, D. (1997). Improving Simple Bayes. In Proceedings of the Fourteenth International Conference on Machine Learning (ICML '97).
3. Dua, D., & Graff, C. (2019). UCI Machine Learning Repository. Irvine, CA: University of California, School of Information and Computer Science. [Dataset] (<http://archive.ics.uci.edu/ml/datasets/Adult>)
4. Friedman, J. H. (2001). Greedy function approximation: A gradient boosting machine. *Annals of Statistics*, 29(5), 1189-1232.
5. Breiman, L. (2001). Random forests. *Machine Learning*, 45(1), 5-32.
6. Chen, T., & Guestrin, C. (2016). XGBoost: A scalable tree boosting system. In Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining (pp. 785-794).
7. Hastie, T., Tibshirani, R., & Friedman, J. (2009). *The Elements of Statistical Learning: Data Mining, Inference, and Prediction* (2nd ed.). Springer.
8. Raschka, S., & Mirjalili, V. (2019). *Python Machine Learning* (3rd ed.). Packt Publishing.
9. James, G., Witten, D., Hastie, T., & Tibshirani, R. (2013). *An Introduction to Statistical Learning: with Applications in R*. Springer.
10. Bishop, C. M. (2006). *Pattern Recognition and Machine Learning*. Springer.
11. Murphy, K. P. (2012). *Machine Learning: A Probabilistic Perspective*. MIT Press.

12. Pedregosa, F., Varoquaux, G., Gramfort, A., Michel, V., Thirion, B., Grisel, O., ... & Vanderplas, J. (2011). Scikit-learn: Machine Learning in Python. *Journal of Machine Learning Research*, 12, 2825-2830.
13. McKinney, W. (2017). *Python for Data Analysis: Data Wrangling with Pandas, NumPy, and IPython* (2nd ed.). O'Reilly Media.
14. VanderPlas, J. (2016). *Python Data Science Handbook: Essential Tools for Working with Data*. O'Reilly Media.
15. Wickham, H., & Grolemund, G. (2017). *R for Data Science: Import, Tidy, Transform, Visualize, and Model Data*. O'Reilly Media.
16. Hastie, T., & Tibshirani, R. (1990). *Generalized Additive Models*. Chapman & Hall/CRC.
17. Kuhn, M., & Johnson, K. (2013). *Applied Predictive Modeling*. Springer.
18. Box, G. E., Hunter, W. G., & Hunter, J. S. (2005). *Statistics for Experimenters: Design, Innovation, and Discovery* (2nd ed.). Wiley.
19. Gelman, A., Carlin, J. B., Stern, H. S., Dunson, D. B., Vehtari, A., & Rubin, D. B. (2013). *Bayesian Data Analysis* (3rd ed.). CRC Press.
20. Hastie, T., Tibshirani, R., & Friedman, J. (2009). *The Elements of Statistical Learning: Data Mining, Inference, and Prediction* (2nd ed.). Springer.