# EC 9560 DATA MINING
# LAB 03

**PRIYATHARSINI S.**
**2020E122**
**SEMESTER 7**
**30 OCT 2024**

## Check Null Values

```
[36]  1  # Check for null values
      2  null_values = data.isnull().sum()
      3  print("Null Values in Each Column:")
      4  print(null_values)
```

```
Null Values in Each Column:
id                          0
person_age                  0
person_income               0
person_home_ownership       0
person_emp_length           0
loan_intent                 0
loan_grade                  0
loan_amnt                   0
loan_int_rate               0
loan_percent_income         0
cb_person_default_on_file   0
cb_person_cred_hist_length  0
loan_status                 0
dtype: int64
```

## Check Duplicate Values

```
[ ]  1  # Check for duplicates
     2  duplicate_rows = data.duplicated().sum()
     3  print("\nTotal Duplicate Rows:", duplicate_rows)
```

```
Total Duplicate Rows: 0
```

```
[9]  1  data.shape
```

```
(58645, 13)
```

```
[11]    1    data.dtypes
0s
```

| | 0 |
|---|---|
| id | int64 |
| person_age | int64 |
| person_income | int64 |
| person_home_ownership | object |
| person_emp_length | float64 |
| loan_intent | object |
| loan_grade | object |
| loan_amnt | int64 |
| loan_int_rate | float64 |
| loan_percent_income | float64 |
| cb_person_default_on_file | object |
| cb_person_cred_hist_length | int64 |
| loan_status | int64 |

## Encode the Categorical Columns

```python
from sklearn.preprocessing import LabelEncoder

# Create DataFrame
df = pd.DataFrame(data)

# Initialize the LabelEncoder
label_encoder = LabelEncoder()

# Columns to be label encoded
categorical_columns = [
    'person_home_ownership',
    'loan_intent',
    'loan_grade',
    'cb_person_default_on_file'
]

# Convert person_emp_length to int
df['person_emp_length'] = pd.to_numeric(df['person_emp_length'],
errors='coerce').fillna(0).astype(int)

# Apply label encoding to categorical columns
for column in categorical_columns:
    df[column] = label_encoder.fit_transform(df[column])

# Display the updated DataFrame
df.head()
```

```
1  df.dtypes
```

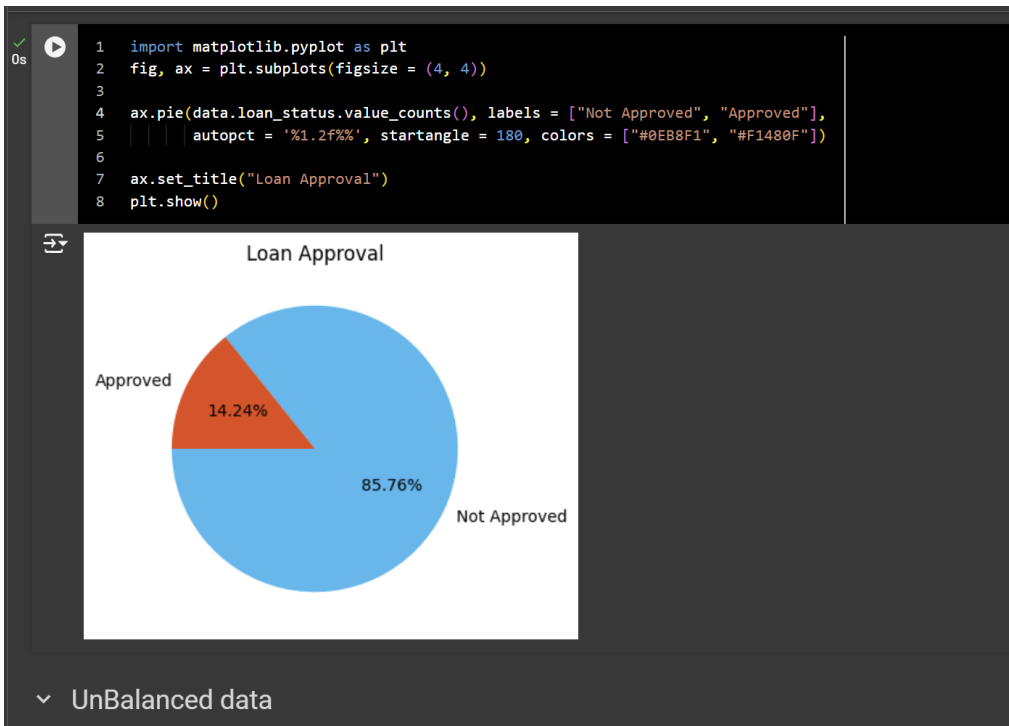| | 0 |
|---|---|
| id | int64 |
| person_age | int64 |
| person_income | int64 |
| person_home_ownership | int64 |
| person_emp_length | int64 |
| loan_intent | int64 |
| loan_grade | int64 |
| loan_amnt | int64 |
| loan_int_rate | float64 |
| loan_percent_income | float64 |
| cb_person_default_on_file | int64 |
| cb_person_cred_hist_length | int64 |
| loan_status | int64 |

dtype: object

```
[12]  1  #print target classes and its count
      2  data['loan_status'].value_counts()
```

| | count |
|---|---|
| loan_status | |
| 0 | 50295 |
| 1 | 8350 |

dtype: int64

```python
1  import matplotlib.pyplot as plt
2  fig, ax = plt.subplots(figsize = (4, 4))
3
4  ax.pie(data.loan_status.value_counts(), labels = ["Not Approved", "Approved"],
5        autopct = '%1.2f%%', startangle = 180, colors = ["#0EB8F1", "#F1480F"])
6
7  ax.set_title("Loan Approval")
8  plt.show()
```



Loan Approval

## UnBalanced data

## Handle UnBalanced Data Using Down Sampling

```python
[24]  1  import pandas as pd
      2
      3  # Combine X and y into a single DataFrame for easier sampling
      4  data = pd.concat([X, y], axis=1)
      5
      6  # Separate the classes
      7  minority_class = data[data['loan_status'] == 1]
      8  majority_class = data[data['loan_status'] == 0]
      9
     10  # Randomly sample from the majority class to match the minority class size
     11  majority_class_downsampled = majority_class.sample(n=len(minority_class),
         random_state=42)
     12
     13  # Concatenate the downsampled majority class with the minority class
     14  balanced_data = pd.concat([minority_class, majority_class_downsampled])
     15
     16  # Separate features and target variable
     17  X_resampled = balanced_data.drop('loan_status', axis=1)
     18  y_resampled = balanced_data['loan_status']
     19
```

```python
[25]  1  # Shape of resampled features
      2  print("Shape of X_resampled:", X_resampled.shape)
      3
      4  # Shape of resampled target
      5  print("Shape of y_resampled:", y_resampled.shape)
      6
```

```
Shape of X_resampled: (16700, 12)
Shape of y_resampled: (16700,)
```

```
1  import matplotlib.pyplot as plt
2  fig, ax = plt.subplots(figsize = (4, 4))
3
4  ax.pie(balanced_data.loan_status.value_counts(), labels = ["Not Approved",
   "Approved"],
5         autopct = '%1.2f%%', startangle = 180, colors = ["#0EB8F1", "#F1480F"])
6
7  ax.set_title("Loan Approval")
8  plt.show()
```

Loan Approval

Approved

50.00%

50.00%

Not Approved

## Checking Outliers

```python
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
from scipy import stats

# 1. Check for outliers using the IQR method
Q1 = df.describe().loc['25%']
Q3 = df.describe().loc['75%']
IQR = Q3 - Q1

# Define outlier bounds
lower_bound = Q1 - 1.5 * IQR
upper_bound = Q3 + 1.5 * IQR

# 2. Identify outliers
outliers = df[~((df >= lower_bound) & (df <= upper_bound)).all(axis=1)]

# 3. Plotting
plt.figure(figsize=(15, 10))

# Create box plots for numeric columns to visualize outliers
numeric_columns = df.select_dtypes(include=[np.number]).columns.tolist()

for col in numeric_columns:
    plt.figure(figsize=(10, 5))
    sns.boxplot(x=df[col])
    plt.axvline(x=lower_bound[col], color='red', linestyle='--', label='Lower
Bound')
```
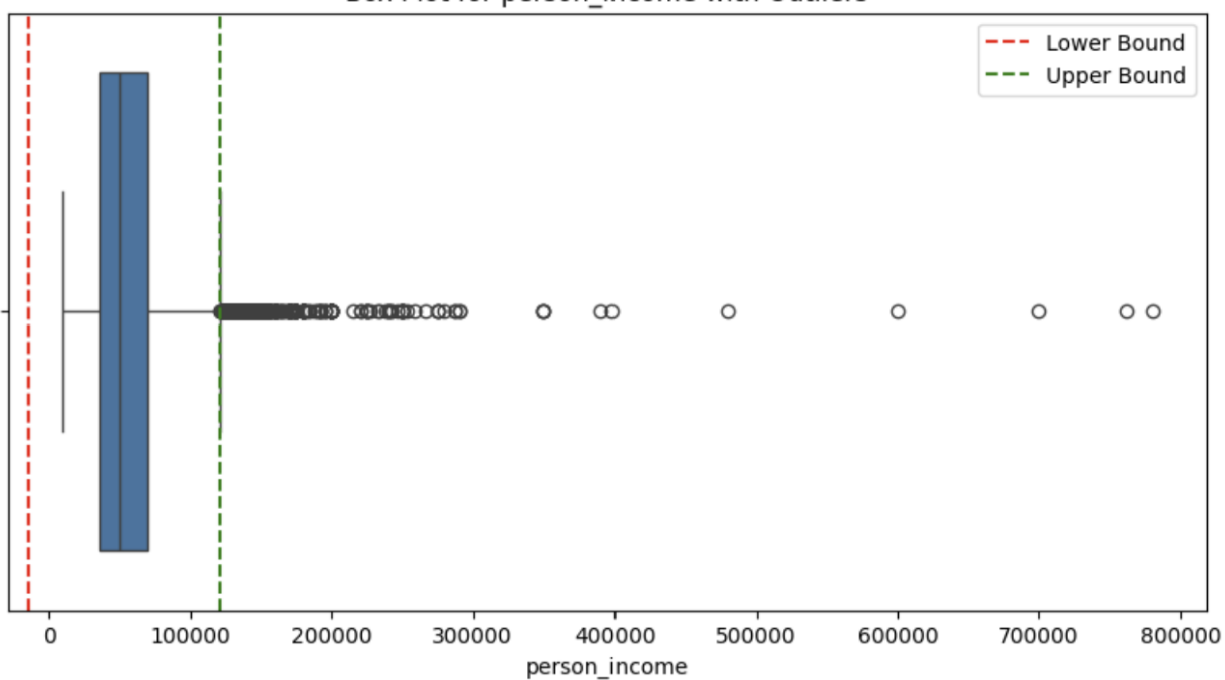
```python
    plt.axvline(x=upper_bound[col], color='green', linestyle='--', label='Upper
Bound')
    plt.title(f'Box Plot for {col} with Outliers')
    plt.xlabel(col)
    plt.legend()
    plt.show()

# Optional: Scatter plot for two specific variables to visualize relationships
and outliers
plt.figure(figsize=(10, 5))
sns.scatterplot(data=df, x='person_income', y='loan_amnt', color='blue',
label='Data Points')
sns.scatterplot(data=outliers, x='person_income', y='loan_amnt', color='red',
label='Outliers')
plt.title('Scatter Plot of Person Income vs Loan Amount')
plt.xlabel('Person Income')
plt.ylabel('Loan Amount')
plt.legend()
plt.show()
```
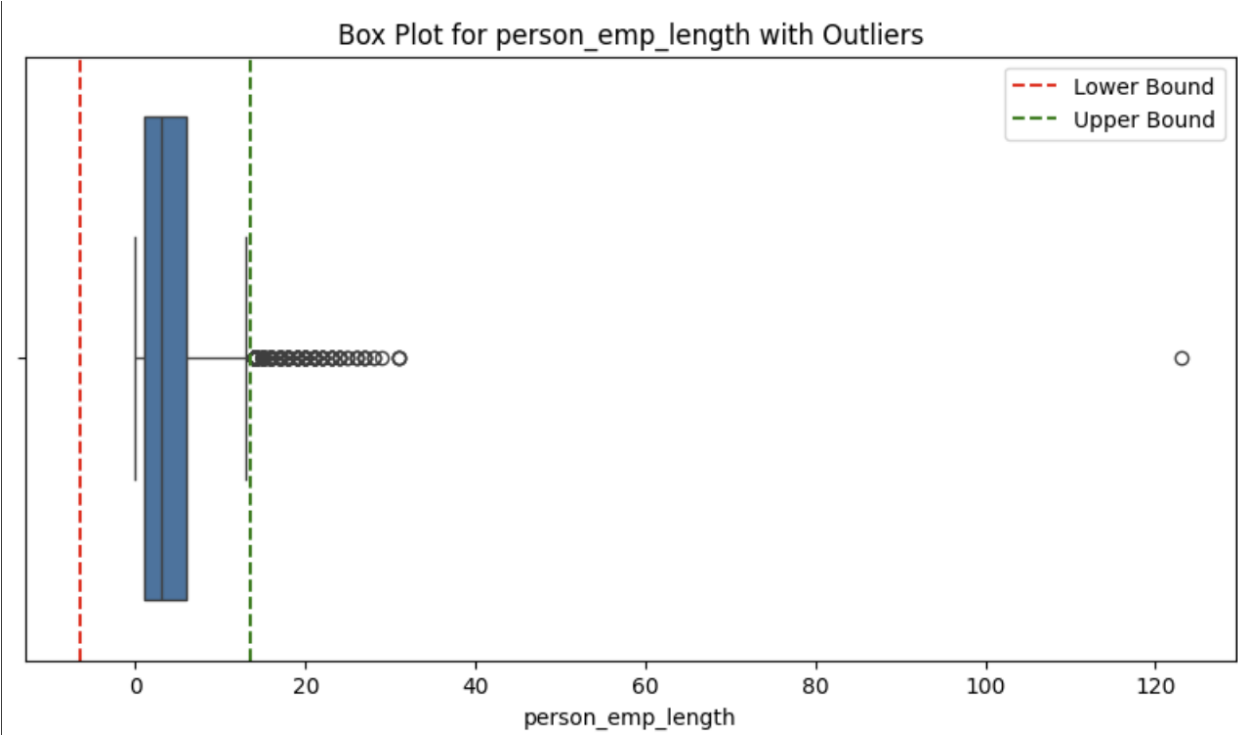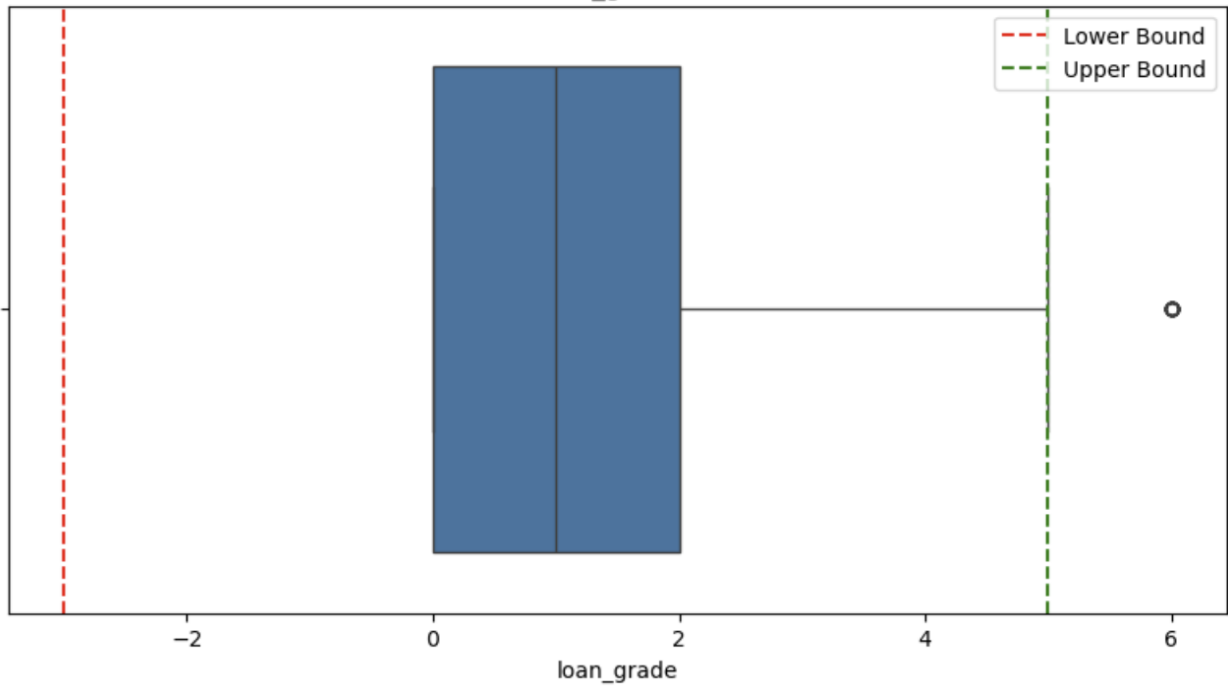

Box Plot for person_age with Outliers

Box Plot for person_income with Outliers
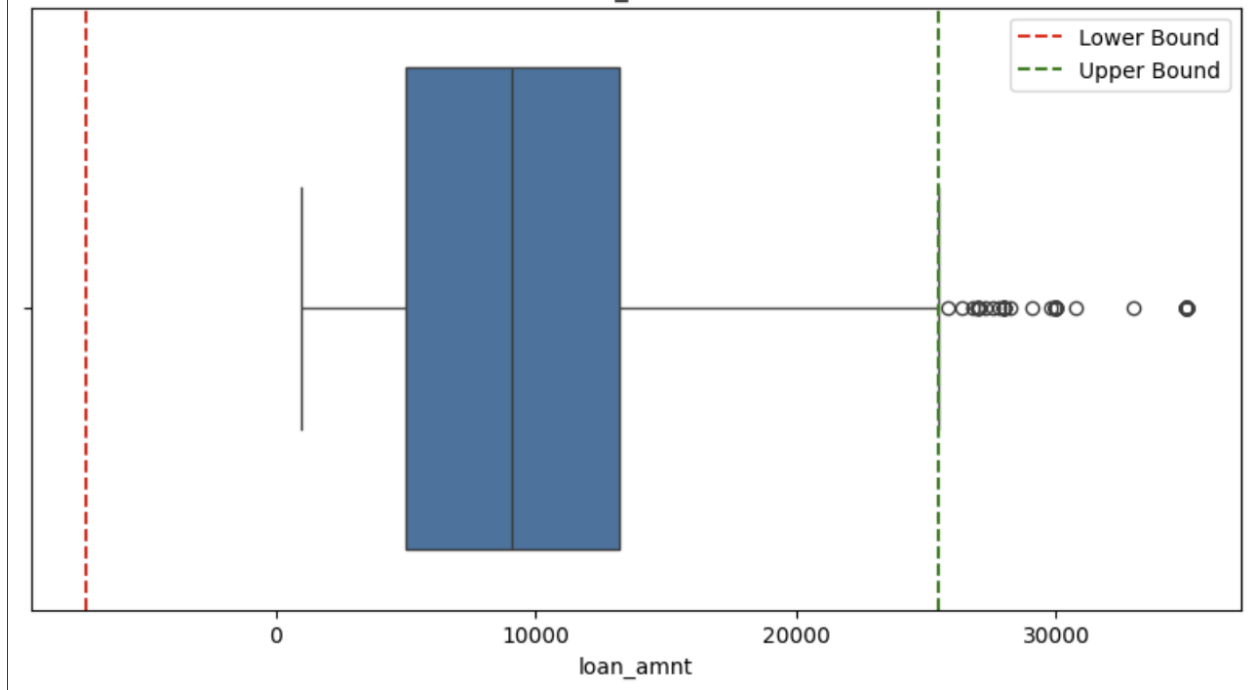


Box Plot for person_home_ownership with Outliers

Box Plot for person_emp_length with Outliers
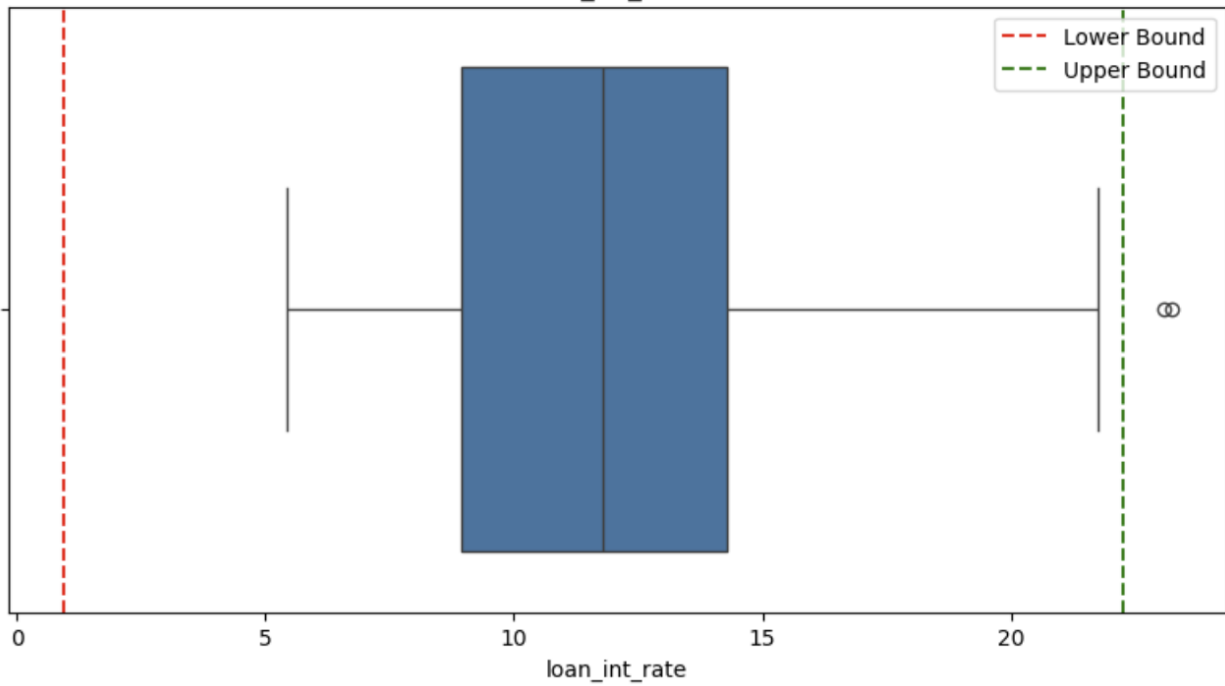

Box Plot for loan_intent with Outliers
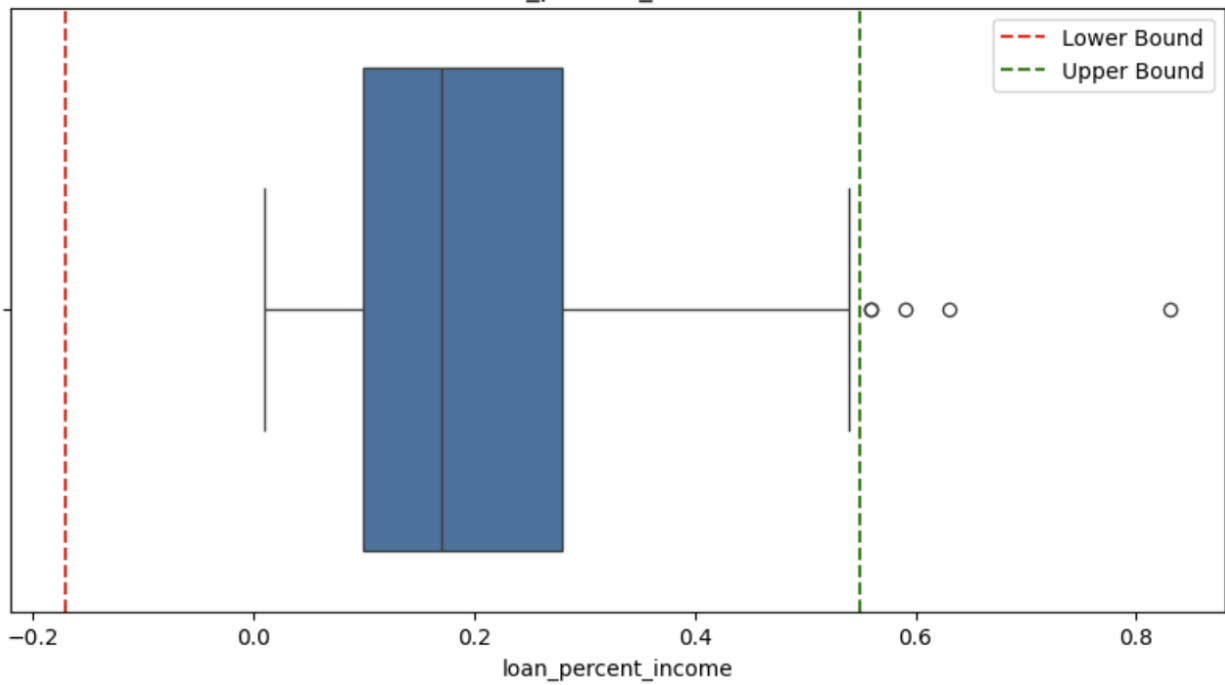
Box Plot for loan_grade with Outliers
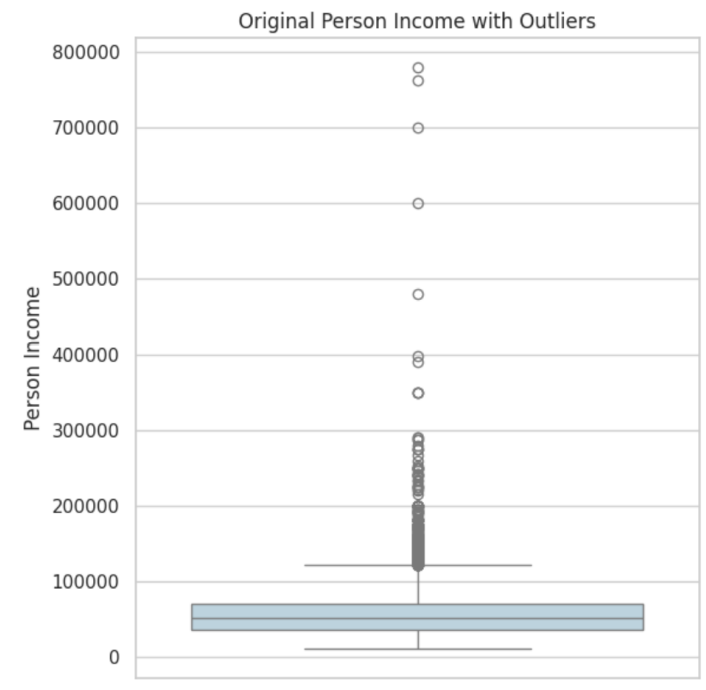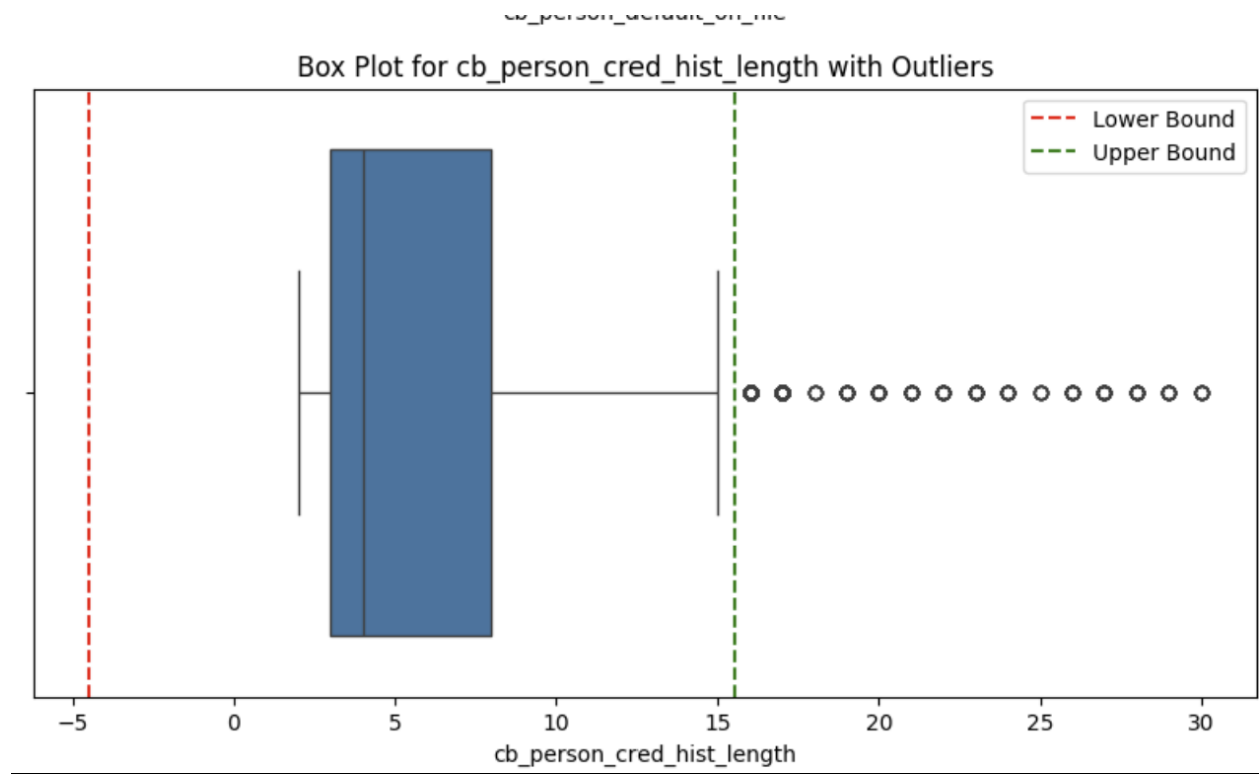


Box Plot for loan_amnt with Outliers

Box Plot for loan_int_rate with Outliers



Box Plot for loan_percent_income with Outliers

Box Plot for cb_person_cred_hist_length with Outliers
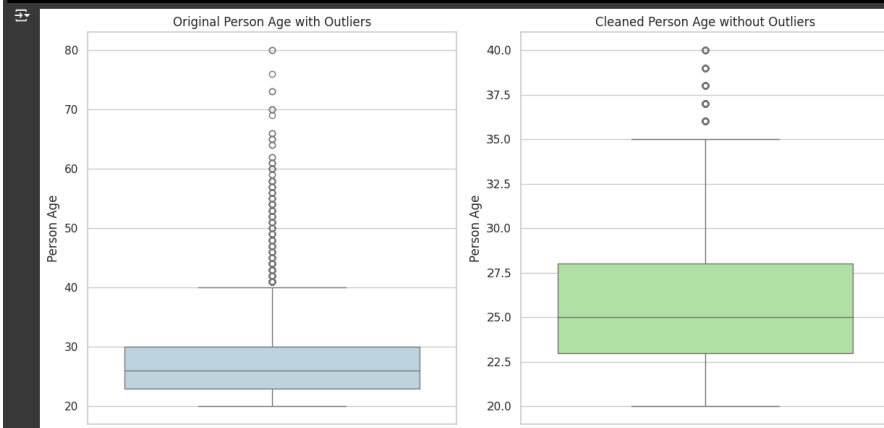


Original Person Income with Outliers

## Remove Outliers

```python
import pandas as pd
import numpy as np

# Assuming df is your DataFrame
# 1. Check for outliers using the IQR method
Q1 = df.describe().loc['25%']
Q3 = df.describe().loc['75%']
IQR = Q3 - Q1

# Define outlier bounds
lower_bound = Q1 - 1.5 * IQR
upper_bound = Q3 + 1.5 * IQR

# 2. Remove outliers
# Keep only the rows that are not outliers
df_cleaned = df[~((df < lower_bound) | (df > upper_bound)).any(axis=1)]

# Display the shape of the original and cleaned DataFrame
print(f"Original shape: {df.shape}")
print(f"Cleaned shape: {df_cleaned.shape}")
```

```
Original shape: (16700, 12)
Cleaned shape: (11540, 12)
```

```python
plt.figure(figsize=(12, 6))

# Box plot for original DataFrame for person_age
plt.subplot(1, 2, 1)
sns.boxplot(y=df['person_age'], color='lightblue')
plt.title("Original Person Age with Outliers")
plt.ylabel("Person Age")

# Box plot for cleaned DataFrame for person_age
plt.subplot(1, 2, 2)
sns.boxplot(y=df_cleaned['person_age'], color='lightgreen')
plt.title("Cleaned Person Age without Outliers")
plt.ylabel("Person Age")

plt.tight_layout()
plt.show()
```



- The left plot ("Original Person Age with Outliers") shows the data before removing outliers, with several data points outside the typical range, especially above 40 years.
- The right plot ("Cleaned Person Age without Outliers") displays the data after outlier removal, resulting in a tighter interquartile range (IQR) and fewer extreme values.

The right plot should ideally have removed the most significant outliers while retaining the core distribution of the data. If this is what you intended, then the visualization effectively demonstrates the cleaned distribution of ages.

```python
import matplotlib.pyplot as plt
import seaborn as sns

# Set the style for the plots
sns.set(style="whitegrid")

# Create a figure
plt.figure(figsize=(12, 6))

# Box plot for original DataFrame for person_income
plt.subplot(1, 2, 1)
sns.boxplot(y=df['person_income'], color='lightblue')
plt.title("Original Person Income with Outliers")
plt.ylabel("Person Income")

# Box plot for cleaned DataFrame for person_income
plt.subplot(1, 2, 2)
sns.boxplot(y=df_cleaned['person_income'], color='lightgreen')
plt.title("Cleaned Person Income without Outliers")
plt.ylabel("Person Income")

# Show the plots
plt.tight_layout()
plt.show()
```



```python
plt.figure(figsize=(12, 6))

# Box plot for original DataFrame for person_emp_length
plt.subplot(1, 2, 1)
sns.boxplot(y=df['person_emp_length'], color='lightblue')
plt.title("Original Person Employment Length with Outliers")
plt.ylabel("Person Employment Length")

# Box plot for cleaned DataFrame for person_emp_length
plt.subplot(1, 2, 2)
sns.boxplot(y=df_cleaned['person_emp_length'], color='lightgreen')
plt.title("Cleaned Person Employment Length without Outliers")
plt.ylabel("Person Employment Length")

plt.tight_layout()
plt.show()
```
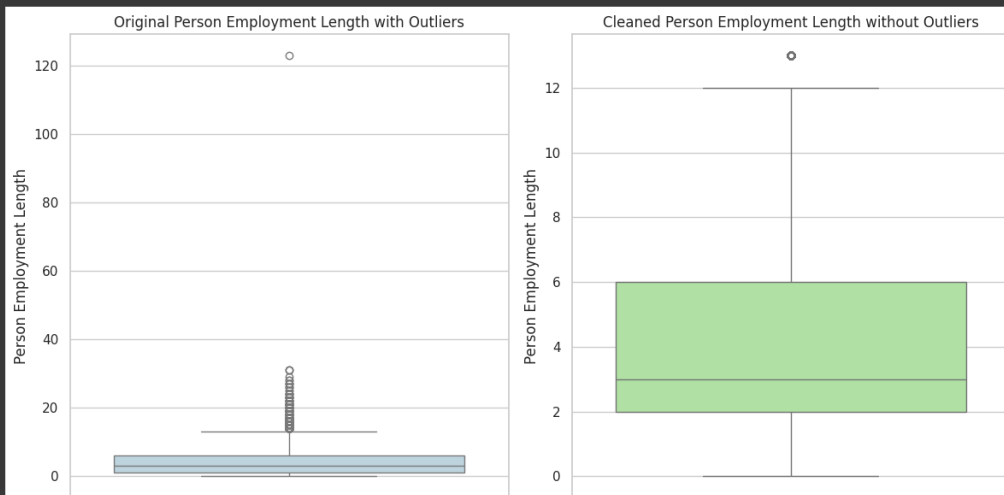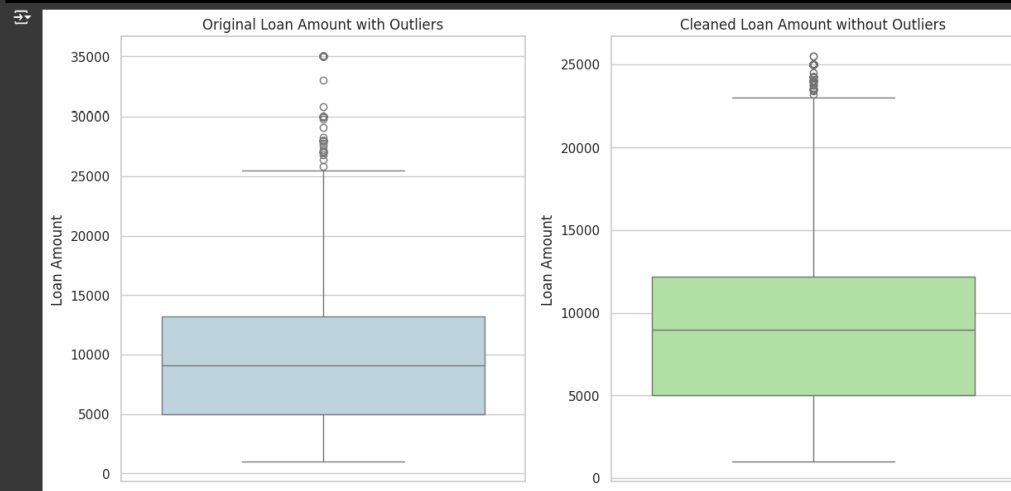
```
[56]  1  plt.figure(figsize=(12, 6))
      2
      3  # Box plot for original DataFrame for loan_amnt
      4  plt.subplot(1, 2, 1)
      5  sns.boxplot(y=df['loan_amnt'], color='lightblue')
      6  plt.title("Original Loan Amount with Outliers")
      7  plt.ylabel("Loan Amount")
      8
      9  # Box plot for cleaned DataFrame for loan_amnt
     10  plt.subplot(1, 2, 2)
     11  sns.boxplot(y=df_cleaned['loan_amnt'], color='lightgreen')
     12  plt.title("Cleaned Loan Amount without Outliers")
     13  plt.ylabel("Loan Amount")
     14
     15  plt.tight_layout()
     16  plt.show()
     17
```



```
[57]  1  plt.figure(figsize=(12, 6))
      2
      3  # Box plot for original DataFrame for loan_int_rate
      4  plt.subplot(1, 2, 1)
      5  sns.boxplot(y=df['loan_int_rate'], color='lightblue')
      6  plt.title("Original Loan Interest Rate with Outliers")
      7  plt.ylabel("Loan Interest Rate")
      8
      9  # Box plot for cleaned DataFrame for loan_int_rate
     10  plt.subplot(1, 2, 2)
     11  sns.boxplot(y=df_cleaned['loan_int_rate'], color='lightgreen')
     12  plt.title("Cleaned Loan Interest Rate without Outliers")
     13  plt.ylabel("Loan Interest Rate")
     14
     15  plt.tight_layout()
     16  plt.show()
     17
```
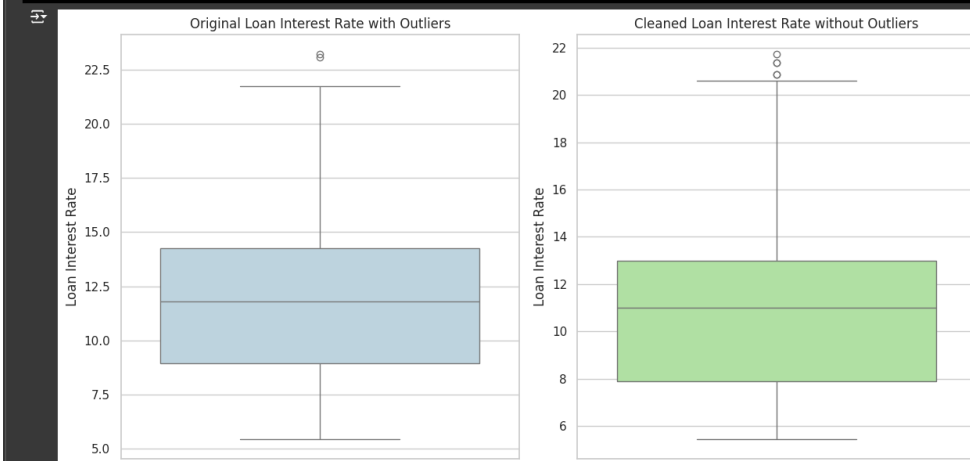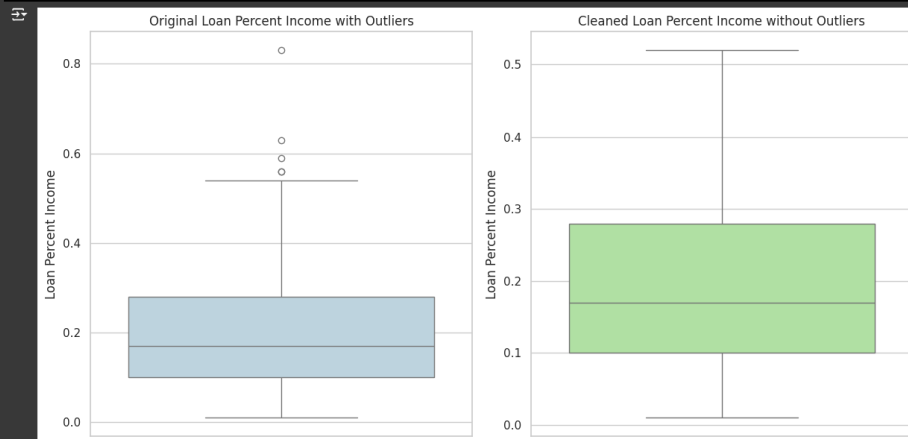
```
[58]   1  plt.figure(figsize=(12, 6))
       2
       3  # Box plot for original DataFrame for loan_percent_income
       4  plt.subplot(1, 2, 1)
       5  sns.boxplot(y=df['loan_percent_income'], color='lightblue')
       6  plt.title("Original Loan Percent Income with Outliers")
       7  plt.ylabel("Loan Percent Income")
       8
       9  # Box plot for cleaned DataFrame for loan_percent_income
      10  plt.subplot(1, 2, 2)
      11  sns.boxplot(y=df_cleaned['loan_percent_income'], color='lightgreen')
      12  plt.title("Cleaned Loan Percent Income without Outliers")
      13  plt.ylabel("Loan Percent Income")
      14
      15  plt.tight_layout()
      16  plt.show()
      17
```



```
[59]   1  plt.figure(figsize=(12, 6))
       2
       3  # Box plot for original DataFrame for cb_person_cred_hist_length
       4  plt.subplot(1, 2, 1)
       5  sns.boxplot(y=df['cb_person_cred_hist_length'], color='lightblue')
       6  plt.title("Original Credit History Length with Outliers")
       7  plt.ylabel("Credit History Length")
       8
       9  # Box plot for cleaned DataFrame for cb_person_cred_hist_length
      10  plt.subplot(1, 2, 2)
      11  sns.boxplot(y=df_cleaned['cb_person_cred_hist_length'], color='lightgreen')
      12  plt.title("Cleaned Credit History Length without Outliers")
      13  plt.ylabel("Credit History Length")
      14
      15  plt.tight_layout()
      16  plt.show()
      17
```
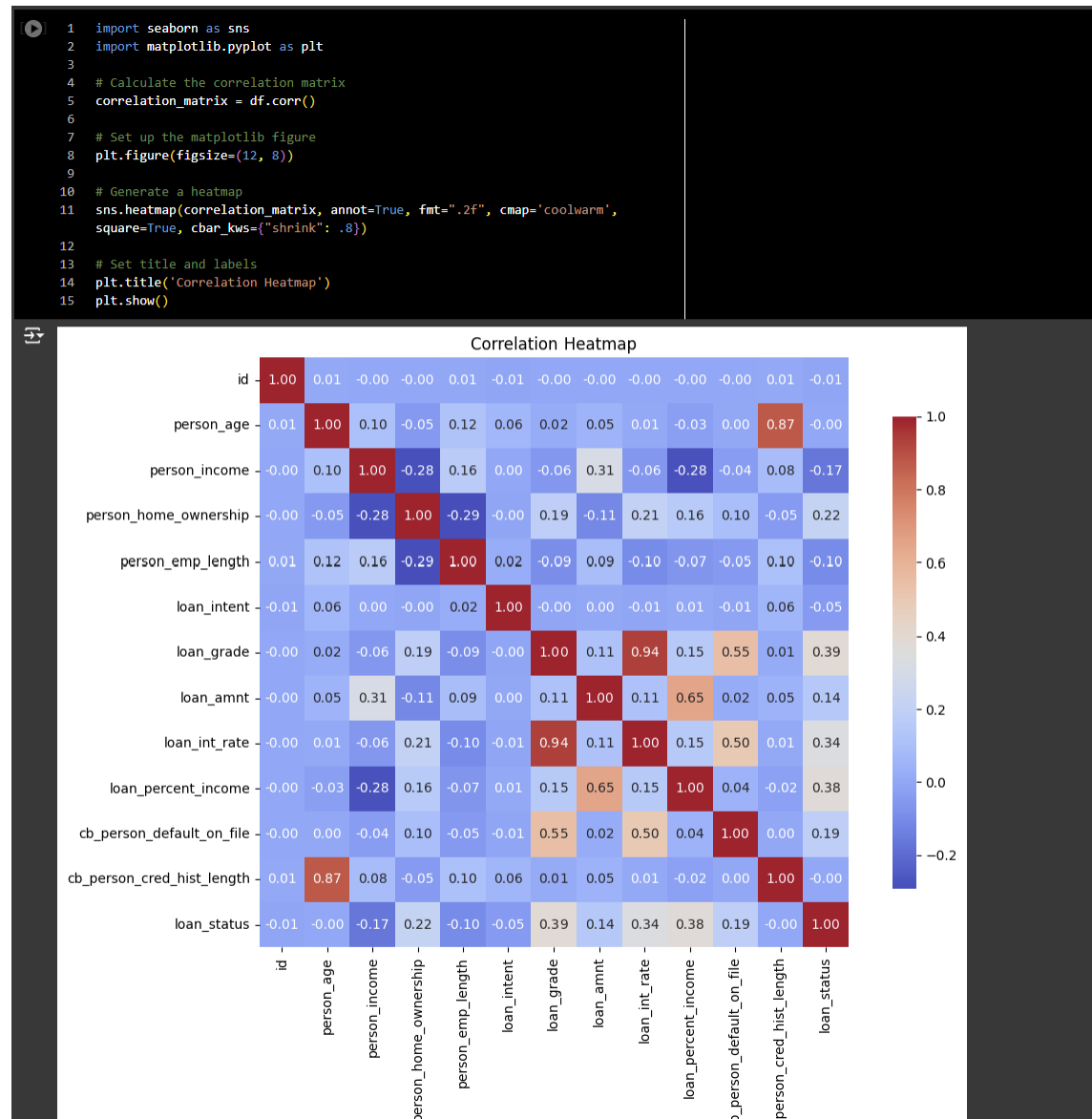
Divide the dataset into target (X) and features (y)

```python
from sklearn.model_selection import train_test_split

X = data_to_train.drop('loan_status', axis=1)  # Replace 'target_column' with
your actual target column name
y = data_to_train['loan_status']  # Replace 'target_column' with your actual
target column name

# Split the dataset into training and testing sets (70% train, 30% test)
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.3,
random_state=42)
```

# **Feature Selection**

```python
import seaborn as sns
import matplotlib.pyplot as plt

# Calculate the correlation matrix
correlation_matrix = df.corr()

# Set up the matplotlib figure
plt.figure(figsize=(12, 8))

# Generate a heatmap
sns.heatmap(correlation_matrix, annot=True, fmt=".2f", cmap='coolwarm',
square=True, cbar_kws={"shrink": .8})

# Set title and labels
plt.title('Correlation Heatmap')
plt.show()
```
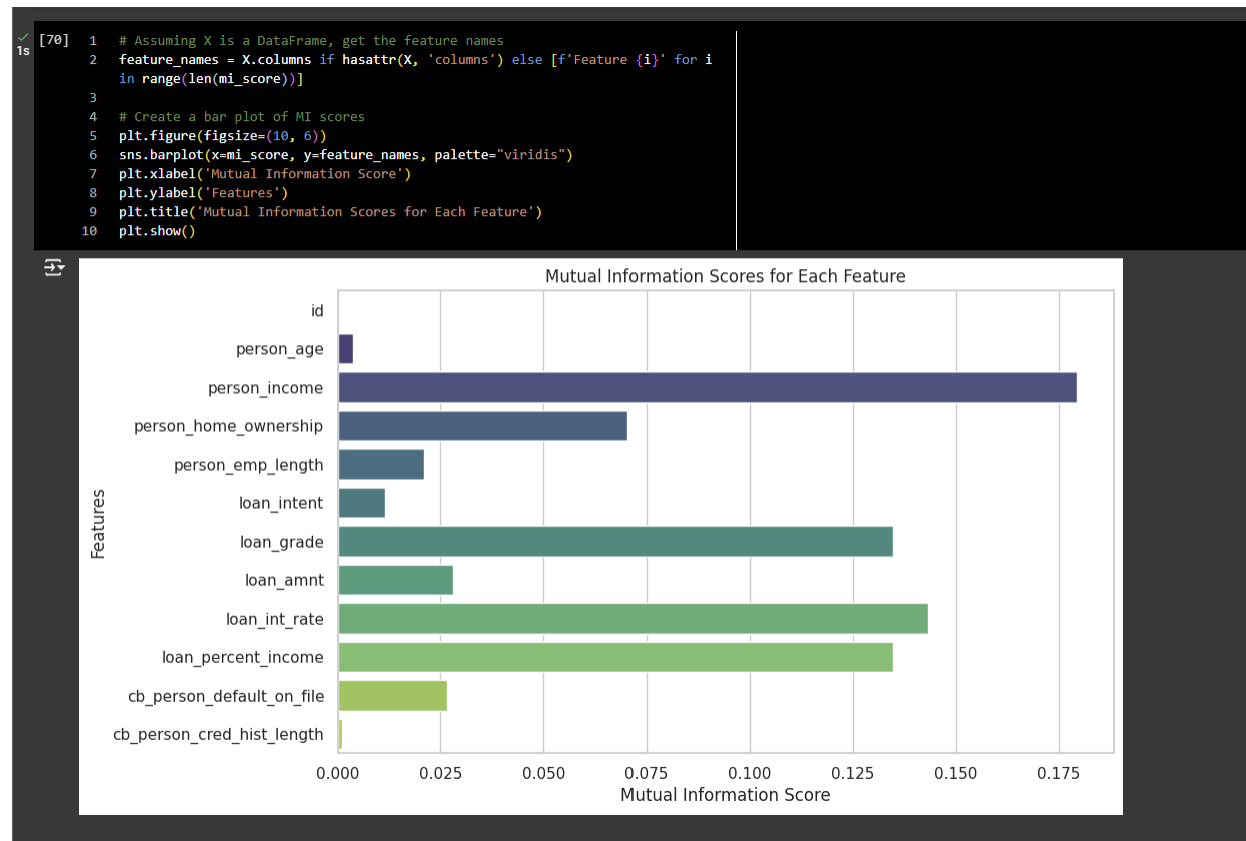


Correlation Heatmap

# Miscore

```
1  from sklearn.feature_selection import mutual_info_classif,
   mutual_info_regression
2
3  # Use mutual_info_classif for classification tasks or mutual_info_regression
   for regression tasks
4  # Replace X and y with your actual feature matrix and target variable
5
6  # For classification tasks
7  mi_score = mutual_info_classif(X, y, random_state=0)
8
9  # For regression tasks
10 # mi_score = mutual_info_regression(X, y, random_state=0)
11
12 # Display MI scores
13 print(mi_score)
14
```

```
[0.         0.00373506 0.17936325 0.07044917 0.02101987 0.01156984
 0.13482432 0.02798012 0.14331738 0.13466079 0.02672561 0.00098097]
```

```
[70] 1  # Assuming X is a DataFrame, get the feature names
     2  feature_names = X.columns if hasattr(X, 'columns') else [f'Feature {i}' for i
        in range(len(mi_score))]
     3
     4  # Create a bar plot of MI scores
     5  plt.figure(figsize=(10, 6))
     6  sns.barplot(x=mi_score, y=feature_names, palette="viridis")
     7  plt.xlabel('Mutual Information Score')
     8  plt.ylabel('Features')
     9  plt.title('Mutual Information Scores for Each Feature')
     10 plt.show()
```

# Correlation with Target



Correlation Matrix with Target Variable (Loan Status)