



Assignment 4

Group 4:

Group names:

1. Amira Abu Issa
2. Aya Metwally
3. Heba Mostafa

Part 1: Numerical Questions

(a) Please build a decision tree by using Gini Children (i.e., $Gini = \sum_{i=1}^n p_i (1 - p_i)$, where NC is the number of classes).

a

our Disition is "Hiking"

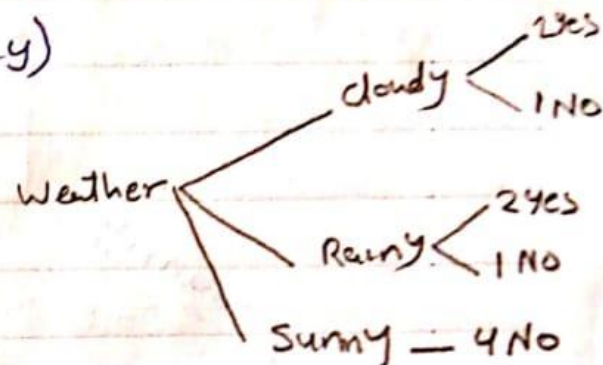
First we calculate Gini of "weather"

- Gini (Hiking, weather = cloudy)

$$= 1 - [(P_{yes})^2 + (P_{no})^2]$$

$$= 1 - \left[\left(\frac{2}{3}\right)^2 + \left(\frac{1}{3}\right)^2 \right]$$

$$= \frac{4}{9}$$



- Gini (Hiking, weather = Sunny)

$$= 1 - \left(\frac{4}{4}\right)^2 = 0$$

- Gini (Hiking, weather = Rainy)

$$= 1 - \left[\left(\frac{2}{3}\right)^2 + \left(\frac{1}{3}\right)^2 \right]$$

$$= \frac{4}{9}$$

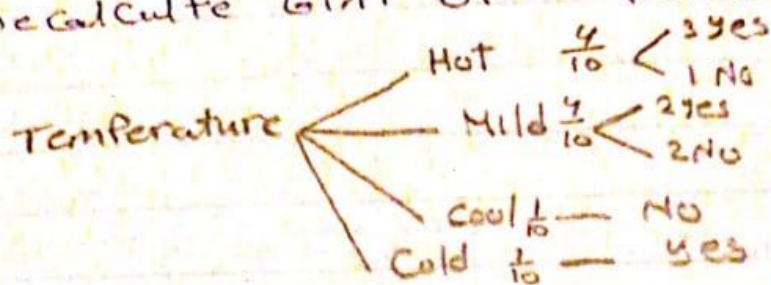
Gini children (weather)

$$= Gini(cloudy) * P_{cloudy} + Gini(Sunny) * P_{sunny} + Gini(Rainy) * P_{Rainy}$$

$$= \left(\frac{4}{9} * \frac{3}{10} \right) + \left(0 * \frac{4}{10} \right) + \left(\frac{4}{9} * \frac{3}{10} \right)$$

$$= 0.267$$

- we calculate Gini of "Temperature"



$$\text{Gini (Hot)} = 1 - \left[\left(\frac{3}{4} \right)^2 + \left(\frac{1}{4} \right)^2 \right] = 0.375$$

$$\text{Gini (Mild)} = 1 - \left[\left(\frac{2}{4} \right)^2 + \left(\frac{2}{4} \right)^2 \right] = 0.5$$

$$\text{Gini (Cool)} = 1 - \left(\frac{1}{1} \right)^2 = 0$$

$$\text{Gini (Cold)} = 1 - \left(\frac{1}{1} \right)^2 = 0$$

Gini children (Temperature)

$$= \left(0.375 * \frac{4}{10} \right) + \left(0.5 * \frac{4}{10} \right) + \left(0 * \frac{1}{10} \right) + \left(0 * \frac{1}{10} \right)$$

$$= \cancel{0.375} \quad 0.35$$

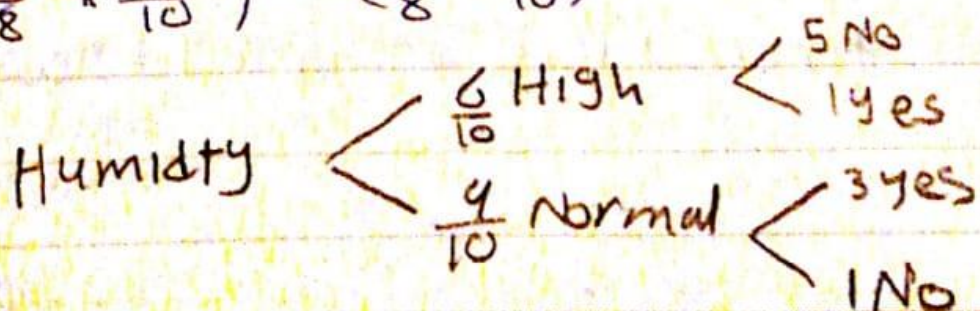
- we calculate Gini of Humidity

$$\text{Gini (High)} = 1 - \left[\left(\frac{5}{6} \right)^2 + \left(\frac{1}{6} \right)^2 \right] = \frac{5}{18}$$

$$\text{Gini (Normal)} = 1 - \left[\left(\frac{3}{4} \right)^2 + \left(\frac{1}{4} \right)^2 \right] = \frac{3}{8}$$

Gini children (Humidity)

$$= \left(\frac{5}{18} * \frac{6}{10} \right) + \left(\frac{3}{8} * \frac{4}{10} \right) = 0.316$$



- we calculate "wind"

```

      Wind
     /  \
  7 strong 5 No
   /  \   /  \
  3 weak 2 yes 2 yes 1 No

```

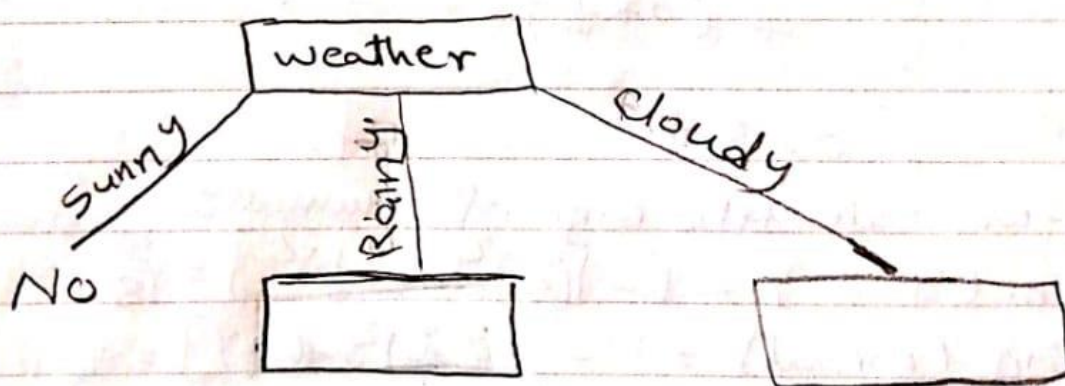
$$\text{Gini (strong)} = 1 - \left[\left(\frac{5}{7} \right)^2 + \left(\frac{2}{7} \right)^2 \right] = \frac{20}{49}$$

$$\text{Gini (weak)} = 1 - \left[\left(\frac{2}{3} \right)^2 + \left(\frac{1}{3} \right)^2 \right] = \frac{4}{9}$$

$$\begin{aligned} \text{Gini Children (wind)} \\ &= \left(\frac{20}{49} * \frac{7}{10} \right) + \left(\frac{4}{9} * \frac{3}{10} \right) \\ &= 0.419 \end{aligned}$$

The minimum Gini Children of "weather"

So the Root is "weather"



We must repeat this again to get new node

When weather = Cloudy

weather	Temperature	Humidity	wind	Hiking
Cloudy	Hot	High	Strong	No
Cloudy	mild	Normal	Strong	Yes
Cloudy	mild	High	Weak	Yes

Gini of "Temperature"

$$\text{Gini (Hot)} = 1 - \left[\left(\frac{1}{3}\right)^2\right] = 0$$

$$\text{Gini (mild)} = 1 - \left[\left(\frac{2}{2}\right)^2\right] = 0$$

$$\text{Gini Children (Temperature)} = \left(0 \times \frac{1}{3}\right) + \left(0 \times \frac{2}{3}\right) = 0$$

Gini of "Humidity"

$$\text{Gini (High)} = 1 - \left[\left(\frac{1}{2}\right)^2 + \left(\frac{1}{2}\right)^2\right] = 0.5$$

$$\text{Gini (Normal)} = 1 - \left[\left(1\right)^2\right] = 0$$

$$\text{Gini Children (Humidity)} = \left(0.5 \times \frac{2}{3}\right) + \left(0 \times \frac{1}{3}\right) = 0.33$$

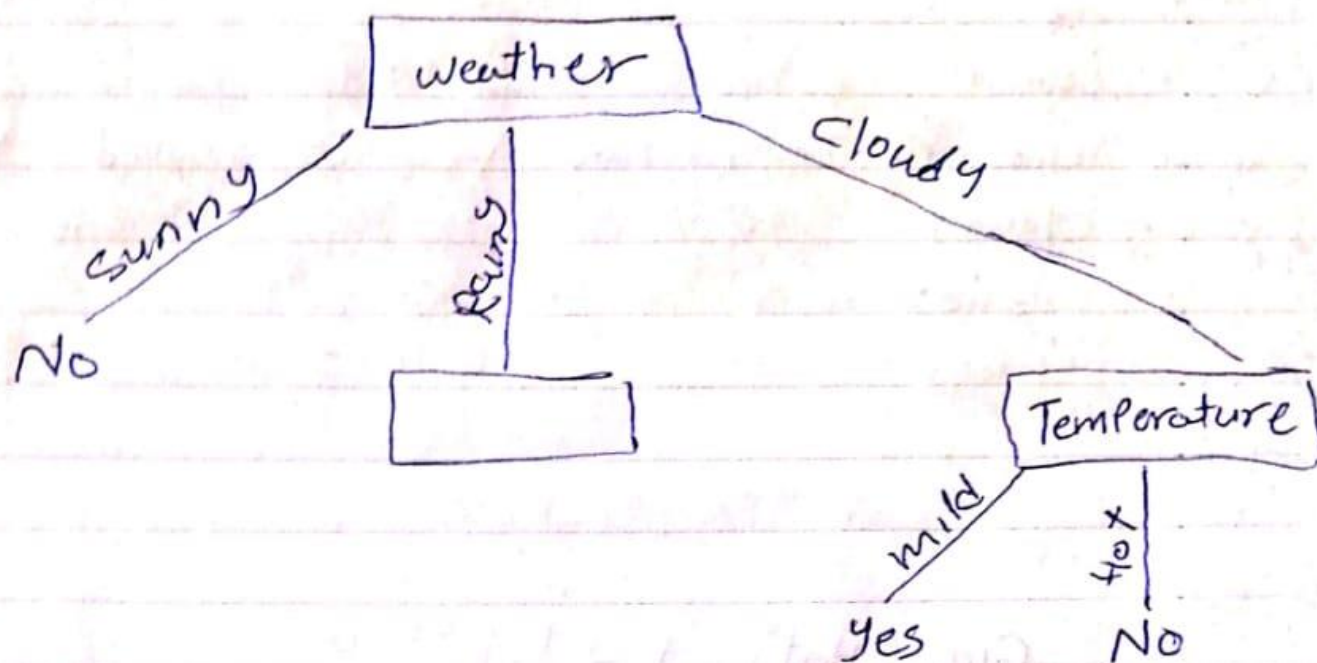
Gini of "Wind"

$$\text{Gini (Strong)} = 1 - \left[\left(\frac{2}{2}\right)^2 + \left(\frac{1}{2}\right)^2\right] = 0.5$$

$$\text{Gini (Weak)} = 1 - \left[\left(1\right)^2\right] = 0$$

$$\text{Gini Children (Wind)} = \left(0.5 \times \frac{2}{3}\right) + \left(0 \times \frac{1}{3}\right) = 0.33$$

The minimum Gini Children of "Temperature"
So the new node is "Temperature"
When (weather = Cloudy)



When weather = Rainy, we will get Gini ("Temperature"), ("Humidity") and ("wind")
 First we write our new Data when weather = Rainy

weather	Temperature	Humidity	wind	Hiking
Rainy	cold	Normal	Strong	yes
Rainy	cool	Normal	Strong	No
Rainy	hot	Normal	Weak	yes

Gini of "Temperature"

$$\text{Gini (Cold)} = 1 - (1)^2 = 0$$

$$\text{Gini (cool)} = 1 - (1)^2 = 0$$

$$\text{Gini (Hot)} = 1 - (1)^2 = 0$$

$$\text{Gini children (Temperature)} = (0 * \frac{1}{3}) + (0 * \frac{1}{3}) + (0 * \frac{1}{3}) = 0$$

Gini of "Humidity"

$$\text{Gini (Normal)} = 1 - \left(\left(\frac{2}{3} \right)^2 + \left(\frac{1}{3} \right)^2 \right) \\ = \frac{4}{9}$$

$$\text{Gini Children (Humidity)} = \left(\frac{4}{9} \times \frac{3}{3} \right) = \frac{4}{9}$$

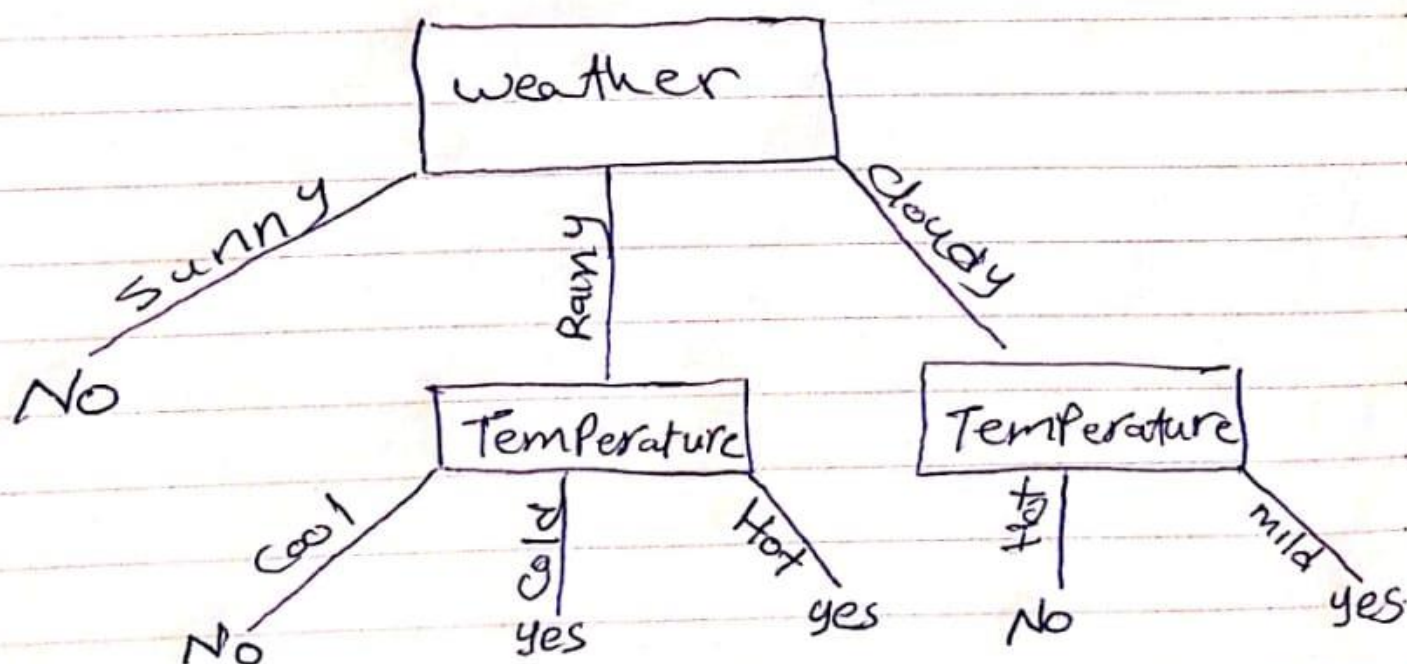
Gini of "wind"

$$\text{Gini (strong)} = 1 - \left(\left(\frac{1}{2} \right)^2 + \left(\frac{1}{2} \right)^2 \right) = 0.5$$

$$\text{Gini (weak)} = 1 - (1)^2 = 0$$

$$\text{Gini children (wind)} = \left(0.5 \times \frac{2}{3} \right) + \left(0 \times \frac{1}{3} \right) \\ = 0.33$$

The minimum Gini Children is
Gini of "Temperature" when
(weather = "Rainy")



This is the decision tree when we use
Gini method

(b) Please build a decision tree by using Information Gain (i.e., $IG(T, a) = \text{Entropy}(T) - \text{Entropy}(T|a)$, More information about IG).

b)

$$\begin{aligned}\text{Entropy}(S) &= -P_{\text{yes}} \log_2 P_{\text{yes}} - P_{\text{no}} \log_2 P_{\text{no}} \\ &= -\frac{4}{10} \log_2 \frac{4}{10} - \frac{6}{10} \log_2 \frac{6}{10} \\ &= 0.971\end{aligned}$$

$$\begin{aligned}\text{GAIN}(S, \text{Temperature}) &= \text{Entropy}(S) - P_{\text{hot}} [\text{Entropy}(\text{hot})] \\ &\quad - P_{\text{mild}} [\text{Entropy}(\text{mild})] \\ &\quad - P_{\text{cool}} [\text{Entropy}(\text{cool})] \\ &\quad - P_{\text{cold}} [\text{Entropy}(\text{cold})] \\ &= 0.971 - \frac{4}{10} \left[-\frac{2}{4} \log_2 \frac{2}{4} - \frac{2}{4} \log_2 \frac{2}{4} \right] \\ &\quad - \frac{4}{10} \left[-\frac{3}{4} \log_2 \frac{3}{4} - \frac{1}{4} \log_2 \frac{1}{4} \right] \\ &\quad - \frac{1}{10} [-\log_2 1] - \frac{1}{10} [-\log_2 1] \\ &= 0.247\end{aligned}$$

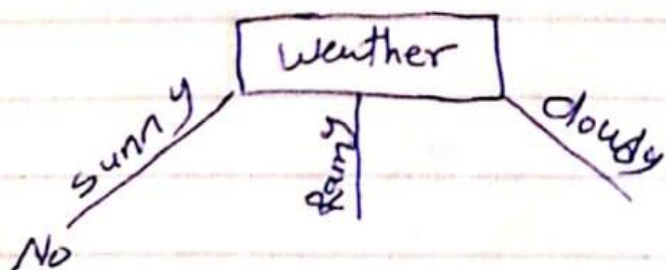
$$\begin{aligned}\text{GAIN}(S, \text{Weather}) &= \text{Entropy}(S) - P_{\text{cloudy}} [-\text{Entropy}(\text{cloudy})] \\ &\quad - P_{\text{rainy}} [-\text{Entropy}(\text{rainy})] \\ &\quad - P_{\text{sunny}} [-\text{Entropy}(\text{sunny})] \\ &= 0.971 - \frac{3}{10} \left[-\frac{1}{3} \log_2 \frac{1}{3} - \frac{2}{3} \log_2 \frac{2}{3} \right] \\ &\quad - \frac{4}{10} \left[-\frac{4}{4} \log_2 \frac{4}{4} \right] \\ &\quad - \frac{3}{10} \left[-\frac{2}{3} \log_2 \frac{2}{3} - \frac{1}{3} \log_2 \frac{1}{3} \right]\end{aligned}$$

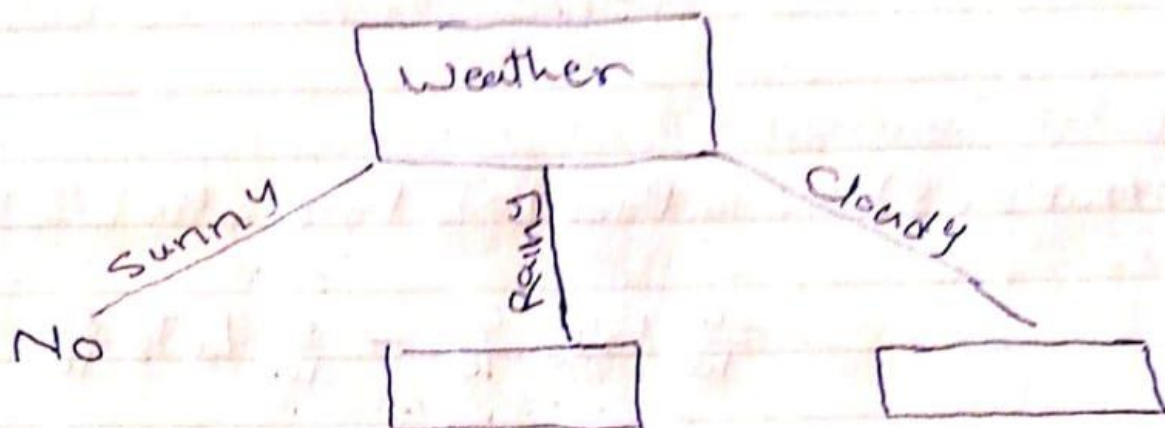
$$\text{Gain}(S, \text{weather}) = 0.42$$

$$\begin{aligned} \text{Gain}(S, \text{Humidity}) &= [-P_{\text{high}} [\text{Entropy high}] - P_{\text{normal}} [\text{Entropy normal}]] \\ &\quad + \text{Entropy}(S) \\ &= 0.971 - \frac{6}{10} \left[-\frac{5}{8} \log_2 \frac{5}{8} - \frac{1}{8} \log_2 \frac{1}{8} \right] \\ &\quad - \frac{4}{10} \left[-\frac{3}{4} \log_2 \frac{3}{4} - \frac{1}{4} \log_2 \frac{1}{4} \right] \\ &= 0.257 \end{aligned}$$

$$\begin{aligned} \text{Gain}(S, \text{Wind}) &= \text{Entropy}(S) - P_{\text{strong}} [\text{Entropy strong}] \\ &\quad - P_{\text{weak}} [\text{Entropy weak}] \\ &= 0.971 - \frac{7}{10} \left[-\frac{2}{7} \log_2 \frac{2}{7} - \frac{5}{7} \log_2 \frac{5}{7} \right] \\ &\quad - \frac{3}{10} \left[-\frac{1}{3} \log_2 \frac{1}{3} - \frac{2}{3} \log_2 \frac{2}{3} \right] \\ &= 0.0913 \end{aligned}$$

The maximum information Gain is the Root
 So the the Root is
 "weather"





our new data when (weather = cloudy)

weather	Temperature	Humidity	wind	Hiking
cloudy	Hot	High	Strong	No
cloudy	mild	Normal	strong	yes
Cloudy	mild	High	weak	yes

Entropy (S, weather = cloudy)

$$\begin{aligned}
 &= -P_{\text{yes}} \log_2 P_{\text{yes}} - P_{\text{no}} \log_2 P_{\text{no}} \\
 &= -\frac{2}{3} \log_2 \frac{2}{3} - \frac{1}{3} \log_2 \frac{1}{3} \\
 &= 0.918
 \end{aligned}$$

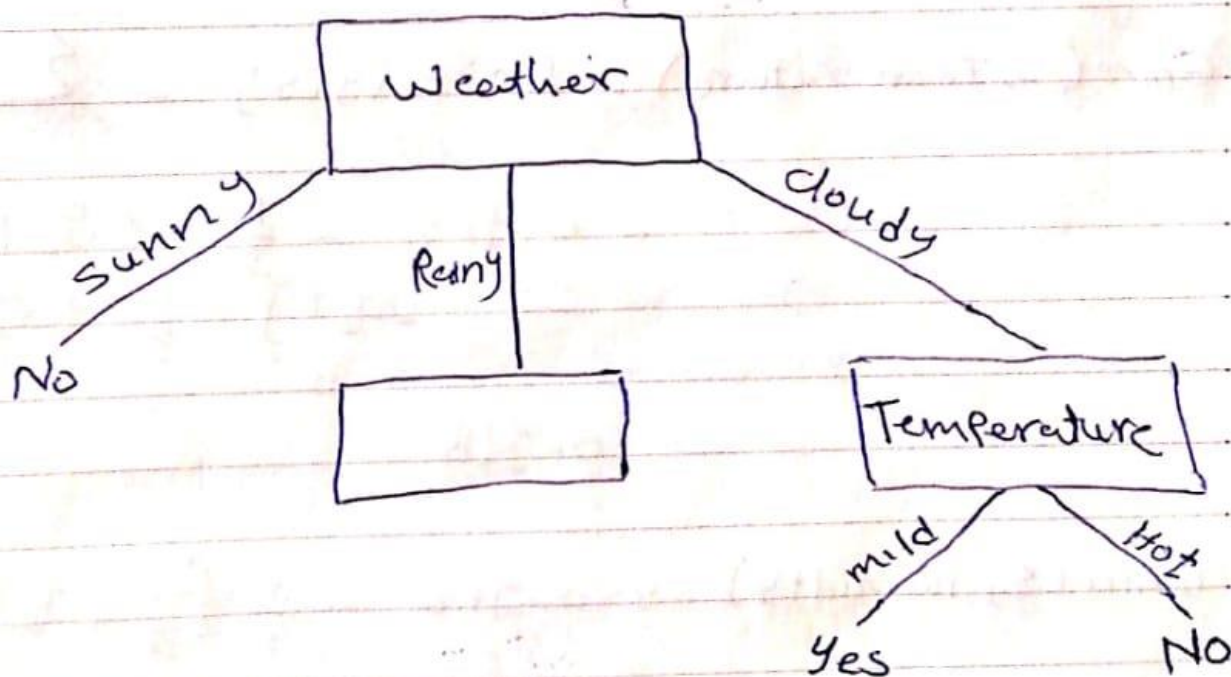
$$\begin{aligned}
 \text{Gain (S, Temperature)} &= \text{Entropy (S)} - \sum \frac{n_i}{n} \text{Entropy}(S_i) \\
 &= 0.918 - \frac{2}{3} \left[-\frac{2}{2} \log_2 \frac{2}{2} \right] \\
 &\quad - \frac{1}{3} \left[1 \log_2 1 \right] \\
 &= 0.918
 \end{aligned}$$

$$\begin{aligned}
 \text{Gain (S, Humidity)} &= \text{Entropy (S)} - \sum \frac{n_i}{n} \text{Entropy}(S_i) \\
 &= 0.918 - \frac{2}{3} \left[-\frac{1}{2} \log_2 \frac{1}{2} - \frac{1}{2} \log_2 \frac{1}{2} \right] \\
 &\quad - \frac{1}{3} \left[-1 \log_2 1 \right]
 \end{aligned}$$

$$\text{GAIN}(S, \text{Humidity}) = 0.251$$

$$\begin{aligned} \text{GAIN}(S, \text{wind}) &= \text{Entropy}(S) - \sum \frac{n_i}{n} \text{Entropy}(L_i) \\ &= 0.918 - \frac{2}{3} \left[-\frac{1}{2} \log_2 \frac{1}{2} - \frac{1}{2} \log_2 \frac{1}{2} \right] \\ &\quad - \frac{1}{3} [-\log_2 1] \\ &= 0.251 \end{aligned}$$

The new node For (weather = Cloudy)
is "Temperature"



our new Data (Weather = Rainy)

Weather	Temperature	Humidity	wind	Hiking
Rainy	Cold	Normal	strong	yes
Rainy	Cool	Normal	Strong	No
Rainy	Hot	Normal	weak	yes

(Weather = Rainy)

$$\text{Entropy}(S) = -P_{\text{yes}} \log_2 P_{\text{yes}} - P_{\text{no}} \log_2 P_{\text{no}}$$

$$= -\frac{2}{3} \log_2 \frac{2}{3} - \frac{1}{3} \log_2 \frac{1}{3}$$

$$= 0.918$$

$$\text{Gain}(S, \text{Temperature}) = \text{Entropy}(S) - \sum \frac{n_i}{n} \text{Entropy}(U_i)$$

$$= 0.918 - \frac{1}{3} [-\log_2 1] - \frac{1}{3} [-\log_2 1] - \frac{1}{3} [-\log_2 1]$$

$$= 0.918$$

$$\text{Gain}(S, \text{Humidity}) = 0.918 - \frac{3}{3} \left[-\frac{2}{3} \log_2 \frac{2}{3} - \frac{1}{3} \log_2 \frac{1}{3} \right]$$

$$= -0.0003$$

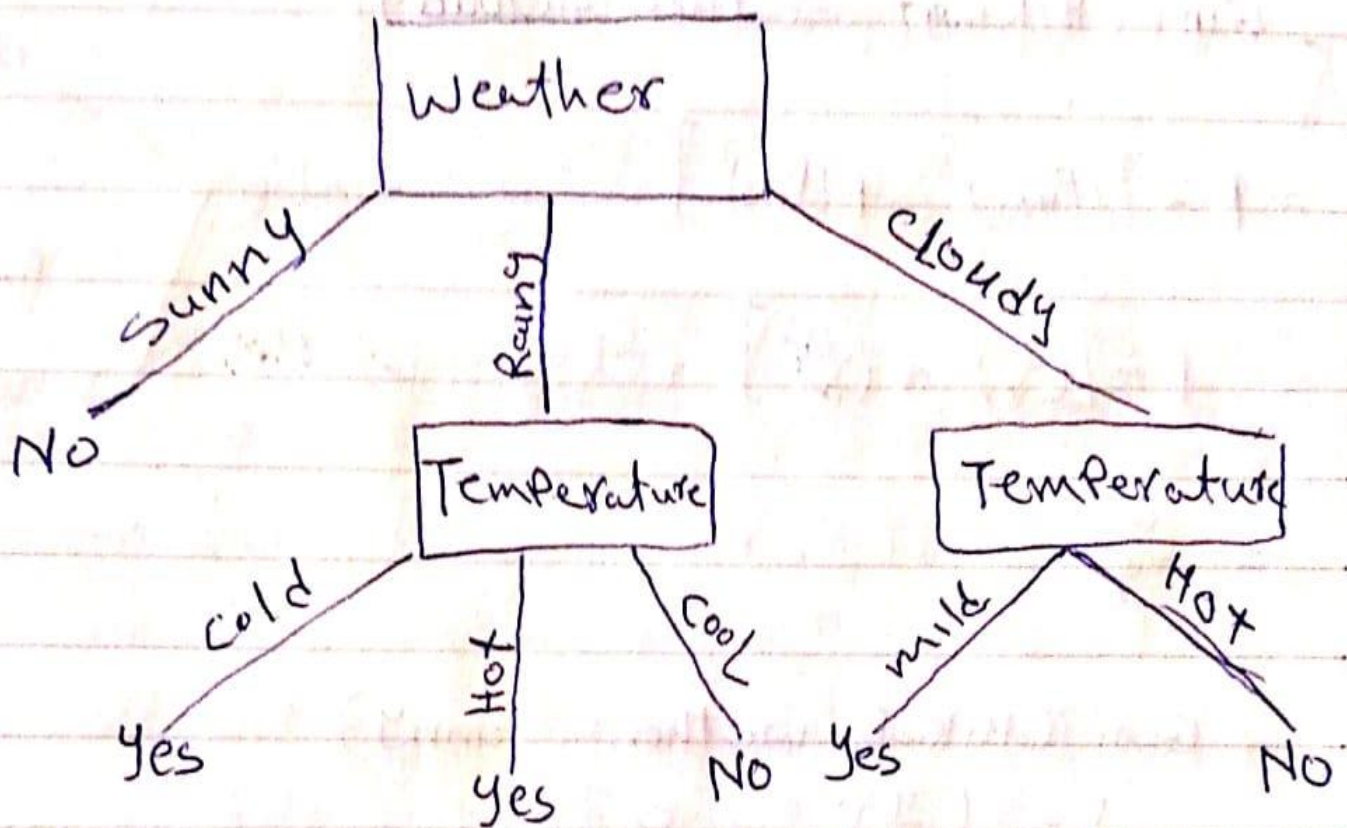
$$\text{Gain}(S, \text{wind}) = \text{Entropy}(S) - \sum \frac{n_i}{n} \text{Entropy}(U_i)$$

$$= 0.918 - \frac{2}{3} \left[-\frac{1}{2} \log_2 \frac{1}{2} - \frac{1}{2} \log_2 \frac{1}{2} \right] - \frac{1}{3} [-\log_2 1]$$

$$\text{Gain}(s, \text{wind}) = 0.251$$

The new node for (weather = Rainy)

is "Temperature"



This is the decision tree when we calculate by information Gain

(c) Comparison of Gini Index and Information Gain:

Both Gini index and information gain are commonly used metrics for building decision trees.

Advantages of Gini Index:

- Gini index is computationally faster than information gain, especially when dealing with large datasets.
- Gini index is less prone to overfitting than information gain, since it does not rely on the number of instances in each class.

Disadvantages of Gini Index:

- Gini index is not as informative as information gain, since it only considers the impurity of the current node and not the potential future nodes.
- Gini index is biased towards features with many categories, since they tend to have lower Gini index values.

Advantages of Information Gain:

- Information gain provides more information about the relationship between the features and the label, since it considers the entropy of both the current node and the potential future nodes.
- Information gain is less biased towards features with many categories, since it takes into account the number of instances in each category.

Disadvantages of Information Gain:

- Information gain is computationally more expensive than Gini index, especially when dealing with large datasets.
- Information gain is more prone to overfitting than Gini index, especially when dealing with noisy or irrelevant features.

Part 2: Programming Questions

(a)

View the dataset:

```
# Load dataset
dataset = pd.read_csv('/content/KDD.csv')
# print shape of dataset
print("Shape of dataset", dataset.shape)
# show the first five rows
dataset
```

Shape of dataset (494021, 39)

	duration	src_bytes	dst_bytes	land	wrong_fragment	urgent	hot	num_failed_logins	logged_in	num_compromised	...	dst_host_srv_count	dst_host_same_srv_rate	dst_host_diff_srv_rate
0	0	181	5450	0	0	0	0	0	1	0	...	9	1.0	0.0
1	0	239	486	0	0	0	0	0	1	0	...	19	1.0	0.0
2	0	235	1337	0	0	0	0	0	1	0	...	29	1.0	0.0
3	0	219	1337	0	0	0	0	0	1	0	...	39	1.0	0.0
4	0	217	2032	0	0	0	0	0	1	0	...	49	1.0	0.0
...
494016	0	310	1881	0	0	0	0	0	1	0	...	255	1.0	0.0
494017	0	282	2286	0	0	0	0	0	1	0	...	255	1.0	0.0
494018	0	203	1200	0	0	0	0	0	1	0	...	255	1.0	0.0
494019	0	291	1200	0	0	0	0	0	1	0	...	255	1.0	0.0
494020	0	219	1234	0	0	0	0	0	1	0	...	255	1.0	0.0

494021 rows x 39 columns

Function to preprocess the dataset:

```
# function to process dataset
def preprocess_data(df):
    # Drop rows with missing values as ValueError: Input X contains NaN.
    #df = df.dropna()
    # Separate the input features (X) and target variable (Y)
    # Select all columns except the last one
    X = df.iloc[:, :-1]
    # Select the last column
    Y = df.iloc[:, -1]
    # view 38 input feature variables and 1 target
    print("The feature variables:",X.shape[1])
    print("The shape of target:",Y.shape)
    # Normalize the input features using MinMaxScaler
    scaler = MinMaxScaler()
    X_normalized = scaler.fit_transform(X)
    # Perform filter-based feature selection
    # Select top 9 features
    selector = SelectKBest(mutual_info_classif, k=9)
    X_selected = selector.fit_transform(X_normalized, Y)
    # Create a new DataFrame with selected features
    selected_features = selector.get_support(indices=True)
    selected_columns = df.columns[selected_features].tolist()
    my_data = pd.DataFrame(X_selected, columns=selected_columns)
    # Add the target column
    my_data['target'] = Y
    return my_data, X_selected, Y, selected_columns
# Preprocess data
my_data, X_selected, Y,selected_columns = preprocess_data(dataset)
```

Display shape of train dataset (features) and tests dataset (target):

```
The feature variables: 38
The shape of target: (494021,)
```

Display first 5 rows after preprocessing, normalization & feature selection:

	src_bytes	dst_bytes	logged_in	count	srv_count	dst_host_count	dst_host_srv_count	dst_host_same_src_port_rate	dst_host_srv_diff_host_rate	target
0	2.610418e-07	0.001057	1.0	0.015656	0.015656	0.035294	0.035294	0.11	0.0	0
1	3.446905e-07	0.000094	1.0	0.015656	0.015656	0.074510	0.074510	0.05	0.0	0
2	3.389216e-07	0.000259	1.0	0.015656	0.015656	0.113725	0.113725	0.03	0.0	0
3	3.158461e-07	0.000259	1.0	0.011742	0.011742	0.152941	0.152941	0.03	0.0	0
4	3.129617e-07	0.000394	1.0	0.011742	0.011742	0.192157	0.192157	0.02	0.0	0

(b)

Function to train DT and compute performance for 3 different subsets:

```
# function to compute performance of DT in terms of Classification report using 3 subsets
def compute_DT_performance():
    # Define the train and test ratios
    train_ratios = [0.7, 0.6, 0.5]
    test_ratios = [0.3, 0.4, 0.5]
    i=1
    for train_ratio, test_ratio in zip(train_ratios, test_ratios):
        # Split the data into training and testing subsets
        X_train, X_test, Y_train, Y_test = split_data(X_selected,Y,train_ratio, test_ratio)
        # Perform further operations with the subsets, such as training and evaluating the decision tree
        clf = DecisionTreeClassifier()
        clf.fit(X_train, Y_train)
        # Make predictions on the test set and evaluate performance
        Y_pred = clf.predict(X_test)
        report, accuracy, confusion = evaluate_classification(clf, X_test, Y_test)
        print(f"my_data_{i}")
        # Print the results
        print(f"Train Ratio: {train_ratio}, Test Ratio: {test_ratio}")
        print("Classification Report:")
        print(report)
        print()
        i+=1
    # call the compute_DT_performance
    compute_DT_performance()
```

Display classification report for 3 different subsets:

my_data_1

Train Ratio: 0.7, Test Ratio: 0.3

Classification Report:

	precision	recall	f1-score	support
0	1.00	1.00	1.00	29163
1	1.00	1.00	1.00	119044
accuracy			1.00	148207
macro avg	1.00	1.00	1.00	148207
weighted avg	1.00	1.00	1.00	148207

my_data_2

Train Ratio: 0.6, Test Ratio: 0.4

Classification Report:

	precision	recall	f1-score	support
0	1.00	1.00	1.00	38962
1	1.00	1.00	1.00	158647
accuracy			1.00	197609
macro avg	1.00	1.00	1.00	197609
weighted avg	1.00	1.00	1.00	197609

my_data_3

Train Ratio: 0.5, Test Ratio: 0.5

Classification Report:

	precision	recall	f1-score	support
0	1.00	1.00	1.00	48796
1	1.00	1.00	1.00	198215
accuracy			1.00	247011
macro avg	1.00	1.00	1.00	247011
weighted avg	1.00	1.00	1.00	247011

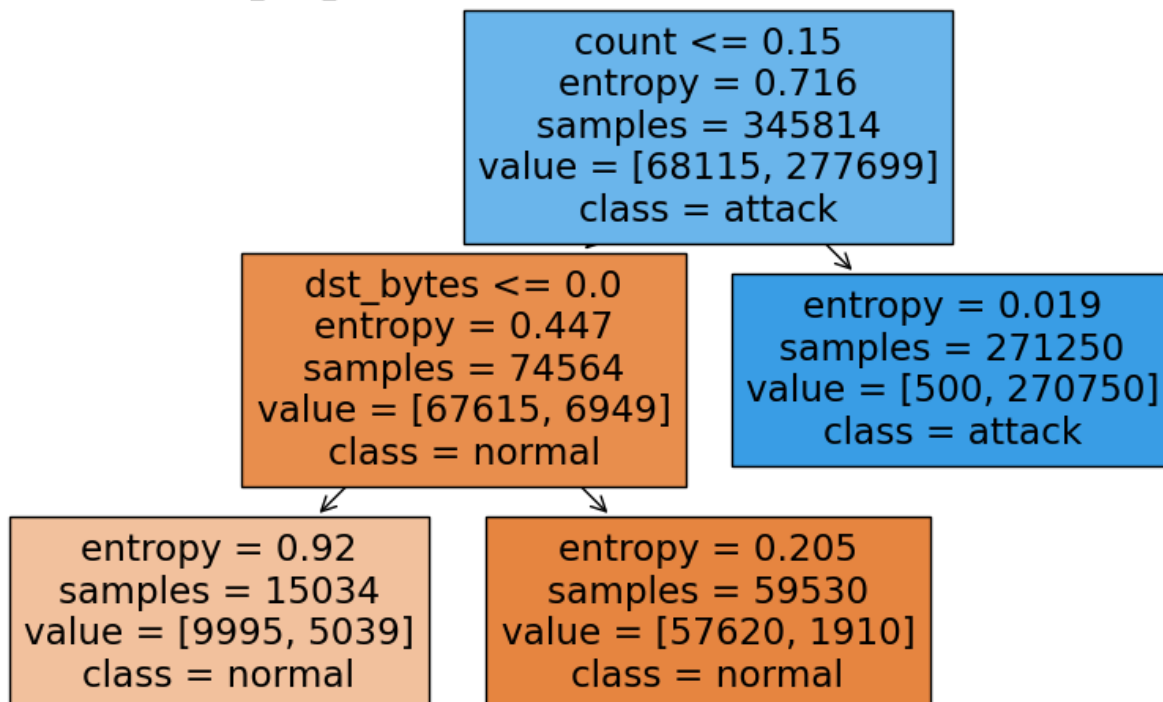
(c)

Function to visualize the best split of the Decision tree:

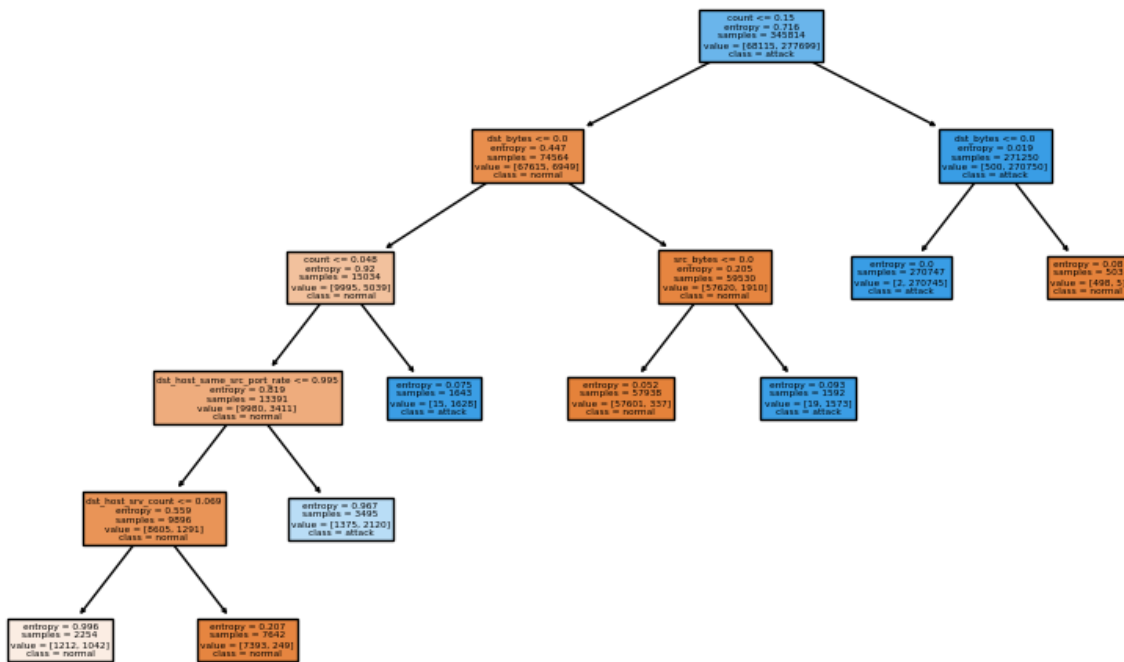
```
# function to train the decision tree model and visualize the best split of the Decision tree
def train_decision_tree():
    # Define the train and test ratios
    train_ratios = [0.7, 0.6, 0.5]
    test_ratios = [0.3, 0.4, 0.5]
    # Define max depths as assumed
    max_depths = [4, 6, 8]
    i = 0
    for max_depth in max_depths:
        for j in range(3):
            train_ratio = train_ratios[j]
            test_ratio = test_ratios[j]
            # Split the data into training and testing subsets
            X_train, X_test, Y_train, Y_test = split_data(X_selected, Y, train_ratio, test_ratio)
            # Initialize and train the Decision Tree classifier
            clf = DecisionTreeClassifier(max_depth=max_depth, criterion='entropy', max_leaf_nodes=ceil((2**max_depth)*0.1+1))
            clf.fit(X_train, Y_train)
            # Plot the decision tree
            plt.figure(figsize=(10, 6))
            plot_tree(clf, filled=True, feature_names=my_data.columns[:-1], class_names=['normal', 'attack'])
            n = j + 1
            plt.title(f"my_data_{n}, Max Depth: {max_depth}, Train Ratio: {train_ratio}, Test Ratio: {test_ratio}")
            plt.show()
            i += 1

train_decision_tree()
```

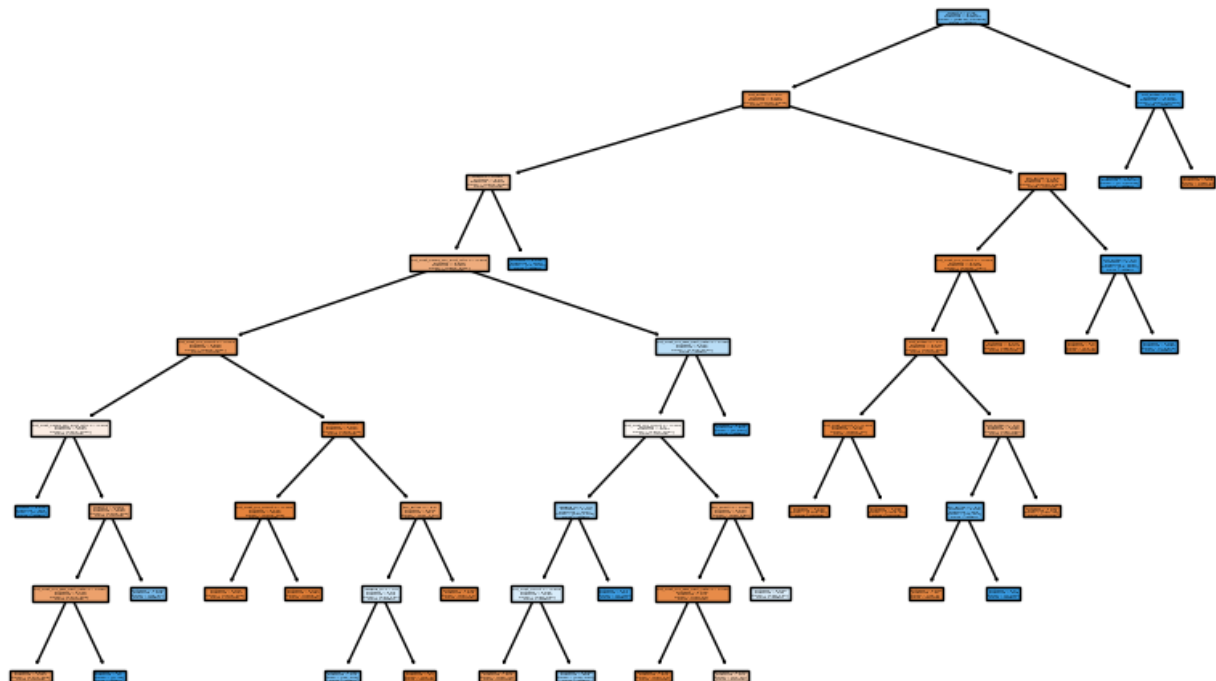
my_data_1, Max Depth: 4, Train Ratio: 0.7, Test Ratio: 0.3



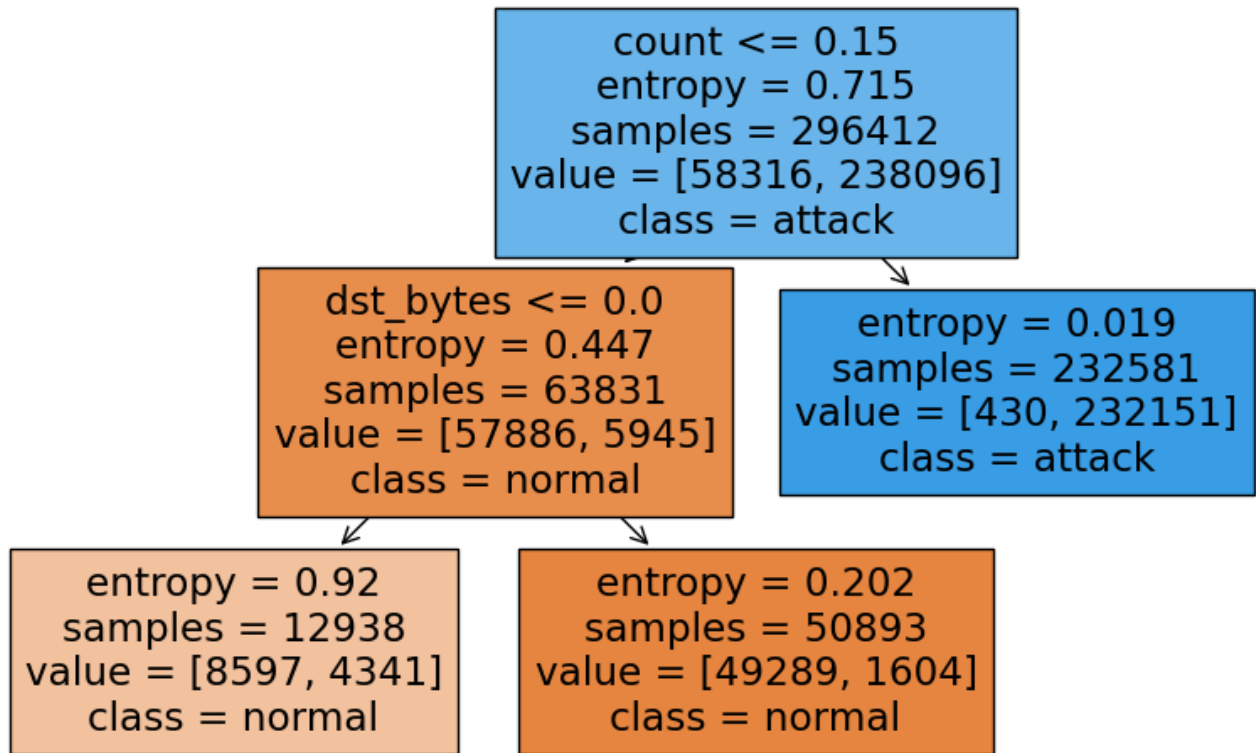
my_data_1, Max Depth: 6, Train Ratio: 0.7, Test Ratio: 0.3



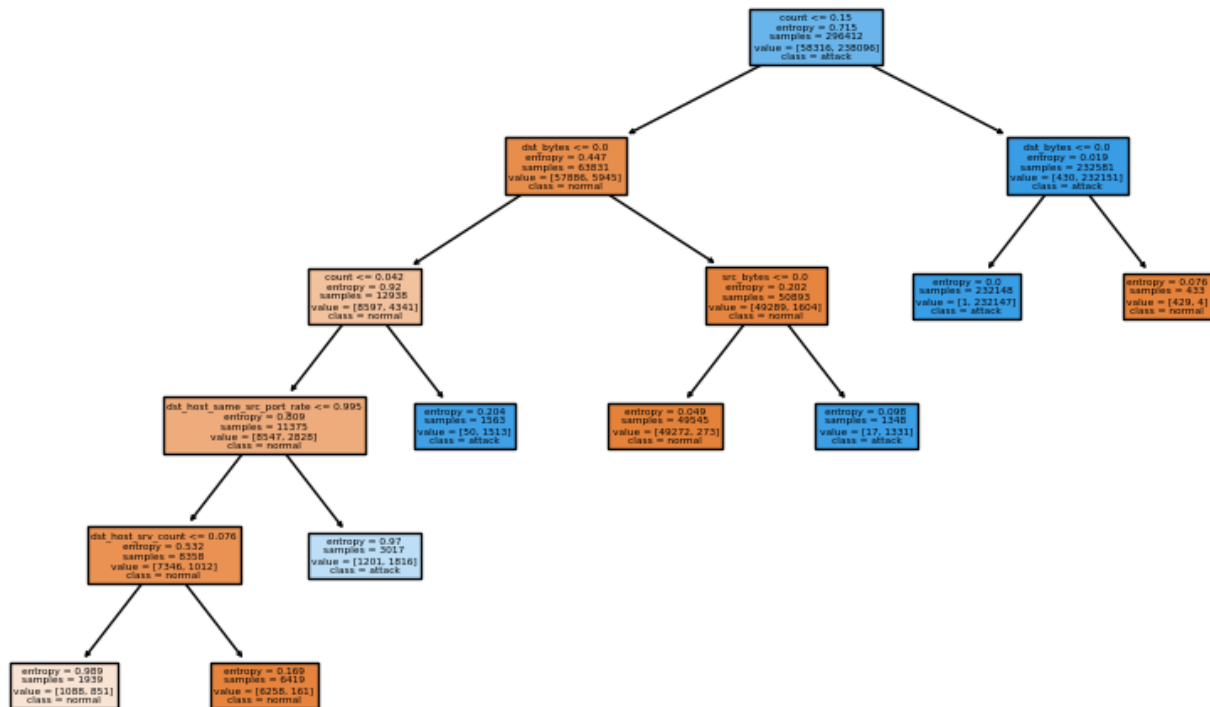
my_data_1, Max Depth: 8, Train Ratio: 0.7, Test Ratio: 0.3



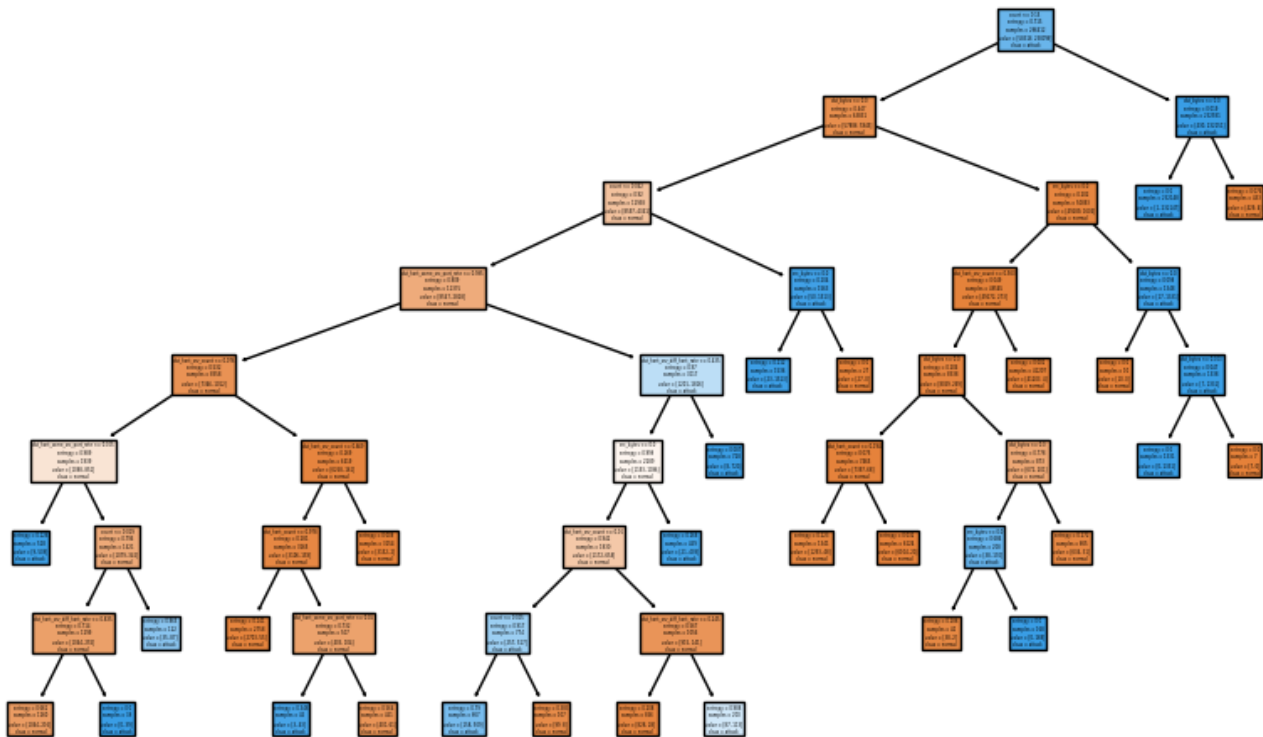
my_data_2, Max Depth: 4, Train Ratio: 0.6, Test Ratio: 0.4



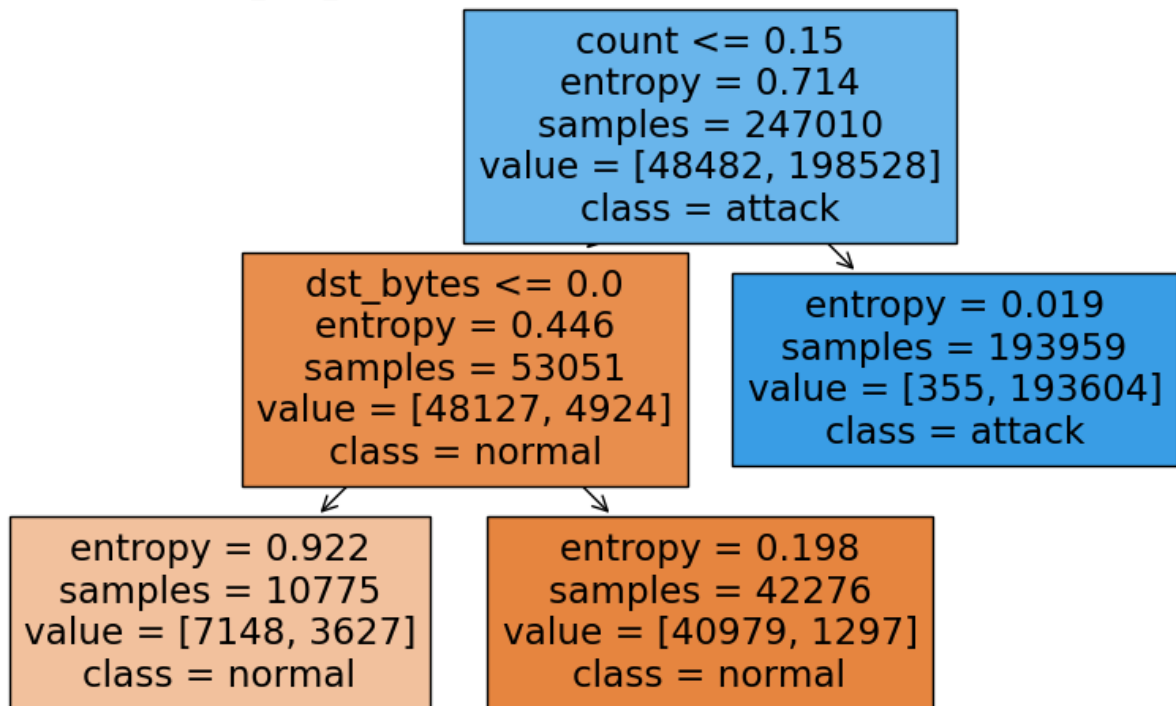
my_data_2, Max Depth: 6, Train Ratio: 0.6, Test Ratio: 0.4



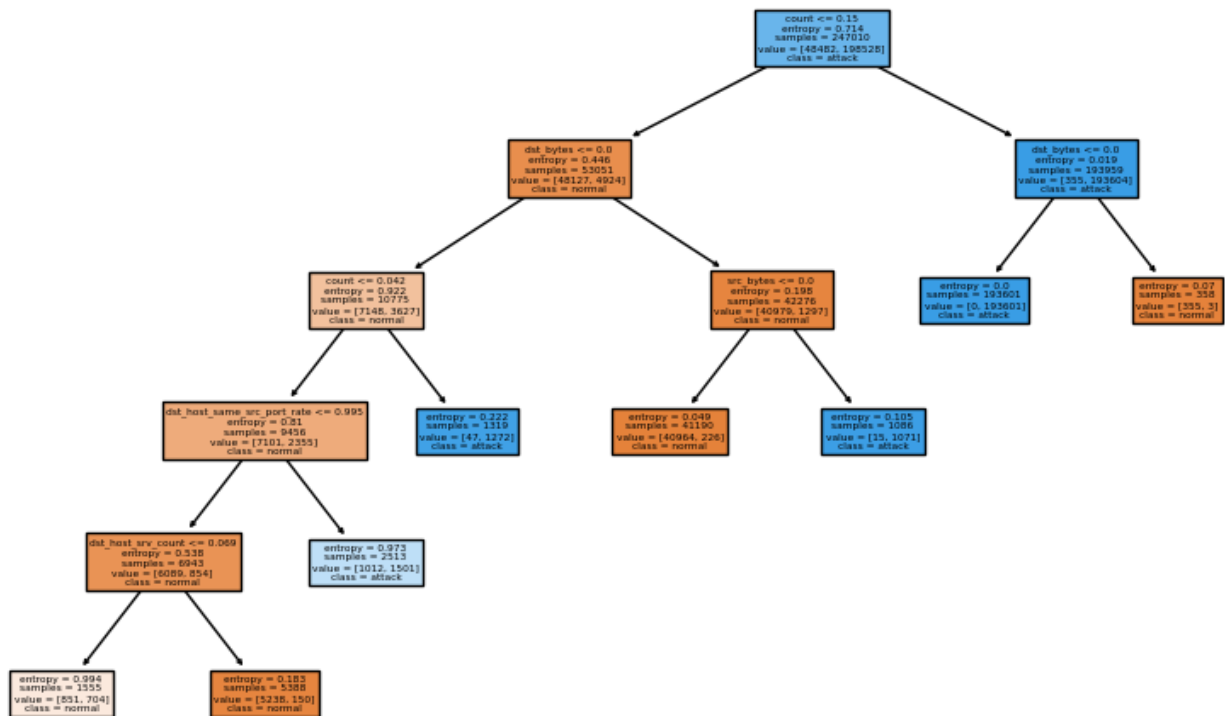
my_data_2, Max Depth: 8, Train Ratio: 0.6, Test Ratio: 0.4



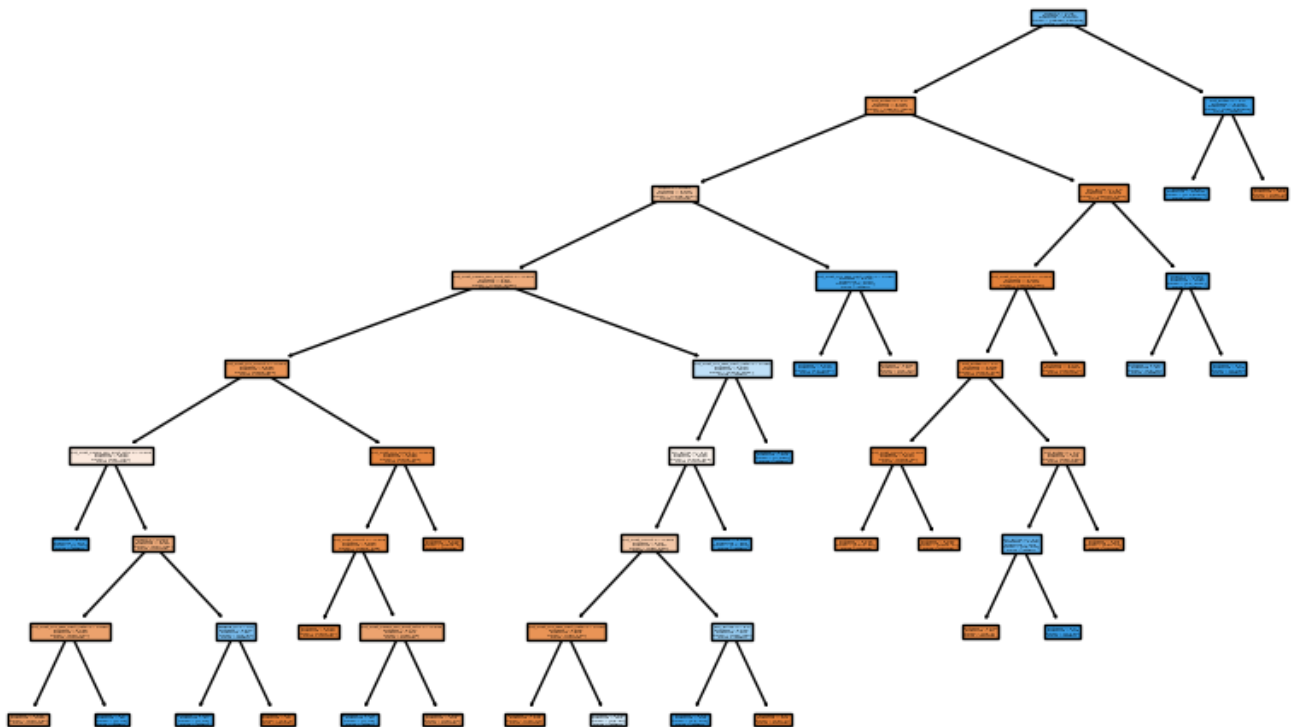
my_data_3, Max Depth: 4, Train Ratio: 0.5, Test Ratio: 0.5



my_data_3, Max Depth: 6, Train Ratio: 0.5, Test Ratio: 0.5



my_data_3, Max Depth: 8, Train Ratio: 0.5, Test Ratio: 0.5



(d)

Function to compute performance of tuned DT:

```
# function to compute and compare tuned DT
def compare_tuned_decision_tree_performance(train_ratio, test_ratio):
    X_train, X_test, Y_train, Y_test = split_data(X_selected, Y, train_ratio, test_ratio)
    max_depths = [4, 6, 8]
    data_labels = {0.3: "my_data_1", 0.4: "my_data_2", 0.5: "my_data_3"}
    data_label = data_labels[test_ratio]
    for max_depth in max_depths:
        # Initialize and train the Decision Tree classifier
        clf = DecisionTreeClassifier(max_depth=max_depth, criterion='entropy', max_leaf_nodes=ceil((2**max_depth)*0.1+1))
        clf.fit(X_train, Y_train)
        # Make predictions on the test set
        Y_pred = clf.predict(X_test)
        report, accuracy, confusion = evaluate_classification(clf, X_test, Y_test)
        # Generate classification report and confusion matrix
        print(f"{data_label}, Max Depth: {max_depth}, Test Ratio: {test_ratio}")
        print("Accuracy:", accuracy)
        print("Classification Report:")
        print(report)
        print(f"Confusion Matrix with: Max Depth: {max_depth}, {test_ratio} test data")
        print()
        plot_confusion_matrix(confusion)
        print()
```

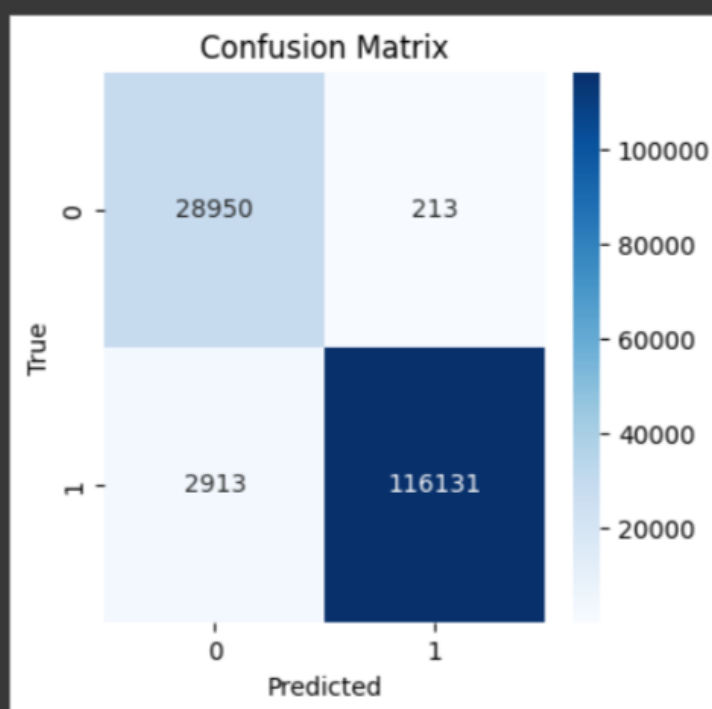
Display the accuracy scores, classification report, and confusion matrix of tuned classifier for each subset:

```
my_data_1, Max Depth: 4, Test Ratio: 0.3
Accuracy: 0.9789078788451288
Classification Report:
              precision    recall  f1-score   support

     0       0.91      0.99      0.95      29163
     1       1.00      0.98      0.99     119044

   accuracy          0.98      148207
  macro avg       0.95      0.98      0.97     148207
 weighted avg     0.98      0.98      0.98     148207
```

Confusion Matrix with: Max Depth: 4, 0.3 test data



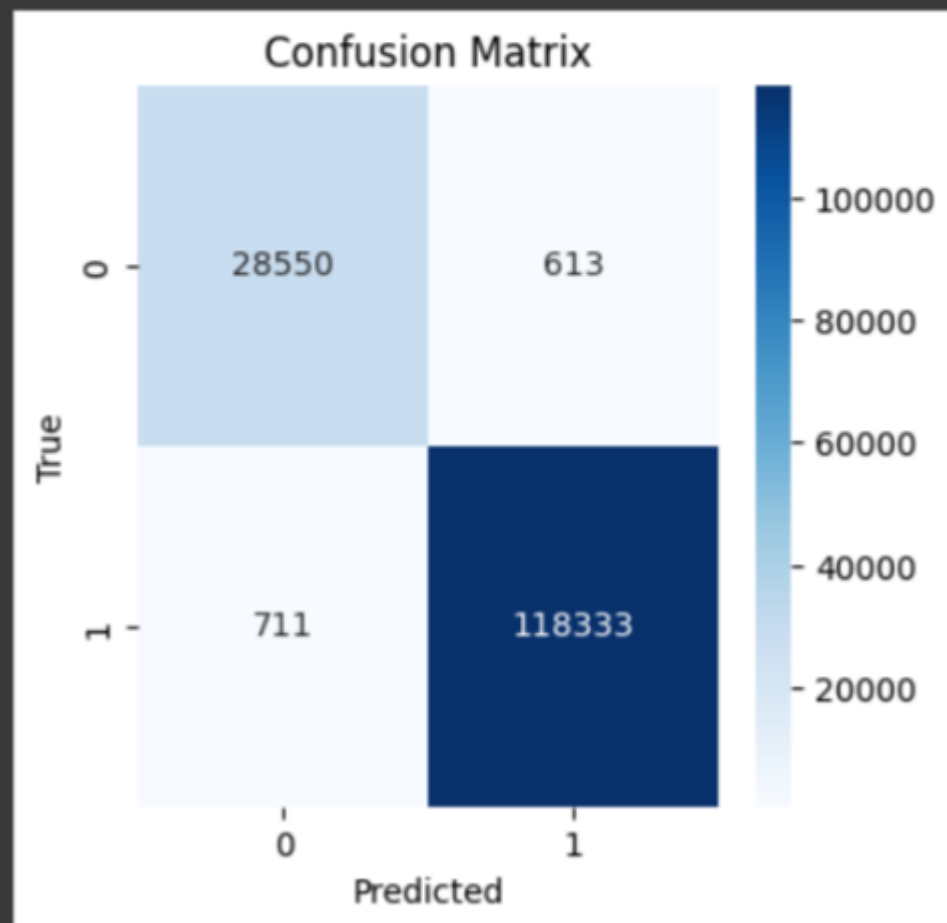
my_data_1, Max Depth: 6, Test Ratio: 0.3

Accuracy: 0.9910665488134839

Classification Report:

	precision	recall	f1-score	support
0	0.98	0.98	0.98	29163
1	0.99	0.99	0.99	119044
accuracy			0.99	148207
macro avg	0.99	0.99	0.99	148207
weighted avg	0.99	0.99	0.99	148207

Confusion Matrix with: Max Depth: 6, 0.3 test data



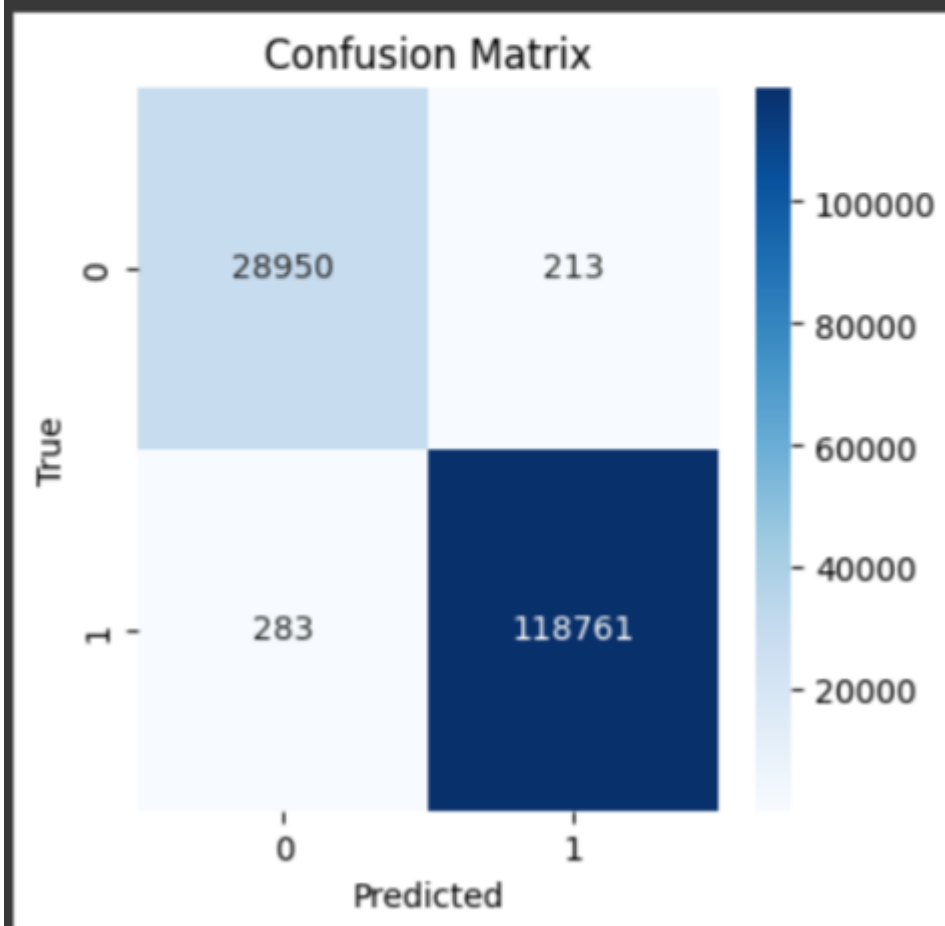
my_data_1, Max Depth: 8, Test Ratio: 0.3

Accuracy: 0.9966533294648701

Classification Report:

	precision	recall	f1-score	support
0	0.99	0.99	0.99	29163
1	1.00	1.00	1.00	119044
accuracy			1.00	148207
macro avg	0.99	1.00	0.99	148207
weighted avg	1.00	1.00	1.00	148207

Confusion Matrix with: Max Depth: 8, 0.3 test data



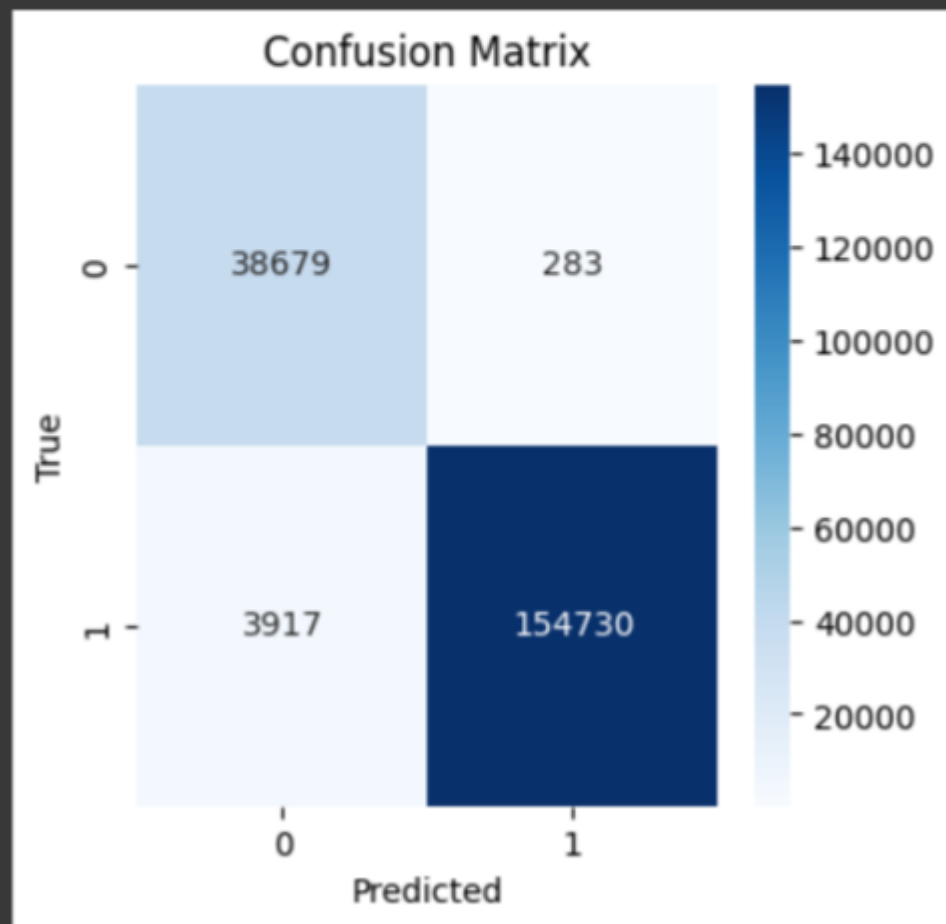
my_data_2, Max Depth: 4, Test Ratio: 0.4

Accuracy: 0.9787459073220349

Classification Report:

	precision	recall	f1-score	support
0	0.91	0.99	0.95	38962
1	1.00	0.98	0.99	158647
accuracy			0.98	197609
macro avg	0.95	0.98	0.97	197609
weighted avg	0.98	0.98	0.98	197609

Confusion Matrix with: Max Depth: 4, 0.4 test data



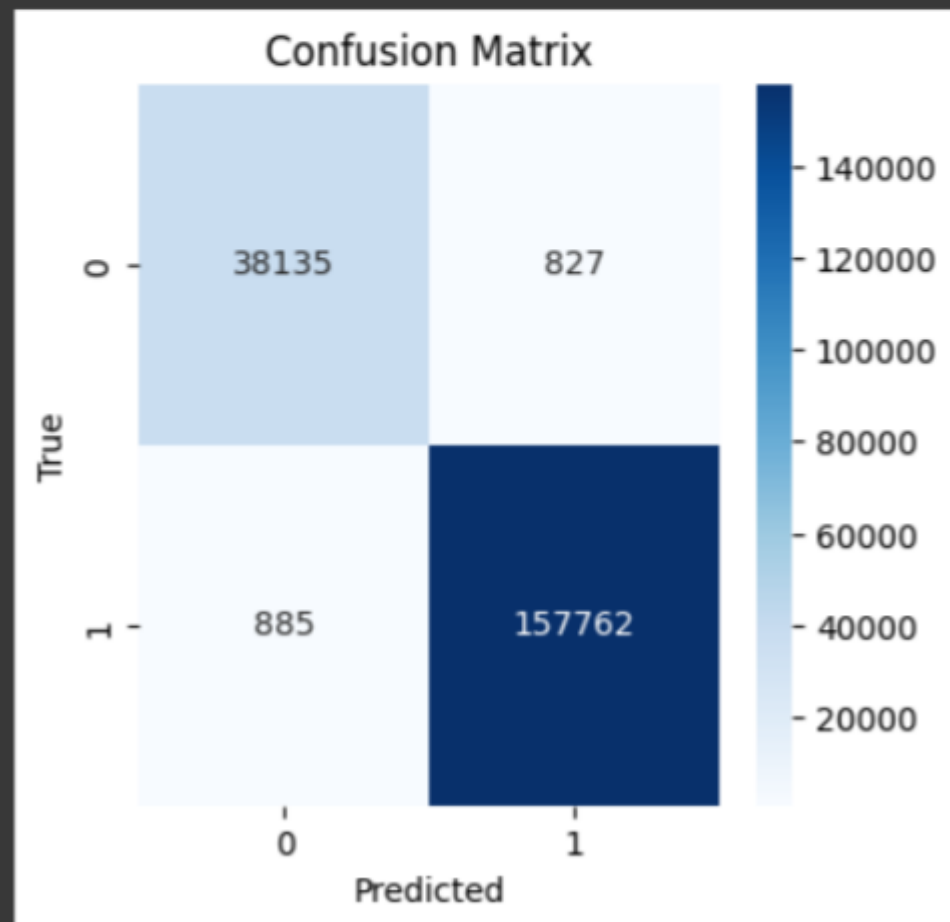
my_data_2, Max Depth: 6, Test Ratio: 0.4

Accuracy: 0.9913364269846009

Classification Report:

	precision	recall	f1-score	support
0	0.98	0.98	0.98	38962
1	0.99	0.99	0.99	158647
accuracy			0.99	197609
macro avg	0.99	0.99	0.99	197609
weighted avg	0.99	0.99	0.99	197609

Confusion Matrix with: Max Depth: 6, 0.4 test data



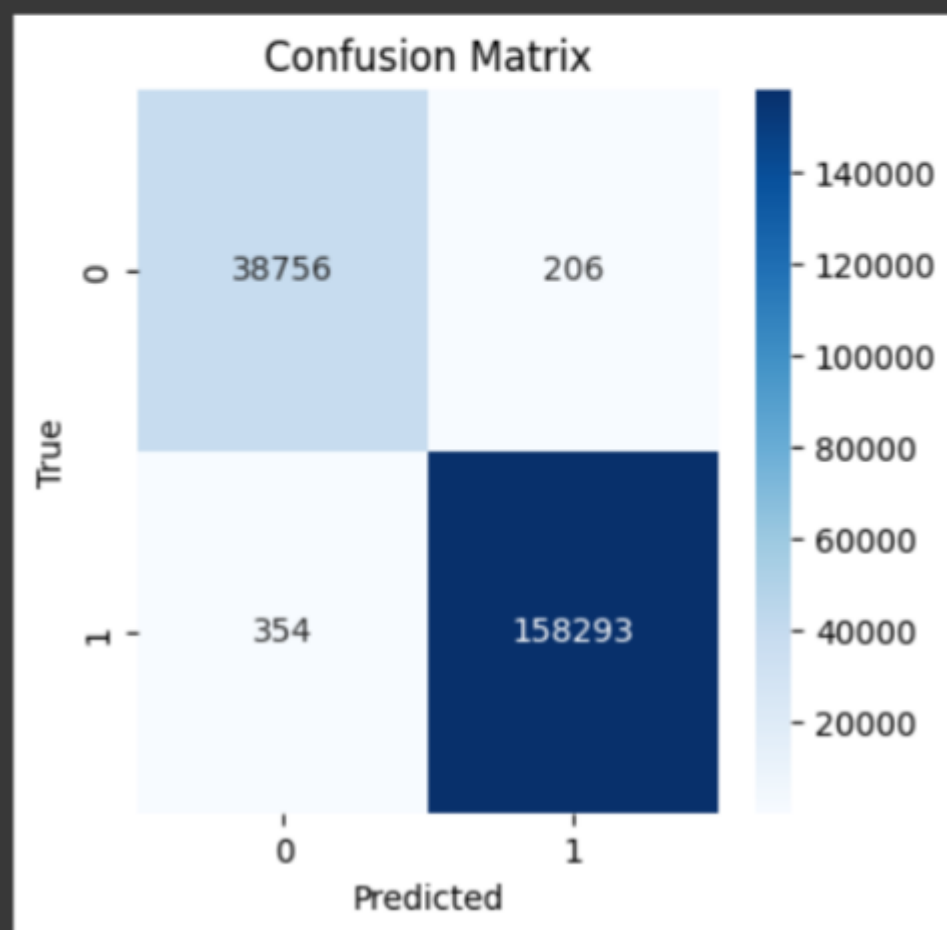
my_data_2, Max Depth: 8, Test Ratio: 0.4

Accuracy: 0.9971661209762713

Classification Report:

	precision	recall	f1-score	support
0	0.99	0.99	0.99	38962
1	1.00	1.00	1.00	158647
accuracy			1.00	197609
macro avg	0.99	1.00	1.00	197609
weighted avg	1.00	1.00	1.00	197609

Confusion Matrix with: Max Depth: 8, 0.4 test data



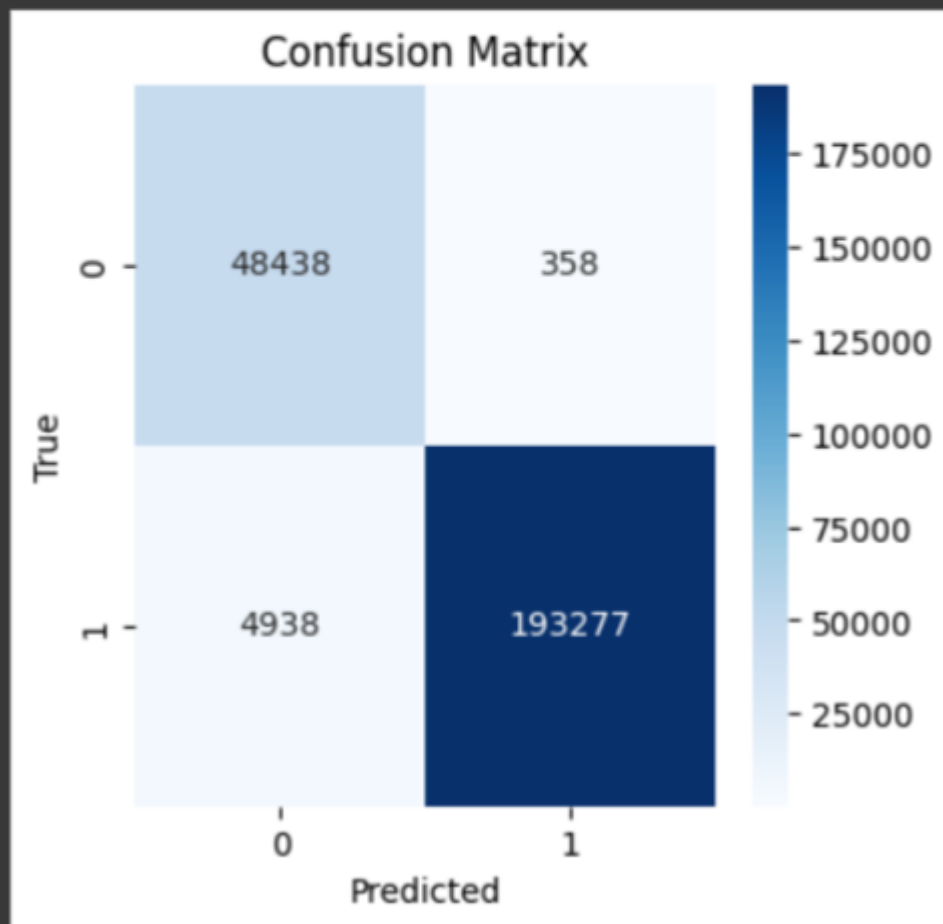
my_data_3, Max Depth: 4, Test Ratio: 0.5

Accuracy: 0.9785596592864285

Classification Report:

	precision	recall	f1-score	support
0	0.91	0.99	0.95	48796
1	1.00	0.98	0.99	198215
accuracy			0.98	247011
macro avg	0.95	0.98	0.97	247011
weighted avg	0.98	0.98	0.98	247011

Confusion Matrix with: Max Depth: 4, 0.5 test data



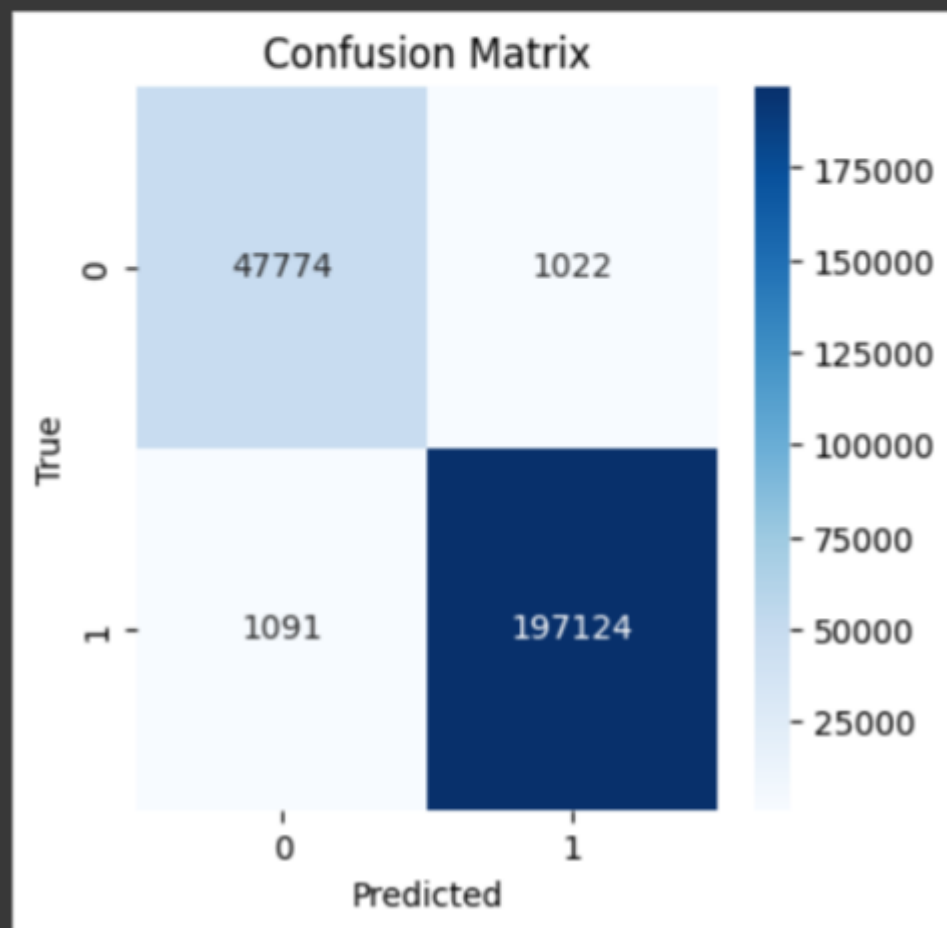
my_data_3, Max Depth: 6, Test Ratio: 0.5

Accuracy: 0.9914457250891661

Classification Report:

	precision	recall	f1-score	support
0	0.98	0.98	0.98	48796
1	0.99	0.99	0.99	198215
accuracy			0.99	247011
macro avg	0.99	0.99	0.99	247011
weighted avg	0.99	0.99	0.99	247011

Confusion Matrix with: Max Depth: 6, 0.5 test data



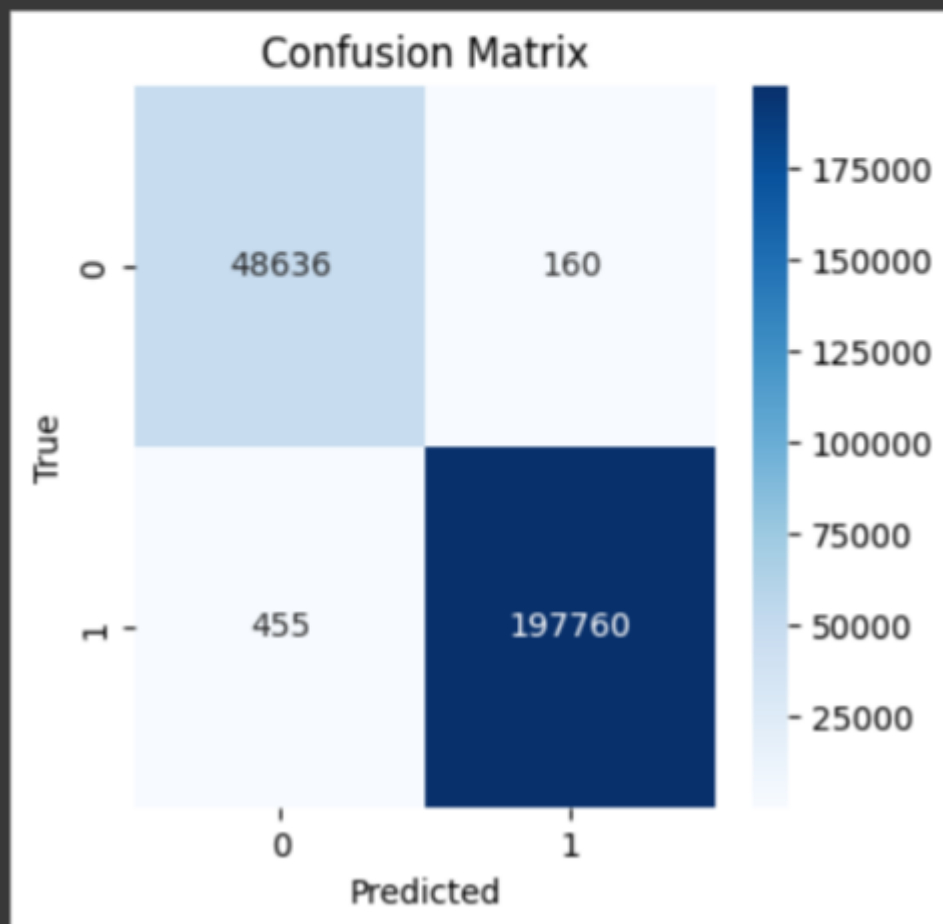
my_data_3, Max Depth: 8, Test Ratio: 0.5

Accuracy: 0.9975102323378311

Classification Report:

	precision	recall	f1-score	support
0	0.99	1.00	0.99	48796
1	1.00	1.00	1.00	198215
accuracy			1.00	247011
macro avg	0.99	1.00	1.00	247011
weighted avg	1.00	1.00	1.00	247011

Confusion Matrix with: Max Depth: 8, 0.5 test data



- (d) Before mitigation, I didn't find overfitting problem although I tried several parameters for Decision Tree algorithm. But I had to do the task as explained and perform the 3 mitigation strategies as asked.

