

```
import pandas as pd

# Path to your CSV file
csv_file_path = '/content/data (1).csv'

# Read the CSV file directly, trying a different encoding
try:
    df = pd.read_csv(csv_file_path, encoding='ISO-8859-1')
    # Display the data
    print(df)
except FileNotFoundError:
    print(f"Error: The file {csv_file_path} was not found.")
except Exception as e:
    print(f"An error occurred while reading the CSV file: {e}")
```

	InvoiceNo	StockCode	Description	Quantity	\
0	536365	85123A	WHITE HANGING HEART T-LIGHT HOLDER	6.0	
1	536365	71053	WHITE METAL LANTERN	6.0	
2	536365	84406B	CREAM CUPID HEARTS COAT HANGER	8.0	
3	536365	84029G	KNITTED UNION FLAG HOT WATER BOTTLE	6.0	
4	536365	84029E	RED WOOLLY HOTTIE WHITE HEART.	6.0	
...
75101	542541	22124	SET OF 2 TEA TOWELS PING MICROWAVE	1.0	
75102	542541	22131	FOOD CONTAINER SET 3 LOVE HEART	1.0	
75103	542541	22135	MINI LADLE LOVE HEART PINK	1.0	
75104	542541	22148	EASTER CRAFT 4 CHICKS	3.0	
75105	542541	22149	FELTCRAFT 6 FLOWER FRIENDS	NaN	

	InvoiceDate	UnitPrice	CustomerID	Country
0	12/1/2010 8:26	2.55	17850.0	United Kingdom
1	12/1/2010 8:26	3.39	17850.0	United Kingdom
2	12/1/2010 8:26	2.75	17850.0	United Kingdom
3	12/1/2010 8:26	3.39	17850.0	United Kingdom
4	12/1/2010 8:26	3.39	17850.0	United Kingdom
...
75101	1/28/2011 14:25	2.46	NaN	United Kingdom
75102	1/28/2011 14:25	4.13	NaN	United Kingdom
75103	1/28/2011 14:25	0.83	NaN	United Kingdom
75104	1/28/2011 14:25	4.13	NaN	United Kingdom
75105	NaN	NaN	NaN	NaN

[75106 rows x 8 columns]

```
import pandas as pd

# Path to your CSV file
csv_file_path = '/content/data (1).csv'

# Read the CSV file directly
try:
    df = pd.read_csv(csv_file_path, encoding='ISO-8859-1')
    # Display the data
    print(df)
except FileNotFoundError:
    print(f"Error: The file {csv_file_path} was not found.")
except Exception as e:
    print(f"An error occurred while reading the CSV file: {e}")

# Drop non-numeric columns that are not suitable for the model
df = df.drop(['InvoiceNo', 'Description'], axis=1)

# Separate X and y
X = df.iloc[:, :-1] # all columns except the last one
y = df.iloc[:, -1] # last column as target

# Display X and y
print("Features (X):")
print(X.head()) # first 5 rows of features

print("\nTarget (y):")
print(y.head()) # first 5 rows of target
```

Drop non-numeric columns that are not suitable for the model

Separate X and y

Display X and y

print(y.head()) # first 5 rows of target

	InvoiceNo	StockCode	Description	Quantity	\
0	536365	85123A	WHITE HANGING HEART T-LIGHT HOLDER	6.0	
1	536365	71053	WHITE METAL LANTERN	6.0	
2	536365	84406B	CREAM CUPID HEARTS COAT HANGER	8.0	
3	536365	84029G	KNITTED UNION FLAG HOT WATER BOTTLE	6.0	
4	536365	84029E	RED WOOLLY HOTTIE WHITE HEART.	6.0	
...
75101	542541	22124	SET OF 2 TEA TOWELS PING MICROWAVE	1.0	
75102	542541	22131	FOOD CONTAINER SET 3 LOVE HEART	1.0	
75103	542541	22135	MINI LADLE LOVE HEART PINK	1.0	
75104	542541	22148	EASTER CRAFT 4 CHICKS	3.0	

```

75105    542541    22149    FELTCRAFT 6 FLOWER FRIENDS    NaN

      InvoiceDate  UnitPrice  CustomerID    Country
0    12/1/2010 8:26      2.55    17850.0    United Kingdom
1    12/1/2010 8:26      3.39    17850.0    United Kingdom
2    12/1/2010 8:26      2.75    17850.0    United Kingdom
3    12/1/2010 8:26      3.39    17850.0    United Kingdom
4    12/1/2010 8:26      3.39    17850.0    United Kingdom
...      ...      ...      ...      ...
75101  1/28/2011 14:25      2.46      NaN    United Kingdom
75102  1/28/2011 14:25      4.13      NaN    United Kingdom
75103  1/28/2011 14:25      0.83      NaN    United Kingdom
75104  1/28/2011 14:25      4.13      NaN    United Kingdom
75105      NaN      NaN      NaN      NaN

```

[75106 rows x 8 columns]

Features (X):

```

      StockCode  Quantity    InvoiceDate  UnitPrice  CustomerID
0    85123A      6.0    12/1/2010 8:26      2.55    17850.0
1    71053      6.0    12/1/2010 8:26      3.39    17850.0
2    84406B      8.0    12/1/2010 8:26      2.75    17850.0
3    84029G      6.0    12/1/2010 8:26      3.39    17850.0
4    84029E      6.0    12/1/2010 8:26      3.39    17850.0

```

Target (y):

```

0    United Kingdom
1    United Kingdom
2    United Kingdom
3    United Kingdom
4    United Kingdom

```

Name: Country, dtype: object

```

from sklearn.model_selection import train_test_split

# Split data into training and testing sets (80% train, 20% test)
X_train, X_test, y_train, y_test = train_test_split(
    X, y, test_size=0.2, random_state=42
)

# Display the shapes of the splits
print("X_train shape:", X_train.shape)
print("X_test shape:", X_test.shape)
print("y_train shape:", y_train.shape)
print("y_test shape:", y_test.shape)

```

```

X_train shape: (60084, 5)
X_test shape: (15022, 5)
y_train shape: (60084,)
y_test shape: (15022,)

```

SHAP VALUES

```

import pandas as pd
import numpy as np
from sklearn.model_selection import train_test_split
from sklearn.ensemble import RandomForestClassifier
from sklearn.metrics import accuracy_score
import shap
from sklearn.preprocessing import LabelEncoder
import matplotlib.pyplot as plt

# Load dataset with encoding to handle special characters
data = pd.read_csv('/content/data (1).csv', encoding='ISO-8859-1')

# Inspect column names to find target
print("Columns in dataset:", data.columns)

# Set the target column
target_column = 'Country'

# Separate features and target
X = data.drop(target_column, axis=1)
y = data[target_column]

# Data cleaning and preprocessing
X['InvoiceDate'] = pd.to_datetime(X['InvoiceDate'], errors='coerce')
combined = pd.concat([X, y], axis=1)
combined = combined.dropna(subset=['InvoiceDate'])
X = combined.drop(target_column, axis=1)
y = combined[target_column]
X['InvoiceDate'] = X['InvoiceDate'].astype(np.int64) // 10**9

```

```

combined = pd.concat([X, y], axis=1)
combined = combined.dropna(subset=['CustomerID'])
X = combined.drop(target_column, axis=1)
y = combined[target_column]

combined = pd.concat([X, y], axis=1)
combined = combined.dropna()
X = combined.drop(target_column, axis=1)
y = combined[target_column]

# Encode categorical columns
for col in X.select_dtypes(include=['object']).columns:
    X[col] = LabelEncoder().fit_transform(X[col].astype(str))

# Split into train and test sets
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.3, random_state=42)

assert not X_train.isnull().values.any(), "X_train contains NaNs!"
assert not X_test.isnull().values.any(), "X_test contains NaNs!"
assert not y_train.isnull().values.any(), "y_train contains NaNs!"
assert not y_test.isnull().values.any(), "y_test contains NaNs!"

# Train Random Forest model
model = RandomForestClassifier(n_estimators=100, random_state=42)
model.fit(X_train, y_train)

# Evaluate model
y_pred = model.predict(X_test)
print("Model Accuracy on test data:", accuracy_score(y_test, y_pred))

# SHAP explanation starts here
subset_size = 100
X_test_subset = X_test.sample(n=min(subset_size, len(X_test)), random_state=42)

explainer = shap.TreeExplainer(model)
shap_values = explainer.shap_values(X_test_subset)

feature_values = X_test_subset.values
feature_names = X_test_subset.columns.tolist()
base_values = explainer.expected_value

shap_explanation = shap.Explanation(values=shap_values,
                                  base_values=base_values,
                                  data=feature_values,
                                  feature_names=feature_names)

# --- Global Explanation: Bar Plot ---
print("\n--- SHAP Summary Bar Plot ---")
print("This bar plot shows which features contributed the most, on average, to the model's predictions for the first class. Longer bars mean higher importance. It helps identify which features the model relies on the most.")
shap.summary_plot(shap_explanation[:, :, 0], plot_type="bar")
plt.show()

# --- Global Explanation: Beeswarm Plot ---
print("\n--- SHAP Summary Beeswarm Plot ---")
print("This beeswarm plot shows the distribution of impacts each feature had across the test samples.")
print("Each dot represents one sample. Color shows feature value: red = high, blue = low.")
print("This helps you understand how feature values influence predictions, both positively and negatively.")
shap.summary_plot(shap_explanation[:, :, 0])
plt.show()

# --- Local Explanation: Force Plot ---
bar = 0
print(f"\n--- SHAP Force Plot for Test Instance {bar} ---")
print("This force plot explains how the individual features pushed the model's prediction away from the base value for this instance. Red arrows push the prediction towards the class, blue arrows push it away. The base value is the average prediction for the test set.")
shap.initjs()
shap.force_plot(shap_explanation[bar, :, 0])

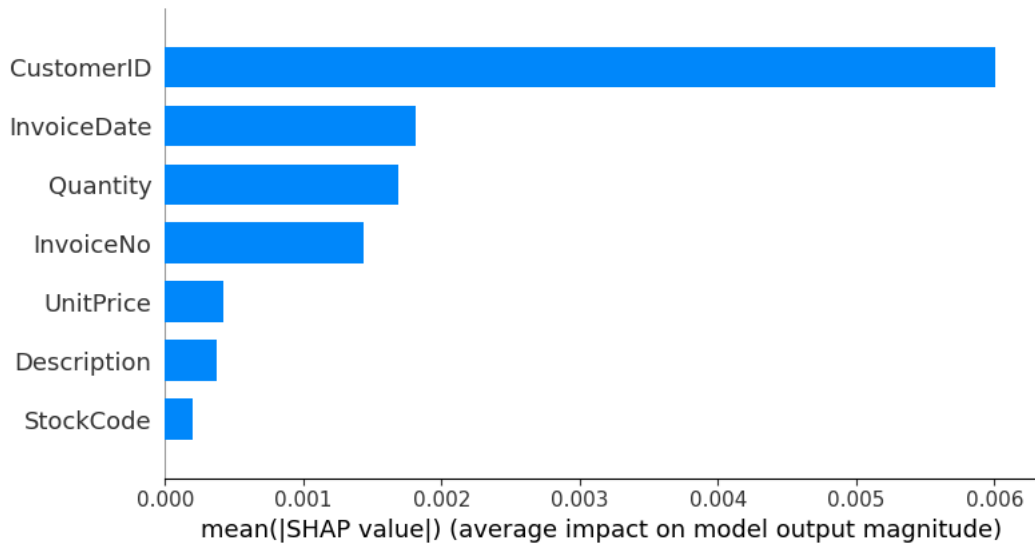
```

```
Columns in dataset: Index(['InvoiceNo', 'StockCode', 'Description', 'Quantity', 'InvoiceDate',
                          'UnitPrice', 'CustomerID', 'Country'],
                          dtype='object')
```

Model Accuracy on test data: 0.9972951811516834

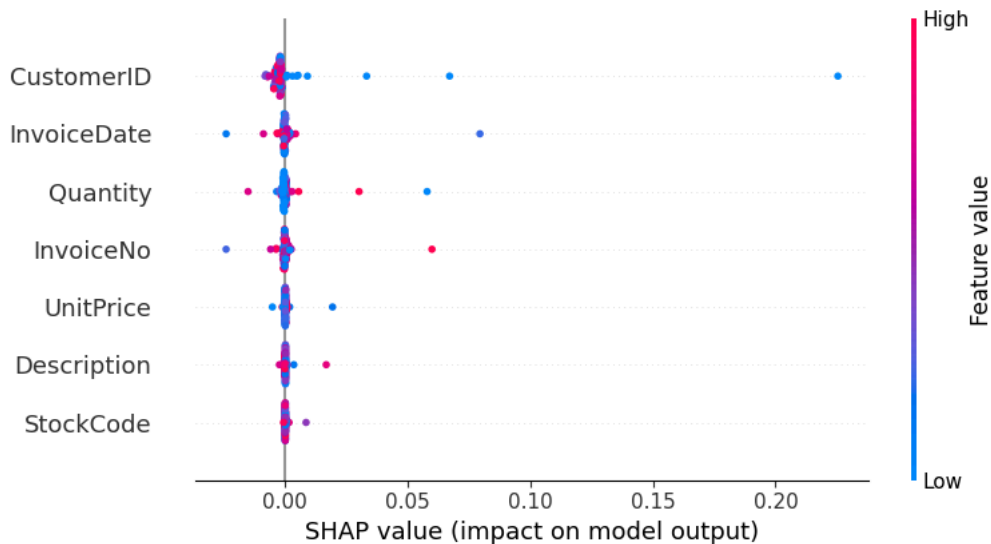
--- SHAP Summary Bar Plot ---

This bar plot shows which features contributed the most, on average, to the model's predictions for the first class. Longer bars mean higher importance. It helps identify which features the model relies on the most.



--- SHAP Summary Beeswarm Plot ---

This beeswarm plot shows the distribution of impacts each feature had across the test samples. Each dot represents one sample. Color shows feature value: red = high, blue = low. This helps you understand how feature values influence predictions, both positively and negatively.



--- SHAP Force Plot for Test Instance 0 ---

This force plot explains how the individual features pushed the model's prediction away from the base value for this sample. Red arrows push the prediction towards the class, blue arrows push it away. The base value is the average prediction.



PARTIAL DEPENDENCE

```
import pandas as pd
import numpy as np
from sklearn.model_selection import train_test_split
from sklearn.ensemble import RandomForestClassifier
from sklearn.metrics import accuracy_score
from sklearn.preprocessing import LabelEncoder
from sklearn.inspection import PartialDependenceDisplay
import matplotlib.pyplot as plt
```

```

# Load dataset with proper encoding
try:
    data = pd.read_csv('/content/data (1).csv', encoding='ISO-8859-1')
except Exception as e:
    print(f"Error loading data: {e}")
    exit()

print("Columns:", data.columns)

# Set the target column name here
target_column = 'Country'

# Separate features and target
X = data.drop(columns=[target_column, 'InvoiceNo', 'Description'])
y = data[target_column]

# --- Data Cleaning and Preprocessing ---
X['InvoiceDate'] = pd.to_datetime(X['InvoiceDate'], errors='coerce')
combined = pd.concat([X, y], axis=1)
combined = combined.dropna(subset=['InvoiceDate'])
X = combined.drop(target_column, axis=1)
y = combined[target_column]
X['InvoiceDate'] = X['InvoiceDate'].astype(np.int64) // 10**9

combined = pd.concat([X, y], axis=1)
combined = combined.dropna(subset=['CustomerID'])
X = combined.drop(target_column, axis=1)
y = combined[target_column]

combined = pd.concat([X, y], axis=1)
combined = combined.dropna()
X = combined.drop(target_column, axis=1)
y = combined[target_column]

# Encode categorical variables
for col in X.select_dtypes(include=['object']).columns:
    le = LabelEncoder()
    X[col] = le.fit_transform(X[col].astype(str))

if y.dtype == 'object':
    le_target = LabelEncoder()
    y = le_target.fit_transform(y)

for col in X.columns:
    X[col] = X[col].astype(float)

# Split into train and test sets
X_train, X_test, y_train, y_test = train_test_split(
    X, y, test_size=0.2, random_state=42
)

# Train Random Forest Classifier
model = RandomForestClassifier(random_state=42)
model.fit(X_train, y_train)

# Evaluate the model
y_pred = model.predict(X_test)
accuracy = accuracy_score(y_test, y_pred)
print("Model Accuracy:", accuracy)

# Choose features for partial dependence plot
features = [0, 1] # You can adjust this based on your dataset

# Plot Partial Dependence
print("\n--- Partial Dependence Plot ---")
print("This plot shows how changing the values of selected features affects the model's predictions on average.")
print("The x-axis represents values of the feature, and the y-axis shows the average predicted probability or output.")
print("A rising curve means that increasing the feature value leads to a higher prediction for the target class.")
print("A flat curve means that changing the feature has little effect on the prediction.")
print("This helps to understand the relationship between features and predictions, assuming all other features are averaged")

PartialDependenceDisplay.from_estimator(
    model,
    X_train,
    features=features,
    feature_names=X.columns,
    grid_resolution=50,
    target=0 # Example: plot for the first class
)

plt.tight_layout()

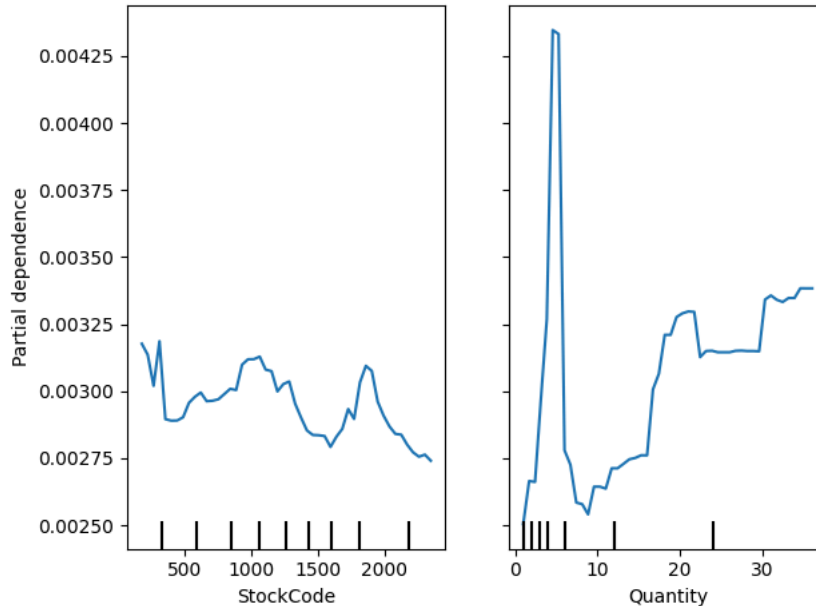
```

```
plt.show()
```

```
Columns: Index(['InvoiceNo', 'StockCode', 'Description', 'Quantity', 'InvoiceDate',
               'UnitPrice', 'CustomerID', 'Country'],
              dtype='object')
Model Accuracy: 0.9981849241938928
```

--- Partial Dependence Plot ---

This plot shows how changing the values of selected features affects the model's predictions on average. The x-axis represents values of the feature, and the y-axis shows the average predicted probability or output. A rising curve means that increasing the feature value leads to a higher prediction for the target class. A flat curve means that changing the feature has little effect on the prediction. This helps to understand the relationship between features and predictions, assuming all other features are averaged out.



ICE PLOT

```
import pandas as pd
import numpy as np
from sklearn.model_selection import train_test_split
from sklearn.ensemble import RandomForestClassifier
from sklearn.preprocessing import LabelEncoder
from sklearn.inspection import PartialDependenceDisplay
import matplotlib.pyplot as plt

# Load dataset with proper encoding
try:
    data = pd.read_csv('/content/data (1).csv', encoding='ISO-8859-1')
except Exception as e:
    print(f"Error loading data: {e}")
    # Exit cell execution if data loading fails
    exit()

print("Columns:", data.columns)

# Set the target column name here
target_column = 'Country'

# Separate features and target
# Drop columns that are not suitable for the model or are the target
X = data.drop(columns=[target_column, 'InvoiceNo', 'Description'])
y = data[target_column]

# --- Data Cleaning and Preprocessing ---
# Convert 'InvoiceDate' to datetime and then to numerical timestamp
X['InvoiceDate'] = pd.to_datetime(X['InvoiceDate'], errors='coerce')

# Drop rows where 'InvoiceDate' conversion failed
combined = pd.concat([X, y], axis=1)
combined = combined.dropna(subset=['InvoiceDate'])
X = combined.drop(target_column, axis=1)
y = combined[target_column]

X['InvoiceDate'] = X['InvoiceDate'].astype(np.int64) // 10**9
```

```

# Handle missing CustomerID by dropping those rows
combined = pd.concat([X, y], axis=1)
combined = combined.dropna(subset=['CustomerID'])
X = combined.drop(target_column, axis=1)
y = combined[target_column]

# Drop any remaining rows with NaN values in X and y
combined = pd.concat([X, y], axis=1)
combined = combined.dropna()
X = combined.drop(target_column, axis=1)
y = combined[target_column]

# Encode categorical variables if any (after dropping non-numeric and NaNs)
for col in X.select_dtypes(include=['object']).columns:
    le = LabelEncoder()
    X[col] = le.fit_transform(X[col].astype(str))

# Encode target if it's categorical (after dropping NaNs in y)
if y.dtype == 'object':
    le_target = LabelEncoder()
    y = le_target.fit_transform(y)

# Explicitly cast all feature columns to float to ensure numerical consistency
for col in X.columns:
    X[col] = X[col].astype(float)

# Split into train and test sets
X_train, X_test, y_train, y_test = train_test_split(
    X, y, test_size=0.4, random_state=42
)

# Train Random Forest Classifier
model = RandomForestClassifier(random_state=42)
model.fit(X_train, y_train)

# Evaluate the model
y_pred = model.predict(X_test)
accuracy = accuracy_score(y_test, y_pred)
print("Accuracy:", accuracy)

# Select the feature(s) for ICE plot
# Example: use the first feature; change as needed
feature_index = 0
feature_name = X.columns[feature_index]

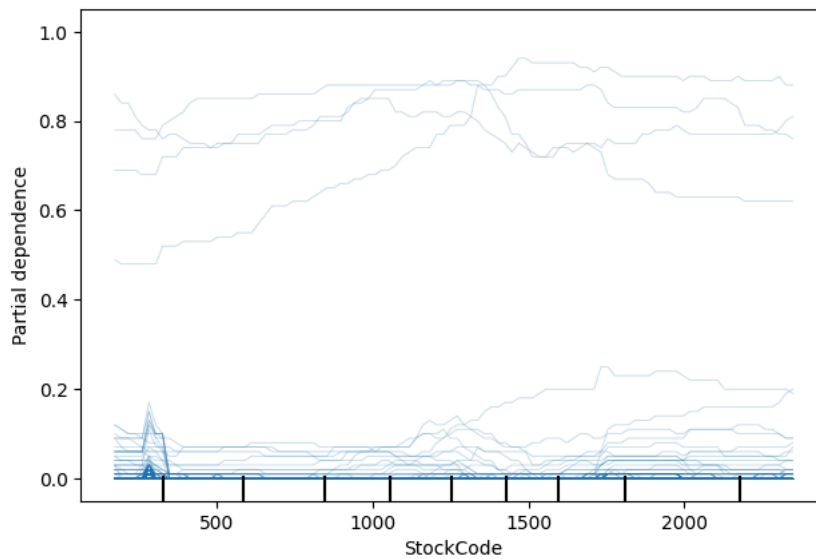
# Plot ICE curves
PartialDependenceDisplay.from_estimator(
    model,
    X_train,
    features=[feature_index],
    kind='individual',
    feature_names=X.columns,
    grid_resolution=100,
    target=0 # Specify the target class for which to plot ICE (e.g., 0 for the first class)
)

plt.suptitle(f'ICE Plot for Feature: {feature_name}', fontsize=18)
plt.tight_layout()
plt.show()

```

```
Columns: Index(['InvoiceNo', 'StockCode', 'Description', 'Quantity', 'InvoiceDate',
               'UnitPrice', 'CustomerID', 'Country'],
              dtype='object')
Accuracy: 0.9970638479607089
```

ICE Plot for Feature: StockCode



decision tree

```
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
from sklearn.model_selection import train_test_split
from sklearn.tree import DecisionTreeClassifier
from sklearn.metrics import accuracy_score, classification_report, confusion_matrix
from sklearn.preprocessing import LabelEncoder

# Load the dataset
data = pd.read_csv('/content/data (1).csv', encoding='ISO-8859-1')

# Handle missing values
data = data.dropna()

# Encode categorical columns if any
label_encoders = {}
for column in data.select_dtypes(include=['object']).columns:
    le = LabelEncoder()
    data[column] = le.fit_transform(data[column])
    label_encoders[column] = le

# Define features and target
X = data.iloc[:, :-1]
y = data.iloc[:, -1]

# Split into train and test sets
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)

# Create and train the Decision Tree classifier
model = DecisionTreeClassifier(random_state=42)
model.fit(X_train, y_train)

# Make predictions
y_pred = model.predict(X_test)

# Evaluate the model
accuracy = accuracy_score(y_test, y_pred)
print(f"Accuracy: {accuracy:.4f}\n")
print("Classification Report:")
print(classification_report(y_test, y_pred))
print("Confusion Matrix:")
cm = confusion_matrix(y_test, y_pred)
print(cm)

# Plot the confusion matrix as a box (heatmap)
plt.figure(figsize=(8, 6))
sns.heatmap(cm, annot=True, fmt='d', cmap='Blues', cbar=True, square=True)
plt.title("Confusion Matrix - Decision Tree")
plt.xlabel("Predicted Labels")
```



```
plt.ylabel("True Labels")
plt.show()
```

Accuracy: 0.9990

Classification Report:

	precision	recall	f1-score	support
0	1.00	1.00	1.00	31
1	1.00	0.67	0.80	3
2	1.00	1.00	1.00	31
3	1.00	1.00	1.00	9
4	1.00	1.00	1.00	19
5	1.00	1.00	1.00	3
6	1.00	0.97	0.99	113
7	1.00	1.00	1.00	4
8	1.00	1.00	1.00	189
9	1.00	1.00	1.00	255
10	1.00	1.00	1.00	5
11	1.00	1.00	1.00	9
13	1.00	1.00	1.00	29
14	1.00	1.00	1.00	18
15	1.00	1.00	1.00	10
16	1.00	1.00	1.00	7
17	0.97	0.97	0.97	70
18	1.00	1.00	1.00	31
19	0.91	1.00	0.95	10
20	1.00	1.00	1.00	47
21	1.00	1.00	1.00	8
22	1.00	1.00	1.00	98
23	1.00	1.00	1.00	14
24	0.97	0.94	0.96	36
25	1.00	1.00	1.00	8317
accuracy			1.00	9366
macro avg	0.99	0.98	0.99	9366
weighted avg	1.00	1.00	1.00	9366

Confusion Matrix:

```
[[ 31  0  0  0  0  0  0  0  0  0  0  0  0  0  0
   0  0  0  0  0  0  0  0  0  0  0  0  0  0
   0  2  0  0  0  0  0  0  0  0  0  0  0  0  0
   0  0  0  0  0  0  0  0  0  0  0  0  1  0  0
   0  0  31  0  0  0  0  0  0  0  0  0  0  0  0
   0  0  0  9  0  0  0  0  0  0  0  0  0  0  0
   0  0  0  0  0  0  0  0  0  0  0  0  0  0  0
   0  0  0  0  19  0  0  0  0  0  0  0  0  0  0
   0  0  0  0  0  0  0  0  0  0  0  0  0  0  0
   0  0  0  0  0  3  0  0  0  0  0  0  0  0  0
   0  0  0  0  0  0  0  0  0  0  0  0  0  0  0
   0  0  0  0  0  0  110  0  0  0  0  0  0  0  0
   0  0  0  0  0  0  0  0  0  0  0  3  0  0  0
   0  0  0  0  0  0  0  0  4  0  0  0  0  0  0
   0  0  0  0  0  0  0  0  0  0  0  0  0  0  0
   0  0  0  0  0  0  0  0  0  189  0  0  0  0  0
   0  0  0  0  0  0  0  0  0  0  0  0  0  0  0
   0  0  0  0  0  0  0  0  0  0  255  0  0  0  0
   0  0  0  0  0  0  0  0  0  0  0  0  0  0  0
   0  0  0  0  0  0  0  0  0  0  0  5  0  0  0
   0  0  0  0  0  0  0  0  0  0  0  0  0  0  0
   0  0  0  0  0  0  0  0  0  0  0  0  9  0  0
   0  0  0  0  0  0  0  0  0  0  0  0  0  0  0
   0  0  0  0  0  0  0  0  0  0  0  0  0  29  0]
```

random forest

```
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
from sklearn.model_selection import train_test_split
from sklearn.ensemble import RandomForestClassifier
from sklearn.metrics import accuracy_score, classification_report, confusion_matrix
from sklearn.preprocessing import LabelEncoder

# Load the dataset
data = pd.read_csv('/content/data (1).csv', encoding='ISO-8859-1')

# Handle missing values
data = data.dropna()

# Encode categorical columns if any
label_encoders = {}
for column in data.select_dtypes(include=['object']).columns:
    le = LabelEncoder()
    data[column] = le.fit_transform(data[column])
    label_encoders[column] = le

# Define features and target
```

```
X = data.iloc[:, :-1]
y = data.iloc[:, -1]

# Split into training and testing sets
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)

# Create and train the Random Forest classifier
model = RandomForestClassifier(n_estimators=100, random_state=42)
model.fit(X_train, y_train)

# Make predictions
y_pred = model.predict(X_test)

# Evaluate the model
accuracy = accuracy_score(y_test, y_pred)
print(f"Accuracy: {accuracy:.4f}\n")
print("Classification Report:")
print(classification_report(y_test, y_pred))

# Confusion matrix
cm = confusion_matrix(y_test, y_pred)
print("Confusion Matrix:")
print(cm)

# Plot the confusion matrix as a box (heatmap)
plt.figure(figsize=(8, 6))
sns.heatmap(cm, annot=True, fmt='d', cmap='Blues', cbar=True, square=True)
plt.title("Confusion Matrix - Random Forest")
plt.xlabel("Predicted Labels")
plt.ylabel("True Labels")
plt.show()
```



```
data[column] = le.fit_transform(data[column])
label_encoders[column] = le

# Define features and target
X = data.iloc[:, :-1]
y = data.iloc[:, -1]

# Split into train and test sets
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)

# Create and train the XGBoost classifier
model = XGBClassifier(use_label_encoder=False, eval_metric='mlogloss', random_state=42)
model.fit(X_train, y_train)

# Make predictions
y_pred = model.predict(X_test)

# Evaluate the model
accuracy = accuracy_score(y_test, y_pred)
print(f"Accuracy: {accuracy:.4f}\n")
print("Classification Report:")
print(classification_report(y_test, y_pred))

# Confusion matrix
cm = confusion_matrix(y_test, y_pred)
print("Confusion Matrix:")
print(cm)

# Plot the confusion matrix as a box (heatmap)
plt.figure(figsize=(8, 6))
sns.heatmap(cm, annot=True, fmt='d', cmap='Blues', cbar=True, square=True)
plt.title("Confusion Matrix - XGBoost")
plt.xlabel("Predicted Labels")
plt.ylabel("True Labels")
plt.show()
```

```
/usr/local/lib/python3.12/dist-packages/xgboost/training.py:183: UserWarning: [18:12:49] WARNING: /workspace/src/learner.c
Parameters: { "use_label_encoder" } are not used.
```

```
bst.update(dtrain, iteration=i, fobj=obj)
Accuracy: 0.9989
```

Classification Report:

	precision	recall	f1-score	support
0	1.00	1.00	1.00	31
1	1.00	0.67	0.80	3
2	1.00	1.00	1.00	31
3	1.00	1.00	1.00	9
4	1.00	1.00	1.00	19
5	1.00	1.00	1.00	3
6	1.00	1.00	1.00	113
7	1.00	0.50	0.67	4
8	0.99	0.98	0.99	189
9	0.99	1.00	1.00	255
10	1.00	1.00	1.00	5
11	0.90	1.00	0.95	9
13	1.00	0.97	0.98	29
14	1.00	0.94	0.97	18
15	1.00	1.00	1.00	10
16	1.00	1.00	1.00	7
17	1.00	0.99	0.99	70
18	0.97	1.00	0.98	31
19	1.00	1.00	1.00	10
20	1.00	1.00	1.00	47
21	1.00	1.00	1.00	8
22	1.00	1.00	1.00	98
23	1.00	1.00	1.00	14
24	1.00	0.97	0.99	36
25	1.00	1.00	1.00	8317
accuracy			1.00	9366
macro avg	0.99	0.96	0.97	9366
weighted avg	1.00	1.00	1.00	9366

Confusion Matrix:

```
[[ 31  0  0  0  0  0  0  0  0  0  0  0  0  0
   0  0  0  0  0  0  0  0  0  0  0  0  0
   0  2  0  0  0  0  0  0  0  0  0  0  0  0
   0  0  0  0  0  0  0  0  0  0  0  1  0  0
   0  0 31  0  0  0  0  0  0  0  0  0  0  0
   0  0  0  0  0  0  0  0  0  0  0  0  0  0
   0  0  0  9  0  0  0  0  0  0  0  0  0  0
   0  0  0  0  0  0  0  0  0  0  0  0  0  0
   0  0  0  0 19  0  0  0  0  0  0  0  0  0
   0  0  0  0  0  0  0  0  0  0  0  0  0  0
   0  0  0  0  0  3  0  0  0  0  0  0  0  0
   0  0  0  0  0  0  0  0  0  0  0  0  0  0
   0  0  0  0  0  0 113  0  0  0  0  0  0  0
   0  0  0  0  0  0  0  0  0  0  0  0  0  0
   0  0  0  0  0  0  0  2  1  0  0  1  0  0
   0  0  0  0  0  0  0  0  0  0  0  0  0  0
   0  0  0  0  0  0  0  0  0 186  2  0  0  0
   0  0  0  1  0  0  0  0  0  0  0  0  0  0
   0  0  0  0  0  0  0  0  0  0 255  0  0  0
   0  0  0  0  0  0  0  0  0  0  0  1  0  0]
```

lightGBM

```
pip install lightgbm
```

```
Requirement already satisfied: lightgbm in /usr/local/lib/python3.12/dist-packages (4.6.0)
Requirement already satisfied: numpy>=1.27.0 in /usr/local/lib/python3.12/dist-packages (from lightgbm) (2.0.2)
Requirement already satisfied: scipy in /usr/local/lib/python3.12/dist-packages (from lightgbm) (1.16.1)
```

```
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
from sklearn.model_selection import train_test_split
from lightgbm import LGBMClassifier
from sklearn.metrics import accuracy_score, classification_report, confusion_matrix
from sklearn.preprocessing import LabelEncoder

# Load the dataset
data = pd.read_csv('/content/data (1).csv', encoding='ISO-8859-1')

# Handle missing values
data = data.dropna()

# Encode categorical columns if any
label_encoders = {}
for column in data.select_dtypes(include='object').columns:
    le = LabelEncoder()
    data[column] = le.fit_transform(data[column])
```

```
label_encoders[column] = le

# Define features and target
X = data.iloc[:, :-1]
y = data.iloc[:, -1]

# Split into train and test sets
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)

# Create and train the LightGBM classifier
model = LGBMClassifier(random_state=42)
model.fit(X_train, y_train)

# Make predictions
y_pred = model.predict(X_test)

# Evaluate the model
accuracy = accuracy_score(y_test, y_pred)
print(f"Accuracy: {accuracy:.4f}\n")
print("Classification Report:")
print(classification_report(y_test, y_pred))

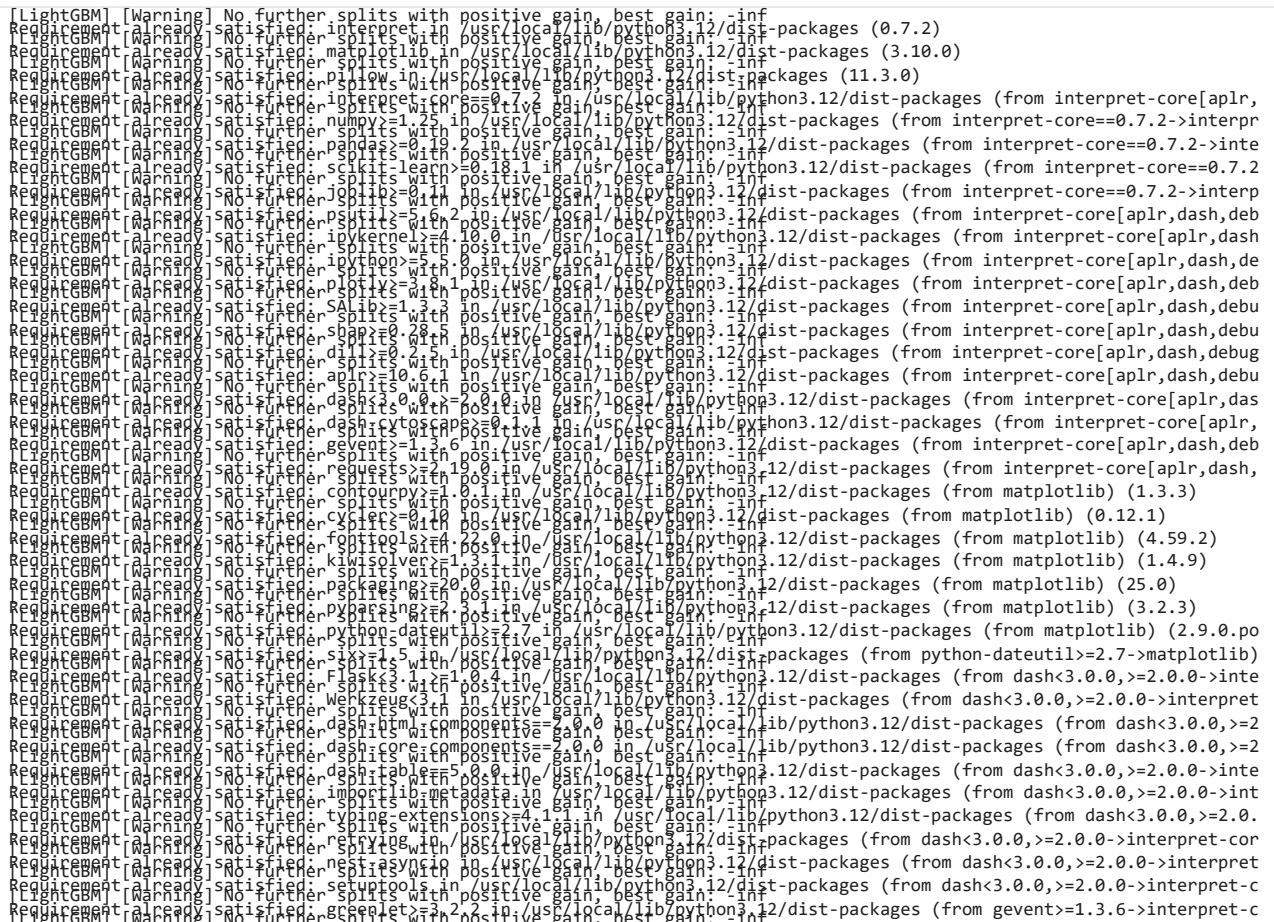
# Confusion matrix
cm = confusion_matrix(y_test, y_pred)
print("Confusion Matrix:")
print(cm)

# Plot the confusion matrix as a box (heatmap)
plt.figure(figsize=(8, 6))
sns.heatmap(cm, annot=True, fmt='d', cmap='Blues', cbar=True, square=True)
plt.title("Confusion Matrix - LightGBM")
plt.xlabel("Predicted Labels")
plt.ylabel("True Labels")
plt.show()
```

linear regression

linear regression

Residuals Plot for Linear Regression



Explainable Boosting

[illegible]

```
import pandas as pd
import numpy as np
from sklearn.model_selection import train_test_split
from sklearn.preprocessing import LabelEncoder
from interpret.glassbox import ExplainableBoostingClassifier
```

[illegible]


```
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns

from sklearn.model_selection import train_test_split
from sklearn.metrics import accuracy_score, confusion_matrix, r2_score
from sklearn.preprocessing import LabelEncoder

# Classifiers
from sklearn.tree import DecisionTreeClassifier
from sklearn.ensemble import RandomForestClassifier
from xgboost import XGBClassifier
from lightgbm import LGBMClassifier

# Regressor
from sklearn.linear_model import LinearRegression

# Load dataset
data = pd.read_csv('/content/data (1).csv', encoding='ISO-8859-1')
data = data.dropna()

# Encode categorical columns if any
for column in data.select_dtypes(include=['object']).columns:
    le = LabelEncoder()
    data[column] = le.fit_transform(data[column])

# Features and target
X = data.iloc[:, :-1]
y = data.iloc[:, -1]

# Split
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)
```

search.google.com/drive/1MIZLCvT8SVBxyKb2D0e61CJTIEAmJDKI#printMode=true

lime

lime

lime

```
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
from sklearn.model_selection import train_test_split
from sklearn.preprocessing import LabelEncoder, StandardScaler
from sklearn.ensemble import RandomForestClassifier
from sklearn.metrics import accuracy_score
import lime
import lime.lime_tabular

# Load dataset
data = pd.read_csv('/content/data (1).csv', encoding='ISO-8859-1')

# Drop non-numeric columns that are not suitable for the model
data = data.drop(['InvoiceNo', 'Description'], axis=1)

# Handle missing values by dropping rows with any NaN values
data = data.dropna()

# Encode categorical columns if any exist after dropping
label_encoders = {}
for column in data.select_dtypes(include=['object']).columns:
    le = LabelEncoder()
    data[column] = le.fit_transform(data[column])
    label_encoders[column] = le

# Define features and target
X = data.iloc[:, :-1].values
y = data.iloc[:, -1].values

# Encode target if it's categorical
if y.dtype == 'object':
    le = LabelEncoder()
    y = le.fit_transform(y)

# Split data
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)

# Scale features
scaler = StandardScaler()
X_train = scaler.fit_transform(X_train)
X_test = scaler.transform(X_test)

# Track accuracy over different numbers of estimators
n_estimators_list = [10, 50, 100, 150, 200, 250, 300]
train_accuracies = []
test_accuracies = []

for n in n_estimators_list:
    model = RandomForestClassifier(n_estimators=n, random_state=42)
    model.fit(X_train, y_train)
    y_train_pred = model.predict(X_train)
    y_test_pred = model.predict(X_test)
    train_acc = accuracy_score(y_train, y_train_pred)
    test_acc = accuracy_score(y_test, y_test_pred)
    train_accuracies.append(train_acc)
    test_accuracies.append(test_acc)

# Plot training and testing accuracy curves
plt.figure(figsize=(10,6))
plt.plot(n_estimators_list, train_accuracies, marker='o', label='Training Accuracy')
plt.plot(n_estimators_list, test_accuracies, marker='s', label='Testing Accuracy')
plt.xlabel('Number of Estimators')
plt.ylabel('Accuracy')
plt.title('Training and Testing Accuracy vs Number of Estimators')
plt.legend()
plt.grid(True)
plt.tight_layout()
plt.show()

# Train final model for LIME explanation
model = RandomForestClassifier(n_estimators=100, random_state=42)
model.fit(X_train, y_train)

# Predict and print final test accuracy
```



```
y_pred = model.predict(X_test)
print(f"Final Test Accuracy: {accuracy_score(y_test, y_pred)*100:.2f}%")

# Create LIME explainer
explainer = lime.lime_tabular.LimeTabularExplainer/
```