


```

1 import pandas as pd
2 import numpy as np
3
4 import seaborn as sns
5 import matplotlib.pyplot as plt
6 import plotly.express as px
7
8 from sklearn.preprocessing import StandardScaler
9 from sklearn.preprocessing import LabelEncoder
10 from sklearn.preprocessing import OneHotEncoder
11 from sklearn.compose import ColumnTransformer
12 from sklearn.model_selection import train_test_split
13 from yellowbrick.classifier import ConfusionMatrix
14 from sklearn.metrics import accuracy_score, classification_report

```

```
1 df=pd.read_csv("/content/credit_risk_dataset.csv")
```

```
1 df.head()
```



	person_age	person_income	person_home_ownership	person_emp_length	loan_intent	loan_grade	loan_amnt	loan_int_rate	loan_status
0	22	59000	RENT	123.0	PERSONAL	D	35000	16.02	PAID
1	21	9600	OWN	5.0	EDUCATION	B	1000	11.14	PAID
2	25	9600	MORTGAGE	1.0	MEDICAL	C	5500	12.87	PAID
3	23	65500	RENT	4.0	MEDICAL	C	35000	15.23	PAID
4	24	54400	RENT	8.0	MEDICAL	C	35000	14.27	PAID

```
1 df.describe()
```

	person_age	person_income	person_emp_length	loan_amnt	loan_int_rate	loan_status	loan_percent_income	cb_person_c
count	32581.000000	3.258100e+04	31686.000000	32581.000000	29465.000000	32581.000000	32581.000000	
mean	27.734600	6.607485e+04	4.789686	9589.371106	11.011695	0.218164	0.170203	
std	6.348078	6.198312e+04	4.142630	6322.086646	3.240459	0.413006	0.106782	
min	20.000000	4.000000e+03	0.000000	500.000000	5.420000	0.000000	0.000000	
25%	23.000000	3.850000e+04	2.000000	5000.000000	7.900000	0.000000	0.090000	
50%	26.000000	5.500000e+04	4.000000	8000.000000	10.990000	0.000000	0.150000	
75%	30.000000	7.920000e+04	7.000000	12200.000000	13.470000	0.000000	0.230000	

```
1 df.info()
```

```

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 32581 entries, 0 to 32580
Data columns (total 12 columns):
#   Column                                Non-Null Count  Dtype
---  -
0   person_age                            32581 non-null  int64
1   person_income                         32581 non-null  int64
2   person_home_ownership                 32581 non-null  object
3   person_emp_length                     31686 non-null  float64
4   loan_intent                           32581 non-null  object
5   loan_grade                            32581 non-null  object
6   loan_amnt                             32581 non-null  int64
7   loan_int_rate                         29465 non-null  float64
8   loan_status                           32581 non-null  int64
9   loan_percent_income                   32581 non-null  float64
10  cb_person_default_on_file              32581 non-null  object
11  cb_person_cred_hist_length             32581 non-null  int64
dtypes: float64(3), int64(5), object(4)
memory usage: 3.0+ MB

```

```

1 missing_values = df.isna().sum()
2 print(missing_values)

```

```

person_age          0
person_income       0
person_home_ownership 0
person_emp_length   895
loan_intent         0
loan_grade         0
loan_amnt          0
loan_int_rate      3116
loan_status        0
loan_percent_income 0
cb_person_default_on_file 0
cb_person_cred_hist_length 0
dtype: int64

```

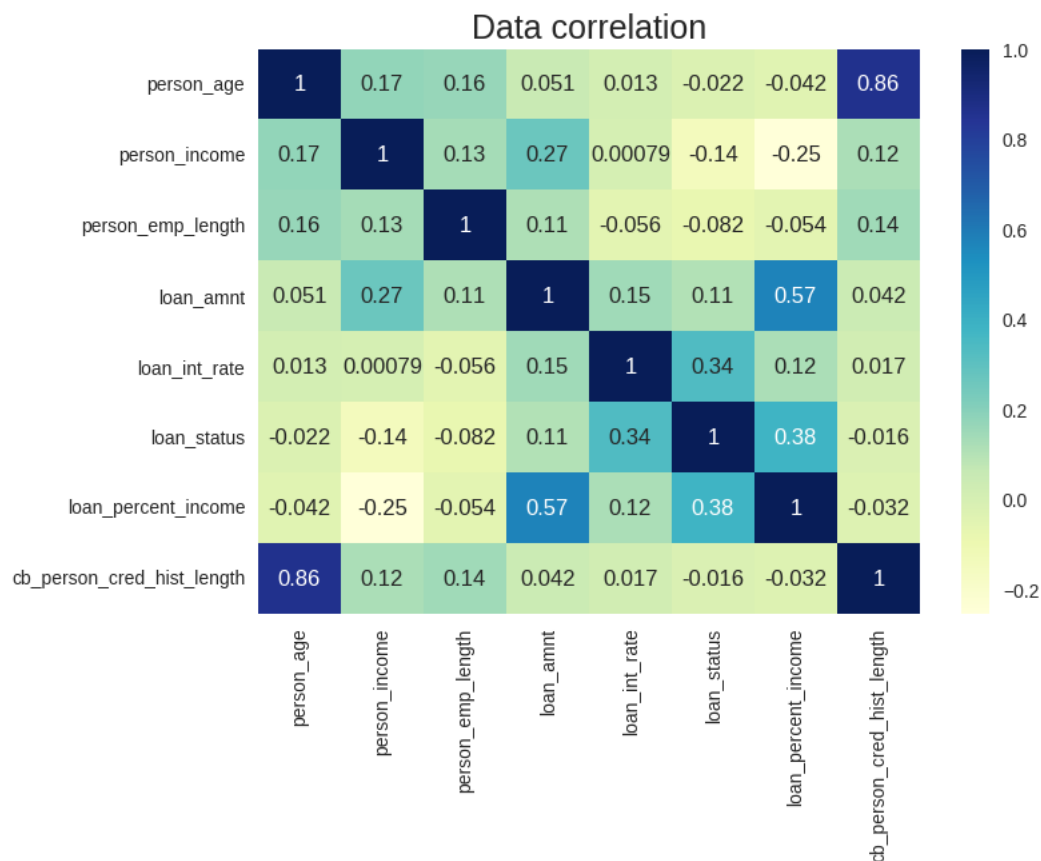
```
1 df.corr(numeric_only=True)
```

	person_age	person_income	person_emp_length	loan_amnt	loan_int_rate	loan_status	loan_percent_inc
person_age	1.000000	0.173202	0.163106	0.050787	0.012580	-0.021629	-0.042
person_income	0.173202	1.000000	0.134268	0.266820	0.000792	-0.144449	-0.254
person_emp_length	0.163106	0.134268	1.000000	0.113082	-0.056405	-0.082489	-0.054
loan_amnt	0.050787	0.266820	0.113082	1.000000	0.146813	0.105376	0.572
loan_int_rate	0.012580	0.000792	-0.056405	0.146813	1.000000	0.335133	0.120
loan_status	-0.021629	-0.144449	-0.082489	0.105376	0.335133	1.000000	0.379
loan_percent_income	-0.042411	-0.254471	-0.054111	0.572612	0.120314	0.379366	1.000
cb_person_cred_hist_length	0.859133	0.117987	0.144699	0.041967	0.016696	-0.015529	-0.031

```

1 dataplot = sns.heatmap(df.corr(numeric_only=True), cmap='YlGnBu', annot=True)
2 plt.title('Data correlation', fontsize=18)
3 plt.show()

```



```

1 sns.countplot(x=df['loan_status'])
2 plt.title('Loan Status', fontsize=18)

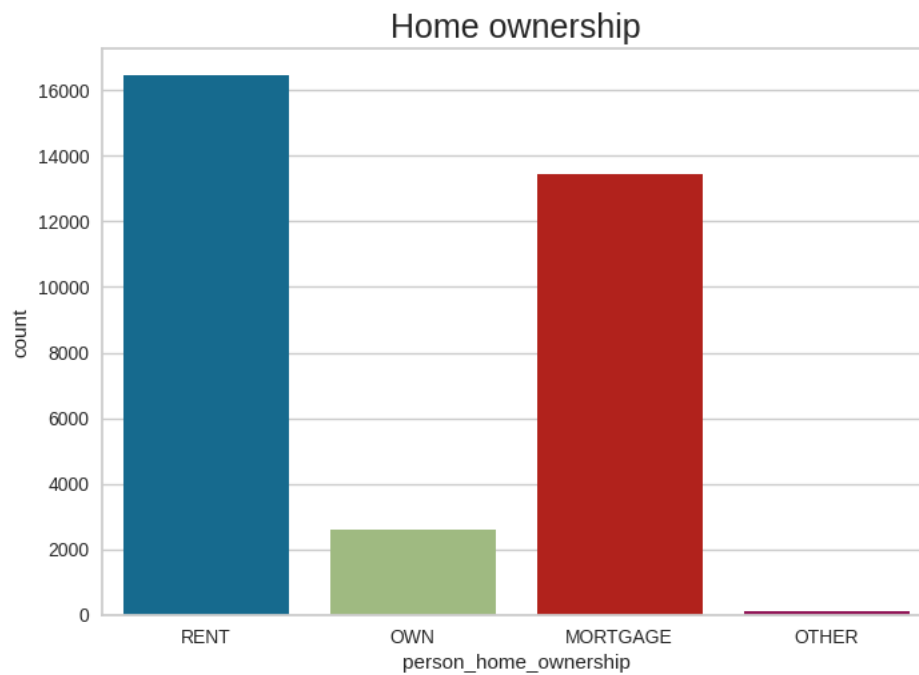
```

```
Text(0.5, 1.0, 'Loan Status')
```



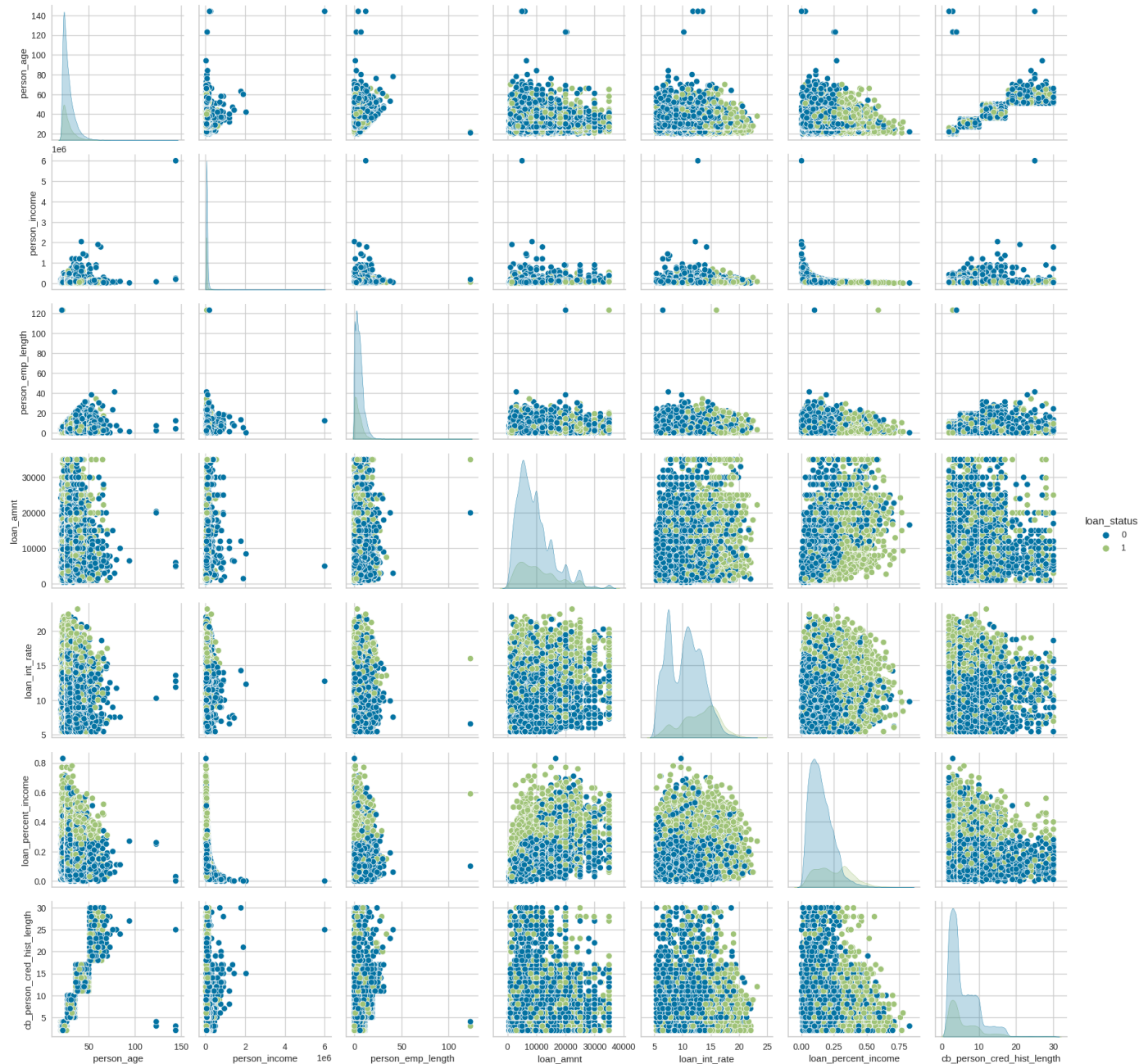
```
1 sns.countplot(x=df['person_home_ownership'])  
2 plt.title('Home ownership', fontsize=18)
```

```
Text(0.5, 1.0, 'Home ownership')
```



```
1 sns.pairplot(df, hue='loan_status')
```

&lt;seaborn.axisgrid.PairGrid at 0x7aae3f5dd060&gt;



```

1 df.isnull().sum()
2 #Filling missing values with mean:
3 df.loc[df['loan_int_rate'].isnull(), 'loan_int_rate'] = df['loan_int_rate'].median()
4 df.loc[df['person_emp_length'].isnull(), 'person_emp_length'] = df['person_emp_length'].median()
5 df.isnull().sum()

```

```

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_status         0
loan_percent_income 0
cb_person_default_on_file 0
cb_person_cred_hist_length 0
dtype: int64

```

```

1 df['person_age'].max()
2
3 #Assuming individuals with age > 90 to be errors
4 df = df.loc[df['person_age'] < 90]
5 df['person_emp_length'].max()
6
7 #Employment cannot be greater than the individual's age (accounting for childhood)
8 df = df.loc[df['person_emp_length'] < df['person_age'] - 10]

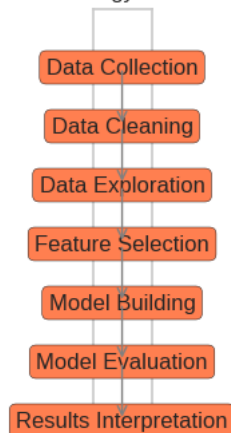
```

```

1 import matplotlib.pyplot as plt
2 # Create the approach flowchart
3 fig, ax = plt.subplots(figsize=(6, 4))
4 steps = [
5     "Data Collection",
6     "Data Cleaning",
7     "Data Exploration",
8     "Feature Selection",
9     "Model Building",
10    "Model Evaluation",
11    "Results Interpretation"
12 ]
13 positions = list(range(len(steps)))
14 step_coordinates = {}
15 for i, step in enumerate(steps):
16     x = 0.5
17     y = len(steps) - 1 - i
18     step_coordinates[step] = (x, y)
19     ax.text(x, y, step, ha='center', va='center', fontsize=12, bbox=dict(boxstyle='round,pad=0.3', facecolor='coral'))
20 for i in range(len(steps) - 1):
21     ax.annotate("", xy=step_coordinates[steps[i]], xytext=step_coordinates[steps[i+1]], arrowprops=dict(arrowstyle='<->', line
22 ax.set_xlim(0, 1)
23 ax.set_ylim(0, len(steps))
24 ax.set_xticks([])
25 ax.set_yticks([])
26 ax.set_aspect('equal')
27 plt.title("Methodology Flowchart")
28 plt.show()

```

Methodology Flowchart



## Data Processing and Encoding

```

1 # Creating groups
2 df['income_group'] = pd.cut(df['person_income'], bins= [0, 25000, 50000, 75000, 100000, float('inf')], labels=['low', 'l-midc
3 df['income_group']
4 df['loan_amnt_group'] = pd.cut(df['loan_amnt'], bins= [0, 10000, 15000, float('inf')], labels=['small', 'medium', 'large'])
5 df['loan_amnt_group']
6 df['loan_to_income'] = df['loan_amnt'] / df['person_income']
7 df['loan_to_income']

```

```

1      0.104167
2      0.572917
3      0.534351
4      0.643382
5      0.252525
...
32576  0.109434
32577  0.146875

```

```

32578    0.460526
32579    0.100000
32580    0.154167
Name: loan_to_income, Length: 32573, dtype: float64

```

```

1 #Splitting dataset in to Training set and Test set
2 y_credit = df['loan_status']
3 X_credit = df.drop(['loan_status'], axis=1)
4 X_credit.columns
5 label_encode_cols = ['person_home_ownership', 'loan_intent', 'loan_grade', 'cb_person_default_on_file', 'income_group', 'loan_status']
6 label_encoder = LabelEncoder()
7
8 for col in label_encode_cols:
9     X_credit[col] = label_encoder.fit_transform(X_credit[col])
10
11 X_credit = pd.get_dummies(X_credit, columns=label_encode_cols)
12 X_credit.head(1)
13 scaler = StandardScaler()
14 X_credit = scaler.fit_transform(X_credit)
15 X_credit[0]
16 X_training, X_test, y_training, y_test = train_test_split(X_credit, y_credit, test_size= 0.2, random_state=0)
17 X_training.shape, y_training.shape

((26058, 35), (26058,))

```

## Modelling

```

1 #Naive Bayes
2 from sklearn.naive_bayes import GaussianNB
3 naive_bayer = GaussianNB()
4 naive_bayer.fit(X_training, y_training)
5 predict_NB = naive_bayer.predict(X_test)

```

```

1 #Decision Tree
2 from sklearn.tree import DecisionTreeClassifier
3 decision_tree = DecisionTreeClassifier(criterion='entropy', random_state = 0)
4 decision_tree.fit(X_training, y_training)
5 predict_decision_tree = decision_tree.predict(X_test)

```

```

1 #Random Forest
2 from sklearn.ensemble import RandomForestClassifier
3 random_forest = RandomForestClassifier(n_estimators=200, criterion='entropy', random_state =0)
4 random_forest.fit(X_training, y_training)
5 predict_random_forest = random_forest.predict(X_test)

```

```

1 #Nearest Neighbors
2 from sklearn.neighbors import KNeighborsClassifier
3 knn = KNeighborsClassifier(n_neighbors=20)
4 knn.fit(X_training, y_training)
5 predict_knn = knn.predict(X_test)

```

```

1 #Logistic Regression
2 from sklearn.linear_model import LogisticRegression
3 logistic = LogisticRegression(random_state=1)
4 logistic.fit(X_training, y_training)
5 predict_logistic = logistic.predict(X_test)
6 logistic.intercept_

```

```
array([-1.91329537])
```

```

1 #SVM
2 from sklearn.svm import SVC
3 svm = SVC(kernel='rbf', random_state=1, C=2)
4 svm.fit(X_training, y_training)
5 predict_svm = svm.predict(X_test)

```

## Valuation Metrics

```

1 #Naive Bayes
2 from sklearn.metrics import accuracy_score, precision_score, recall_score, f1_score
3 cm = ConfusionMatrix(naive_bayer)
4 cm.fit(X_training, y_training)
5 cm.score(X_test, y_test)
6 # Calculate accuracy
7 accuracy = accuracy_score(y_test, predict_NB)
8
9 # Calculate precision
10 precision = precision_score(y_test, predict_NB)
11
12 # Calculate recall
13 recall = recall_score(y_test, predict_NB)
14
15 # Calculate F1-score
16 f1 = f1_score(y_test, predict_NB)
17
18 print("\033[1mNaive Bayes\033[0m")
19 print(f"Accuracy: {accuracy:.2f}")
20 print(f"Precision: {precision:.2f}")
21 print(f"Recall: {recall:.2f}")
22 print(f"F1-Score: {f1:.2f}")

```

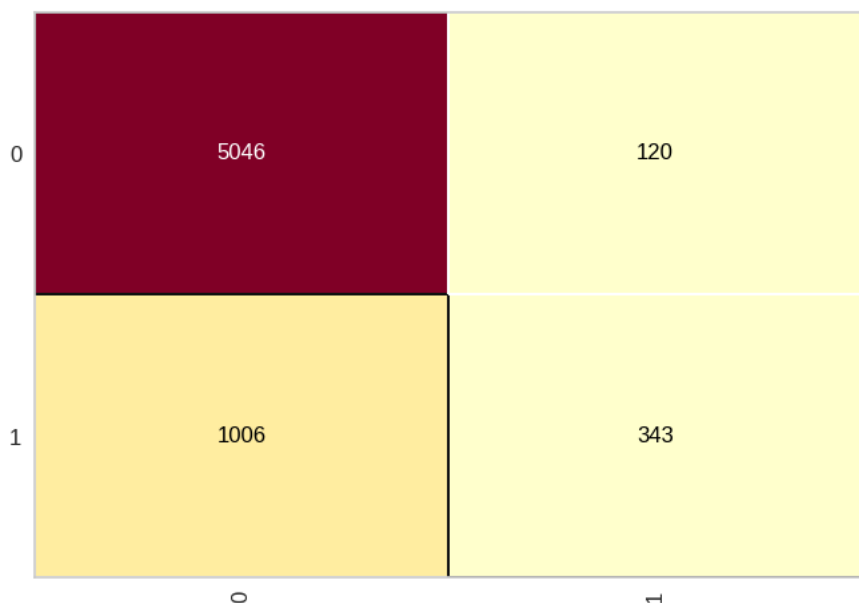
**Naive Bayes**

Accuracy: 0.83

Precision: 0.74

Recall: 0.25

F1-Score: 0.38

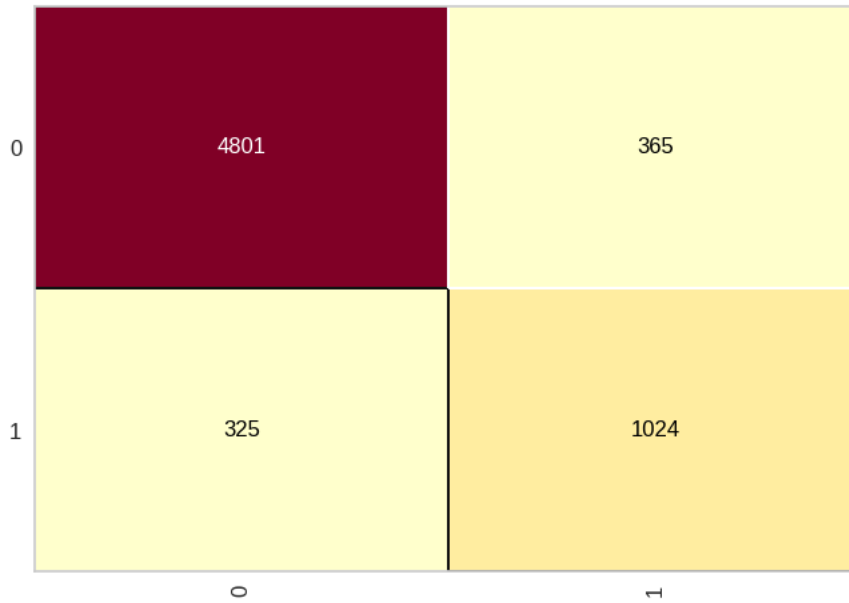


```

1 #Decision Tree
2 cm = ConfusionMatrix(decision_tree)
3 cm.fit(X_training, y_training)
4 cm.score(X_test, y_test)
5 from sklearn.metrics import accuracy_score, precision_score, recall_score, f1_score
6
7 # Calculate accuracy
8 accuracy = accuracy_score(y_test, predict_decision_tree)
9
10 # Calculate precision
11 precision = precision_score(y_test, predict_decision_tree)
12
13 # Calculate recall
14 recall = recall_score(y_test, predict_decision_tree)
15
16 # Calculate F1-score
17 f1 = f1_score(y_test, predict_decision_tree)
18
19 print("\033[1mDecision Tree\033[0m")
20 print(f"Accuracy: {accuracy:.2f}")
21 print(f"Precision: {precision:.2f}")
22 print(f"Recall: {recall:.2f}")
23 print(f"F1-Score: {f1:.2f}")

```

**Decision Tree**  
 Accuracy: 0.89  
 Precision: 0.74  
 Recall: 0.76  
 F1-Score: 0.75



```

1 #Random Forest
2 cm = ConfusionMatrix(random_forest)
3 cm.fit(X_training, y_training)
4 cm.score(X_test, y_test)
5 from sklearn.metrics import accuracy_score, precision_score, recall_score, f1_score
6
7 # Calculate accuracy
8 accuracy = accuracy_score(y_test, predict_random_forest)
9
10 # Calculate precision
11 precision = precision_score(y_test, predict_random_forest)
12
13 # Calculate recall
14 recall = recall_score(y_test, predict_random_forest)
15
16 # Calculate F1-score
17 f1 = f1_score(y_test, predict_random_forest)
18
19 print("\033[1mRandom Forest\033[0m")
20 print(f"Accuracy: {accuracy:.2f}")
21 print(f"Precision: {precision:.2f}")
22 print(f"Recall: {recall:.2f}")
23 print(f"F1-Score: {f1:.2f}")

```



**Random Forest**  
 Accuracy: 0.94  
 Precision: 0.97  
 Recall: 0.72  
 F1-Score: 0.82

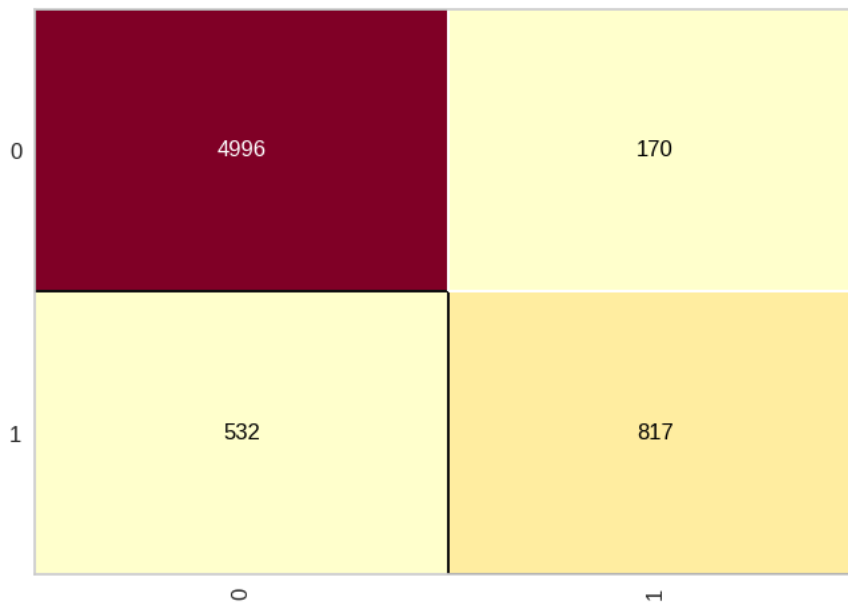


```

1 #Nearest Neighbors
2 cm = ConfusionMatrix(knn)
3 cm.fit(X_training, y_training)
4 cm.score(X_test, y_test)
5 from sklearn.metrics import accuracy_score, precision_score, recall_score, f1_score
6
7 # Calculate accuracy
8 accuracy = accuracy_score(y_test, predict_knn)
9
10 # Calculate precision
11 precision = precision_score(y_test, predict_knn)
12
13 # Calculate recall
14 recall = recall_score(y_test, predict_knn)
15
16 # Calculate F1-score
17 f1 = f1_score(y_test, predict_knn)
18
19 print("\033[1mNearest Neighbors\033[0m")
20 print(f"Accuracy: {accuracy:.2f}")
21 print(f"Precision: {precision:.2f}")
22 print(f"Recall: {recall:.2f}")
23 print(f"F1-Score: {f1:.2f}")

```

**Nearest Neighbors**  
 Accuracy: 0.89  
 Precision: 0.83  
 Recall: 0.61  
 F1-Score: 0.70



```
1 #Logistic Regression
2 cm = ConfusionMatrix(logistic)
3 cm.fit(X_training, y_training)
4 cm.score(X_test, y_test)
5 from sklearn.metrics import accuracy_score, precision_score, recall_score, f1_score
6
7 # Calculate accuracy
8 accuracy = accuracy_score(y_test, predict_logistic)
9
10 # Calculate precision
11 precision = precision_score(y_test, predict_logistic)
12
13 # Calculate recall
14 recall = recall_score(y_test, predict_logistic)
15
16 # Calculate F1-score
17 f1 = f1_score(y_test, predict_logistic)
18
19 print("\033[1mLogistic Regression\033[0m")
20 print(f"Accuracy: {accuracy:.2f}")
21 print(f"Precision: {precision:.2f}")
22 print(f"Recall: {recall:.2f}")
23 print(f"F1-Score: {f1:.2f}")
```

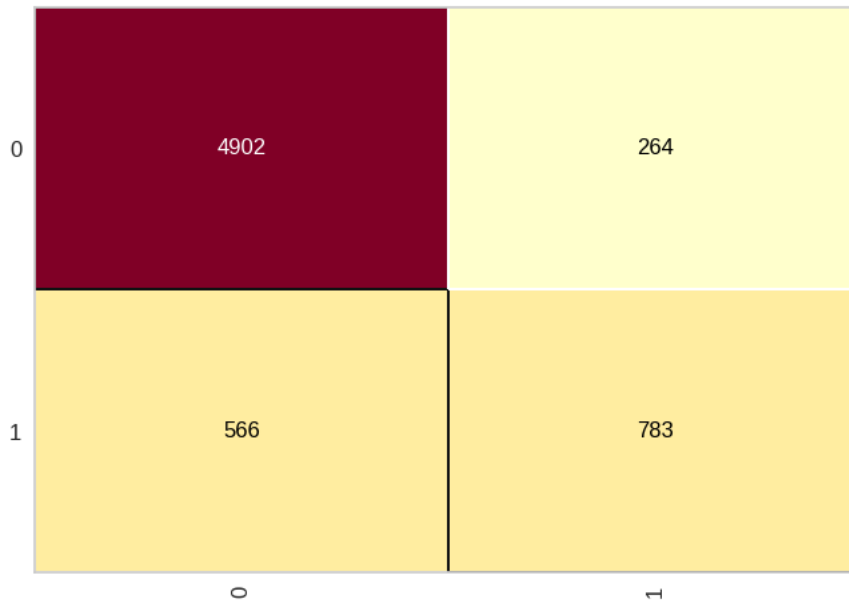
**Logistic Regression**

Accuracy: 0.87

Precision: 0.75

Recall: 0.58

F1-Score: 0.65



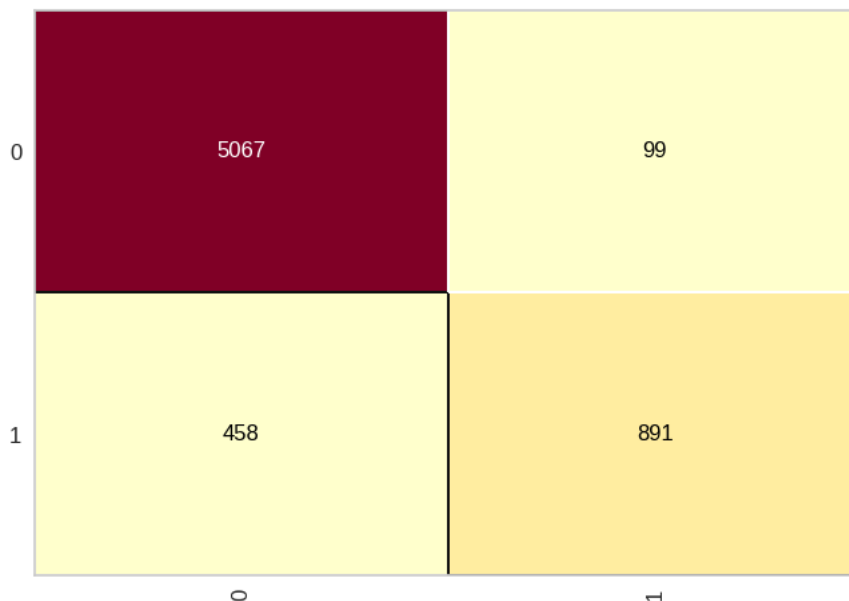
```

1 #SVM
2 cm = ConfusionMatrix(svm)
3 cm.fit(X_training, y_training)
4 cm.score(X_test, y_test)
5
6 from sklearn.metrics import accuracy_score, precision_score, recall_score, f1_score
7
8 # Calculate accuracy
9 accuracy_svm = accuracy_score(y_test, predict_svm)
10
11 # Calculate precision
12 precision_svm = precision_score(y_test, predict_svm)
13
14 # Calculate recall
15 recall_svm = recall_score(y_test, predict_svm)
16
17 # Calculate F1-score
18 f1_svm = f1_score(y_test, predict_svm)
19
20 print("\033[1mSVM\033[0m")
21 print(f"Accuracy: {accuracy_svm:.2f}")
22 print(f"Precision: {precision_svm:.2f}")
23 print(f"Recall: {recall_svm:.2f}")
24 print(f"F1-Score: {f1_svm:.2f}")

```

**SVM**

Accuracy: 0.91  
Precision: 0.90  
Recall: 0.66  
F1-Score: 0.76

**Evaluation of Models**

```

1 from datetime import datetime
2 from sklearn.model_selection import cross_val_score, KFold
3 from sklearn.metrics import accuracy_score, precision_score, recall_score, f1_score
4
5 # Define a function to evaluate and print model performance
6 def evaluate_model(model, X_train, X_test, y_train, y_test):
7     # Time to Train
8     start_time = datetime.now()
9     model.fit(X_train, y_train)
10    end_time = datetime.now()
11    time_to_train = end_time - start_time
12    print(f"Time to Train: {time_to_train}")
13
14    # Time to Test
15    start_time = datetime.now()
16    y_pred = model.predict(X_test)
17    end_time = datetime.now()
18    time_to_test = end_time - start_time
19    print(f"Time to Test: {time_to_test}")

```

```

20
21 # Stability using cross-validation
22 num_folds = 5
23 kf = KFold(n_splits=num_folds, shuffle=True, random_state=0)
24 cross_val_results = cross_val_score(model, X_train, y_train, cv=kf, scoring='accuracy')
25 print(f"Cross-validation results: {cross_val_results}")
26 print(f"Mean accuracy: {cross_val_results.mean()}")
27 print(f"Standard deviation of accuracy: {cross_val_results.std()}")
28
29 # Performance Metrics on Test Set
30 accuracy = accuracy_score(y_test, y_pred)
31 precision = precision_score(y_test, y_pred)
32 recall = recall_score(y_test, y_pred)
33 f1 = f1_score(y_test, y_pred)
34 print(f"Accuracy: {accuracy:.2f}")
35 print(f"Precision: {precision:.2f}")
36 print(f"Recall: {recall:.2f}")
37 print(f"F1-Score: {f1:.2f}")
38
39 # Apply the function for each model
40 models = [naive_bayer, decision_tree, random_forest, knn, logistic, svm]
41 for model in models:
42     print(f"\nEvaluating {model.__class__.__name__}:")
43     evaluate_model(model, X_training, X_test, y_training, y_test)

```

```

Precision: 0.74
Recall: 0.25
F1-Score: 0.38

```

```

Evaluating DecisionTreeClassifier:
Time to Train: 0:00:00.274620
Time to Test: 0:00:00.002521
Cross-validation results: [0.8862241  0.89850345 0.88718342 0.88140472 0.89195932]
Mean accuracy: 0.8890550024573143
Standard deviation of accuracy: 0.005792269139852778
Accuracy: 0.89
Precision: 0.74
Recall: 0.76
F1-Score: 0.75

```

```

Evaluating RandomForestClassifier:
Time to Train: 0:00:07.496751
Time to Test: 0:00:00.334772
Cross-validation results: [0.93016117 0.93937068 0.92689946 0.93033967 0.9337939 ]
Mean accuracy: 0.932112975194306
Standard deviation of accuracy: 0.004234022129317153
Accuracy: 0.94
Precision: 0.97
Recall: 0.72
F1-Score: 0.82

```

```

Evaluating KNeighborsClassifier:
Time to Train: 0:00:00.004485
Time to Test: 0:00:01.256882
Cross-validation results: [0.88660783 0.89140445 0.8823868  0.87622337 0.88677797]
Mean accuracy: 0.8846800844721148

```