```
1 # Task 1: Data Loading & Cleaning
2
3 import pandas as pd
4 import numpy as np
5 import matplotlib.pyplot as plt
6 import seaborn as sns
7
```

```
 1
 2 # Load the Dataset:
 3
 4 file_path = "mini_lending_club.csv"
 5 df = pd.read_csv(file_path)
 6
 7 # Preview the first rows
 8 print("Preview of dataset:")
 9 display(df.head())
10
11 print("\nDataset shape:", df.shape)
```

Preview of dataset:

|   | loan_amnt | term | int_rate | installment | grade | sub_grade | emp_length | home_ownership | annual_inc | purpose | dti | de |
|---|-----------|------|----------|-------------|-------|-----------|------------|----------------|------------|---------|-----|-----|
| 0 | 16795 | 36 months | 23.39% | 55.32 | G | A5 | 2 years | RENT | 33704 | home_improvement | 12.11 | |
| 1 | 1860 | 60 months | 29.96% | 664.55 | F | C1 | 6 years | OWN | 78933 | debt_consolidation | 1.50 | |
| 2 | 12284 | 36 months | 28.33% | 319.87 | E | C1 | 10+ years | OWN | 91101 | home_improvement | 25.18 | |
| 3 | 7265 | 36 months | 21.06% | 455.57 | E | B3 | 10+ years | RENT | 63856 | credit_card | 33.34 | |
| 4 | 17850 | 36 months | 15.53% | 557.98 | D | A1 | 8 years | OWN | 57256 | other | 9.24 | |

Dataset shape: (2000, 16)

```
 1 # Inspect Generic Info About the Dataset
 2 # --------------------------------------
 3
 4 print("\nDataset Info:")
 5 df.info()
 6
 7 # Summary statistics for numerical features
 8 print("\nSummary Statistics:")
 9 display(df.describe())
```

```
Dataset Info:
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 2000 entries, 0 to 1999
Data columns (total 16 columns):
 #   Column           Non-Null Count  Dtype
---  ------           --------------  -----
 0   loan_amnt        2000 non-null   int64
 1   term             2000 non-null   object
 2   int_rate         1900 non-null   object
 3   installment      2000 non-null   float64
 4   grade            2000 non-null   object
 5   sub_grade        2000 non-null   object
 6   emp_length       1900 non-null   object
 7   home_ownership   2000 non-null   object
 8   annual_inc       2000 non-null   int64
 9   purpose          1900 non-null   object
 10  dti              2000 non-null   float64
 11  delinq_2yrs      2000 non-null   int64
 12  revol_util       1900 non-null   object
 13  total_acc        2000 non-null   int64
 14  earliest_cr_line 2000 non-null   int64
 15  loan_status      2000 non-null   object
dtypes: float64(2), int64(5), object(9)
memory usage: 250.1+ KB
```

Summary Statistics:

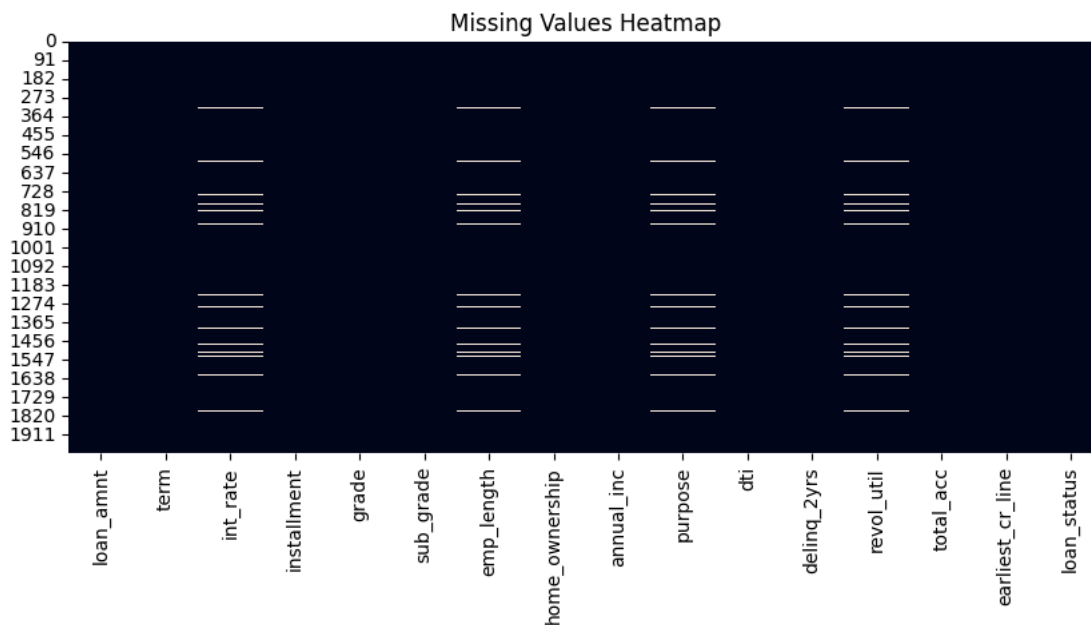|       | loan_amnt    | installment | annual_inc    | dti         | delinq_2yrs | total_acc   | earliest_cr_line |
|-------|--------------|-------------|---------------|-------------|-------------|-------------|------------------|
| count | 2000.000000  | 2000.000000 | 2000.000000   | 2000.000000 | 2000.000000 | 2000.000000 | 2000.000000      |
| mean  | 18031.749500 | 413.815520  | 85211.457000  | 17.889410   | 2.070000    | 40.250000   | 2006.874000      |
| std   | 9788.029854  | 213.172169  | 37139.980987  | 10.114553   | 1.441575    | 22.706041   | 5.899367         |
| min   | 1009.000000  | 50.180000   | 20057.000000  | 0.020000    | 0.000000    | 1.000000    | 1999.000000      |
| 25%   | 9485.500000  | 230.990000  | 52672.250000  | 9.222500    | 1.000000    | 20.000000   | 2001.000000      |
| 50%   | 18330.500000 | 411.110000  | 85122.500000  | 17.765000   | 2.000000    | 40.000000   | 2008.000000      |
| 75%   | 26567.250000 | 589.030000  | 117562.500000 | 26.510000   | 3.000000    | 60.000000   | 2012.000000      |
| max   | 34997.000000 | 799.180000  | 149986.000000 | 35.000000   | 4.000000    | 79.000000   | 2016.000000      |

```
 1  # 3. Check for missing values
 2  # -------------------------------------
 3  print("\nMissing Values:")
 4  display(df.isnull().sum())
 5
 6  # Visualize missing values
 7  plt.figure(figsize=(10,4))
 8  sns.heatmap(df.isnull(), cbar=False)
 9  plt.title("Missing Values Heatmap")
10  plt.show()
11
```

Missing Values:

|                 | 0   |
|-----------------|-----|
| loan_amnt       | 0   |
| term            | 0   |
| int_rate        | 100 |
| installment     | 0   |
| grade           | 0   |
| sub_grade       | 0   |
| emp_length      | 100 |
| home_ownership  | 0   |
| annual_inc      | 0   |
| purpose         | 100 |
| dti             | 0   |
| delinq_2yrs     | 0   |
| revol_util      | 100 |
| total_acc       | 0   |
| earliest_cr_line| 0   |
| loan_status     | 0   |

**dtype:** int64


Missing Values Heatmap

```
 1 # 4. Clean & preprocess raw string features (SAFE VERSION)
 2 # ------------------------------------------------------
 3
 4 # Ensure string conversion before using .str
 5 df['int_rate'] = df['int_rate'].astype(str).str.replace('%','', regex=False)
 6 df['revol_util'] = df['revol_util'].astype(str).str.replace('%','', regex=False)
 7
 8 # Convert to numeric (coerce errors turns bad values into NaN)
 9 df['int_rate'] = pd.to_numeric(df['int_rate'], errors='coerce')
10 df['revol_util'] = pd.to_numeric(df['revol_util'], errors='coerce')
11
12 # emp_length cleaning
13 df['emp_length'] = df['emp_length'].astype(str)
14
15 df['emp_length'] = (
16     df['emp_length']
17     .replace({
18         '< 1 year': '0',
```

```
19        '10+ years': '10'
20    })
21    .str.extract(r'(\d+)', expand=False)   # extract number
22 )
23
24 df['emp_length'] = pd.to_numeric(df['emp_length'], errors='coerce')
25
26 # earliest credit line
27 df['earliest_cr_line'] = pd.to_numeric(df['earliest_cr_line'], errors='coerce')
28
29
30 # Handle missing values
31 num_cols = df.select_dtypes(include=['float64','int64']).columns
32 df[num_cols] = df[num_cols].fillna(df[num_cols].median())
33
34 cat_cols = df.select_dtypes(include=['object']).columns
35 df[cat_cols] = df[cat_cols].fillna(df[cat_cols].mode().iloc[0])
36
37 # -------------------------------------
38 # 6. Convert loan_status to a binary target
39 # -------------------------------------
40
41 df['loan_default'] = df['loan_status'].apply(
42     lambda x: 1 if x == "Charged Off" else 0
43 )
44
45 print("\nLoan Status Conversion (value counts):")
46 print(df['loan_default'].value_counts())
47
48
```

```
Loan Status Conversion (value counts):
loan_default
0    1518
1     482
Name: count, dtype: int64
```

```
1 df.to_parquet("df_cleaned.parquet", index=False)
2
```

```
1 #Task 2 – FEATURE ENGINEERING
2 # --------------------------
3
4 from sklearn.preprocessing import StandardScaler
5 from sklearn.preprocessing import OneHotEncoder
6 from sklearn.compose import ColumnTransformer
7 from sklearn.pipeline import Pipeline
8
9
10 # These help the model understand borrower risk better.
11
12 df["income_to_loan"] = df["annual_inc"] / df["loan_amnt"]          # Ability to repay
13 df["credit_util_ratio"] = df["revol_util"] / 100                  # Revolving credit usage
14 df["dti_income_ratio"] = df["dti"] / (df["annual_inc"] + 1)       # DTI relative to income
15 df["credit_history_years"] = 2024 - df["earliest_cr_line"]        # Borrower's credit experience
16 df["installment_ratio"] = df["installment"] / (df["annual_inc"] + 1)   # Payment pressure
```

```
1 # SEPARATE FEATURES (X) AND TARGET (y)
2 # ---------------------------------------------
3
4 y = df["loan_default"]  # Target is the dependent variable. what I want to predict and it can't be a part of the data that
5
6 # Remove columns that should not be used in modeling, it could lead to data-leakage.
7 drop_cols = ["loan_status","loan_default","earliest_cr_line"]
8 X = df.drop(columns=drop_cols)
9
10 # IDENTIFY NUMERIC & CATEGORICAL COLUMNS
11 # ---------------------------------------------
12 numeric_features = X.select_dtypes(include=["float64", "int64"]).columns
13 categorical_features = X.select_dtypes(include=["object"]).columns
14
15 print("Numeric Columns:\n", list(numeric_features))
16 print("\nCategorical Columns:\n", list(categorical_features))
17
```

```
18 y.to_frame().to_parquet("y.parquet", index=False) # save target variable y
19
```

```
Numeric Columns:
 ['loan_amnt', 'int_rate', 'installment', 'emp_length', 'annual_inc', 'dti', 'delinq_2yrs', 'revol_util', 'total_acc', 'income_t

Categorical Columns:
 ['term', 'grade', 'sub_grade', 'home_ownership', 'purpose']
```

```
1 df.columns
2
```

```
Index(['loan_amnt', 'term', 'int_rate', 'installment', 'grade', 'sub_grade',
       'emp_length', 'home_ownership', 'annual_inc', 'purpose', 'dti',
       'delinq_2yrs', 'revol_util', 'total_acc', 'earliest_cr_line',
       'loan_status', 'income_to_loan', 'credit_util_ratio',
       'dti_income_ratio', 'credit_history_years', 'installment_ratio',
       'loan_default'],
      dtype='object')
```

```
 1 # SET UP TRANSFORMERS
 2 # -----------------------------------------------
 3
 4 # StandardScaler:
 5 # → Makes numeric features comparable by converting them to mean=0, variance=1.
 6
 7 numeric_transformer = StandardScaler()
 8
 9 # OneHotEncoder:
10 # → Converts categories (e.g., "RENT", "OWN") into 0/1 columns the model can understand.
11 categorical_transformer = OneHotEncoder(handle_unknown='ignore')
12
13 # ColumnTransformer:
14 # → Applies transformations to different column groups:
15 #      - scale numeric features
16 #      - one-hot encode categorical features
17 # → Combines everything into ONE clean matrix
18 # In plain sense I am using the in-built libr ColumnTransf to consolidate both col groups (numericla and Category into one
19
20 preprocessor = ColumnTransformer(
21     transformers=[
22         ("num", numeric_transformer, numeric_features),  # scale numerics
23         ("cat", categorical_transformer, categorical_features)  # encode categoricals
24     ]
25 )
26
27 # -----------------------------------------------
28 #  APPLY TRANSFORMATIONS
29 # -----------------------------------------------
30 # fit_transform:
31 # → Learns how to scale from data
32 # → Creates a transformed feature matrix ready for modelling
33
34 X_prepared = preprocessor.fit_transform(X)
35
36 # ===============================
37 # FIX FEATURE NAMES (IMPORTANT FOR SHAP)
38 # ===============================
39
40 # 1. Numeric feature names
41 numeric_feature_names = numeric_features.tolist()
42
43 # 2. One-hot encoded categorical feature names
44 categorical_feature_names = (
45     preprocessor.named_transformers_['cat']
46     .get_feature_names_out(categorical_features)
47     .tolist()
48 )
49
50 # 3. Combine them into one final list
51 all_feature_names = numeric_feature_names + categorical_feature_names
52
53 print("Number of feature names:", len(all_feature_names)) # the len indicate the number of features
54
55 import scipy.sparse # for saving the file jsut i case it crashes
56 scipy.sparse.save_npz("X_prepared.npz", X_prepared)
```
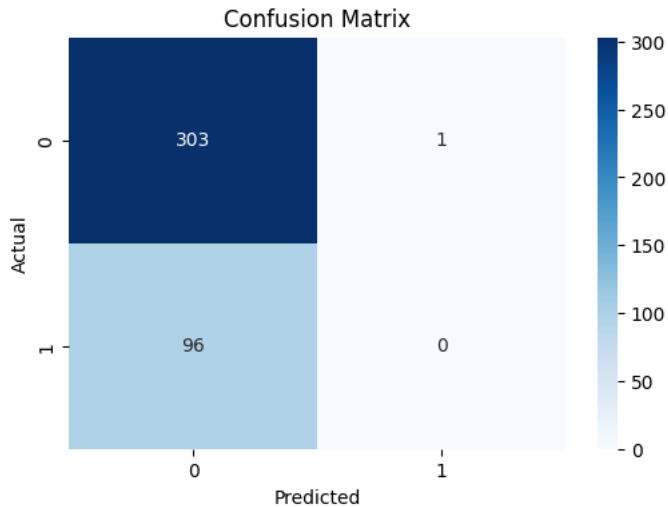
Number of feature names: 66

```python
1  # Task 3: MODEL TRAINING & EVALUATION
2  # =====================================
3
4  # Import Suitable Libraries
5
6  from sklearn.model_selection import train_test_split
7  from sklearn.ensemble import RandomForestClassifier
8  from sklearn.metrics import accuracy_score, confusion_matrix, classification_report
9
10 # 1. TRAIN-TEST SPLIT
11 # ----------------------------------------
12
13 # X_prepared = transformed feature matrix from Task 2 (pipeline output)
14 # y = target variable from Task 2
15
16 X_train, X_test, y_train, y_test = train_test_split(
17     X_prepared, # this replaces the X value for training
18     y, # this represents the target value
19     test_size=0.2,
20     random_state=42,
21     stratify=y  # ensures balanced target distribution
22 )
23
24
25 # TRAIN RANDOM FOREST MODEL
26 # ----------------------------------------
27
28 # Model Selection: I am using RandomForest to determine the prediction for the loan_default prediction
29
30 random_forex = RandomForestClassifier(
31     n_estimators=200,
32     max_depth=None,
33     random_state=42
34 )
35
36 random_forex .fit(X_train, y_train)
37
38 # 3. MAKE PREDICTIONS
39 # ----------------------------------------
40
41 y_pred = random_forex.predict(X_test) # - this is the prediction of the model
42
43
44 # 4. MODEL EVALUATION - this is the process of evaluating that the model has a higher accuracy rate or not.
45 # ----------------------------------------
46
47 print("Accuracy:", accuracy_score(y_test, y_pred))
48 print("\nClassification Report:\n", classification_report(y_test, y_pred))
49
50 # Confusion Matrix
51 cm = confusion_matrix(y_test, y_pred)
52 plt.figure(figsize=(6,4))
53 sns.heatmap(cm, annot=True, fmt="d", cmap="Blues")
54 plt.title("Confusion Matrix")
55 plt.xlabel("Predicted")
56 plt.ylabel("Actual")
57 plt.show()
58
```

```
Accuracy: 0.7575

Classification Report:
              precision    recall  f1-score   support

           0       0.76      1.00      0.86       304
           1       0.00      0.00      0.00        96

    accuracy                           0.76       400
   macro avg       0.38      0.50      0.43       400
weighted avg       0.58      0.76      0.66       400
```
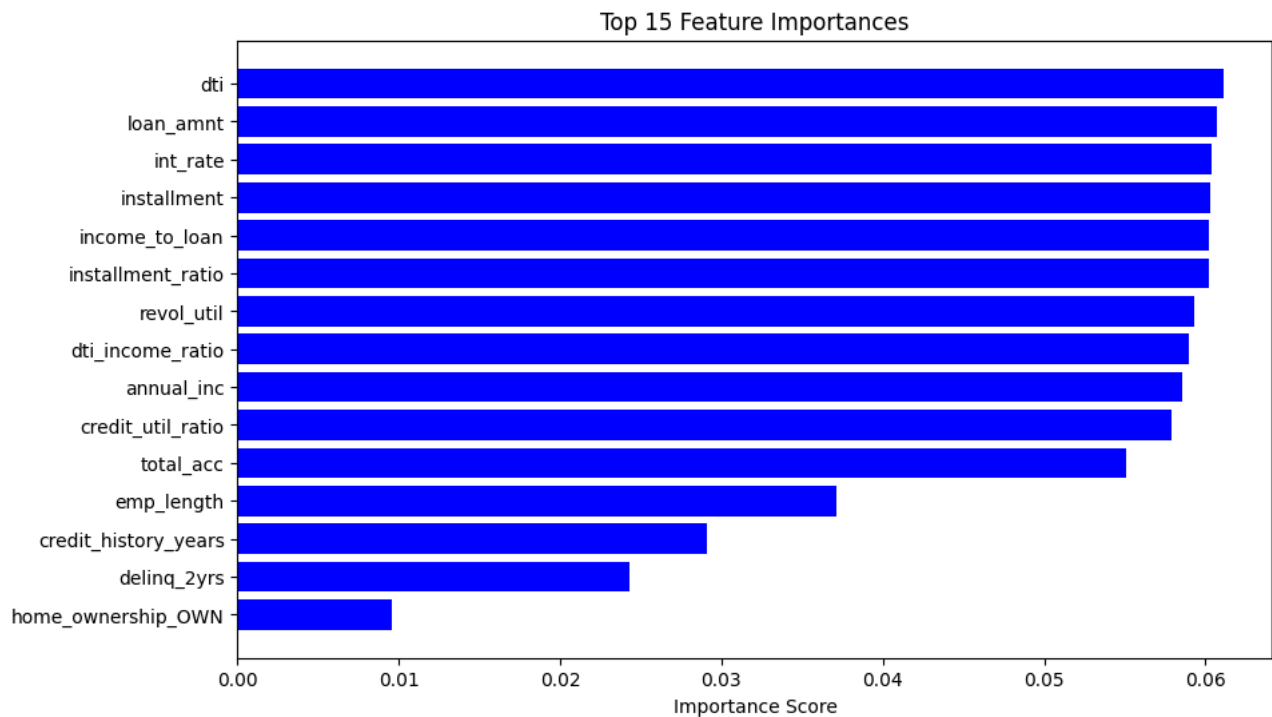


Confusion Matrix

```
1 y.to_frame().to_parquet("y.parquet", index=False)
2
```

```
 1 # 5. FEATURE IMPORTANCE VISUALIZATION
 2 # ----------------------------------------
 3
 4 # Grab feature names from the preprocessor
 5 # (numerical + encoded categorical names)
 6 num_features = numeric_features
 7 cat_features = list(preprocessor.named_transformers_['cat'].get_feature_names_out(categorical_features))
 8
 9 all_feature_names = list(num_features) + cat_features
10
11 importances = random_forex.feature_importances_
12
13 # Sort for visualization
14 indices = np.argsort(importances)[-15:]  # top 15 features
15
16 plt.figure(figsize=(10,6))
17 plt.barh(range(len(indices)), importances[indices], color='blue')
18 plt.yticks(range(len(indices)), [all_feature_names[i] for i in indices])
19 plt.title("Top 15 Feature Importances")
20 plt.xlabel("Importance Score")
21 plt.show()
```

Top 15 Feature Importances

```
 1 # =================================
 2 # TASK 4 — SHAP & FAIRNESS ANALYSIS
 3 # =================================
 4 import shap
 5 shap.initjs()
 6
 7 # ----------------------------------------
 8 # 1. PREP SHAP INPUT (convert sparse → dense)
 9 # I had an issue with this section and request the assistance of Google Gemini
10 # ----------------------------------------
11
12 X_test_sample = X_test[:300]   # smaller sample for speed
13
14 # Convert sparse matrix to dense
15 if hasattr(X_test_sample, "toarray"):
16     X_test_sample_dense = X_test_sample.toarray()
17 else:
18     X_test_sample_dense = np.array(X_test_sample)
19
20 # ----------------------------------------
21 # 2. SHAP EXPLAINER FOR RANDOM FOREST
22 # ----------------------------------------
23
24 # The error indicates a mismatch in the number of features between shap_values
25 # and X_test_sample_dense. This often happens if the explainer is not
26 # perfectly aligned with the data it's explaining, or if the model's
27 # internal feature representation differs.
28 # To try and align them, we explicitly pass the data to the explainer.
29
30 explainer = shap.TreeExplainer(random_forex, data=X_test_sample_dense)
31
32 # RandomForest produces 2 classes → index 1 is "default"
33
34 shap_values = explainer.shap_values(X_test_sample_dense)
35
36 # Diagnostic: Print shapes before plotting to confirm the issue
37 print(f"Shape of X_test_sample_dense: {X_test_sample_dense.shape}")
38
39 # shap_values can be a list of arrays or a 3D array; handle both
40
41 if isinstance(shap_values, list):
42     print(f"Shape of shap_values (list of arrays): {[s.shape for s in shap_values]}")
43     if len(shap_values) > 1:
44         print(f"Shape of shap_values[1]: {shap_values[1].shape}")
45 else: # assuming it's a numpy array
46     print(f"Shape of shap_values (3D array): {shap_values.shape}")
```

```
47      # The previous line shap_values[1] was incorrect for a 3D array; it picked a sample, not a class.
48      # The correct way to get class 1's shap values for all samples is shap_values[:, :, 1]
49      print(f"Shape of shap_values[:, :, 1]: {shap_values[:, :, 1].shape}")
50 print(f"Length of all_feature_names: {len(all_feature_names)}")
51
52
53 # ----------------------------------------
54 # 3. GLOBAL SHAP SUMMARY PLOTS
55 # ----------------------------------------
56
57 # Beeswarm plot — feature impact
58 shap.summary_plot(
59      shap_values[:, :, 1],         # Correctly select SHAP values for class 1 (default) across all samples
60      X_test_sample_dense,
61      feature_names=all_feature_names
62 )
63
64 # Clean bar plot — easy for reports
65 shap.summary_plot(
66      shap_values[:, :, 1],
67      X_test_sample_dense,
68      feature_names=all_feature_names,
69      plot_type="bar"
70 )
71
72 # ----------------------------------------
73 # 4. LOCAL EXPLANATION (single borrower)
74 # ----------------------------------------
75
76 idx = 0
77 shap.force_plot(
78      explainer.expected_value[1],
79      shap_values[idx, :, 1], # Correctly select SHAP values for a single sample (idx) and class 1
80      X_test_sample_dense[idx],
81      feature_names=all_feature_names
82 )
83
84 # ----------------------------------------
85 # 5. FAIRNESS ANALYSIS — GROUP ACCURACY
86 # ----------------------------------------
87
88 # Recreate original X_test (non-transformed) for grouping
89 from sklearn.model_selection import train_test_split
90
91 X_train_df, X_test_df, y_train_df, y_test_df = train_test_split(
92      X, y, test_size=0.2, random_state=42, stratify=y
93 )
94
95 # Attach predictions to original test data
96 fairness_df = X_test_df.copy()
97 fairness_df["y_true"] = y_test_df.values
98 fairness_df["y_pred"] = random_forex.predict(X_test)
99
100 # Choose grouping variable
101 group_col = "home_ownership"
102
103 print("Groups:", fairness_df[group_col].unique())
104
105 group_metrics = {}
106 for group in fairness_df[group_col].unique():
107      sub = fairness_df[fairness_df[group_col] == group]
108      acc = (sub["y_true"] == sub["y_pred"]).mean()
109      group_metrics[group] = acc
110      print(f"Accuracy for {group}: {acc:.3f}")
111
112 # Bar chart for fairness visualization
113 plt.figure(figsize=(6,4))
114 plt.bar(group_metrics.keys(), group_metrics.values(), color="steelblue")
115 plt.title("Model Accuracy by Home Ownership Group")
116 plt.xlabel("Group")
117 plt.ylabel("Accuracy")
118 plt.ylim(0, 1)
119 plt.show()
```