

Code Report Credit Default Prediction

Student ID: 577610

Author: Marcel De Lange

Date:: 26-July-2024 - 28 July 2024

Table of Contents

1. Problem Formulation	2
Problem Statement	2
Motivation	2
2. Features/Variables in the Problem	2
Features Explanation	2
3. Explanation of the AI/ML Algorithm Used	3
Logistic Regression	3
4. Suitable Plots to Show How the AI/ML Algorithm is Applied	3
4.1. Box Plots	3
Results:	4
4.2. Histograms	5
Results:	5
4.3. Correlation Matrix Heatmap	6
Result:	6
4.4. ROC Curve	6
Result:	7
Result:	8
5. Suitable Analytics of the Data	8
5.1. Descriptive Statistics	8
Result:	8
5.2. Feature Correlation	9
Result:	10
6. Evaluation of the Model	10
6.1. Training Accuracy	10
Result:	11
6.2. Testing Accuracy	11
Result:	11
6.3 Prediction Example	11
Result:	11
7. Tools used	11
8. References	12

1. Problem Formulation

Problem Statement

- The primary objective of this project is to predict whether a credit card holder is likely to default on their payment in the next month. Default prediction is crucial for financial institutions to manage risk and make informed decisions about credit issuance and management.

Motivation

- Predicting credit default allows financial institutions to proactively manage credit risk by identifying high-risk customers. This can lead to better financial planning, reduced losses, and targeted interventions to help prevent defaults.

2. Features/Variables in the Problem

Features Explanation

1. **LIMIT (LIMIT_BAL):** Amount of credit given to the cardholder (in dollars). This feature represents the credit limit assigned to the cardholder, which can influence their ability to default.
2. **SEX (SEX):** Gender of the credit card holder.
 - 1 = Male
 - 2 = Female
3. **EDUCATION (EDUCATION):** Educational level of the credit card holder.
 - 1 = Graduate school
 - 2 = University
 - 3 = High school
 - 4 = Other
4. **MARRIAGE (MARRIAGE):** Marital status of the credit card holder.
 - 1 = Married
 - 2 = Single
 - 3 = Others
5. **AGE (AGE):** Age of the credit card holder (in years). Age can affect the likelihood of defaulting due to varying financial responsibilities and stability.
6. **PAY_0 to PAY_6:** History of past payment. This feature shows the repayment status from the last 6 months.
 - -1 = Paid duly
 - 1 = Payment delay for one month
 - 2 = Payment delay for two months
 - ...
 - 9 = Payment delay for nine months and above

7. **BILL_AMT1 to BILL_AMT6:** Amount of bill statement for the past 6 months.
Represents the credit card bill amounts for each month.
8. **PAY_AMT1 to PAY_AMT6:** Amount of previous payment for the past 6 months.
Represents the amount paid towards the credit card bill each month.
9. **Target Variable (default_payment_next_month):** Indicates whether the cardholder defaulted on payment in the next month.
 - 1 = Default
 - 0 = No Default

3. Explanation of the AI/ML Algorithm Used

Logistic Regression

- Logistic Regression is a classification algorithm used to model the probability of a binary outcome based on one or more predictor variables. The key steps in Logistic Regression include:
 1. **Logistic Function:** The model applies the logistic (sigmoid) function to predict probabilities: $\sigma(z) = \frac{1}{1 + e^{-z}}$ where $z = \mathbf{w}^T \mathbf{x} + b$.
 2. **Decision Boundary:** The model classifies an instance as default (1) if the predicted probability is greater than or equal to 0.5, otherwise as no default (0).
 3. **Training:** The model learns optimal weights and bias by minimizing the binary cross-entropy loss function:

$$J(\mathbf{w}, b) = -\frac{1}{m} \sum_{i=1}^m [y(i) \log(\sigma(\mathbf{w}^T \mathbf{x}(i))) + (1 - y(i)) \log(1 - \sigma(\mathbf{w}^T \mathbf{x}(i)))]$$
 4. **Evaluation Metrics:** The performance is assessed using accuracy, precision and other evaluation metrics.

4. Suitable Plots to Show How the AI/ML Algorithm is Applied

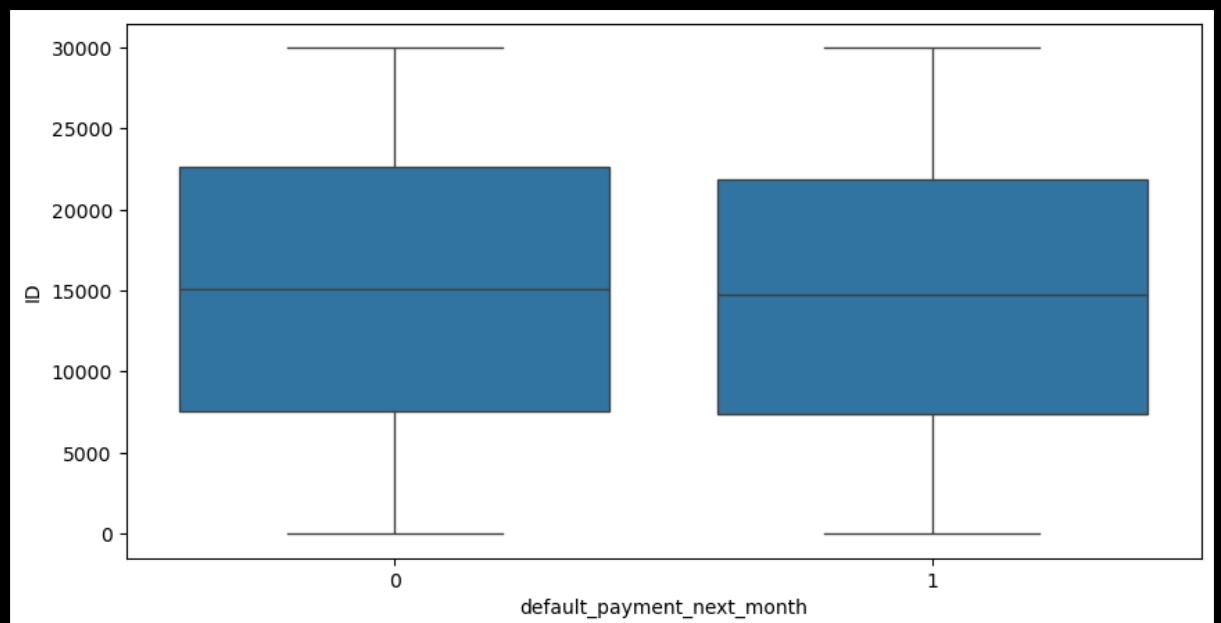
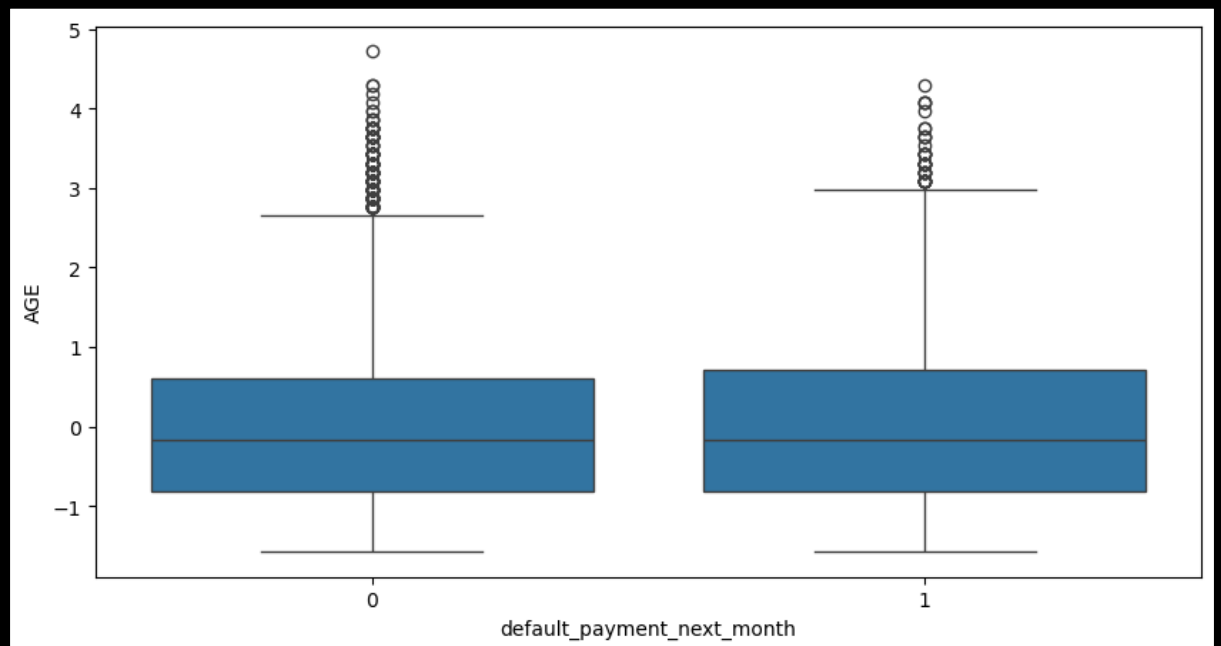
4.1. Box Plots

- Box plots can be used to visualize the distribution of numerical features across the default and no-default classes. This helps in understanding feature variation between the two classes.

python

- `for column in X.columns:`
- `plt.figure(figsize=(10, 5))`
- `sns.boxplot(x=y, y=column, data=X)`
- `plt.title(f'Box Plot of {column}')`
- `plt.show()`

Results:



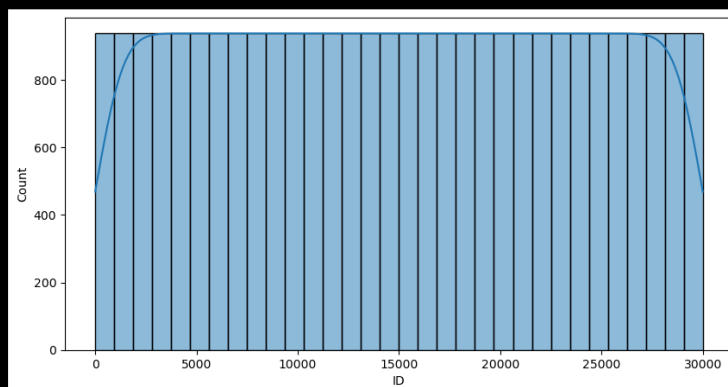
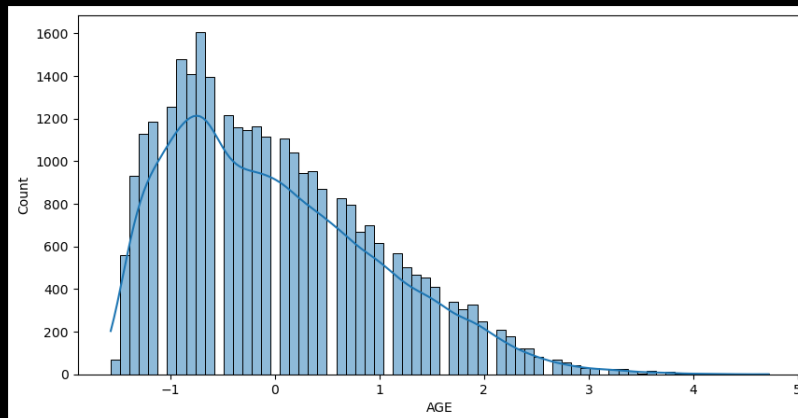
4.2. Histograms

- Histograms show the distribution of individual features across the dataset, highlighting the spread and skewness of the data.

python

- `for column in X.columns:`
- `plt.figure(figsize=(10, 5))`
- `sns.histplot(X[column], kde=True)`
- `plt.show()`

Results:



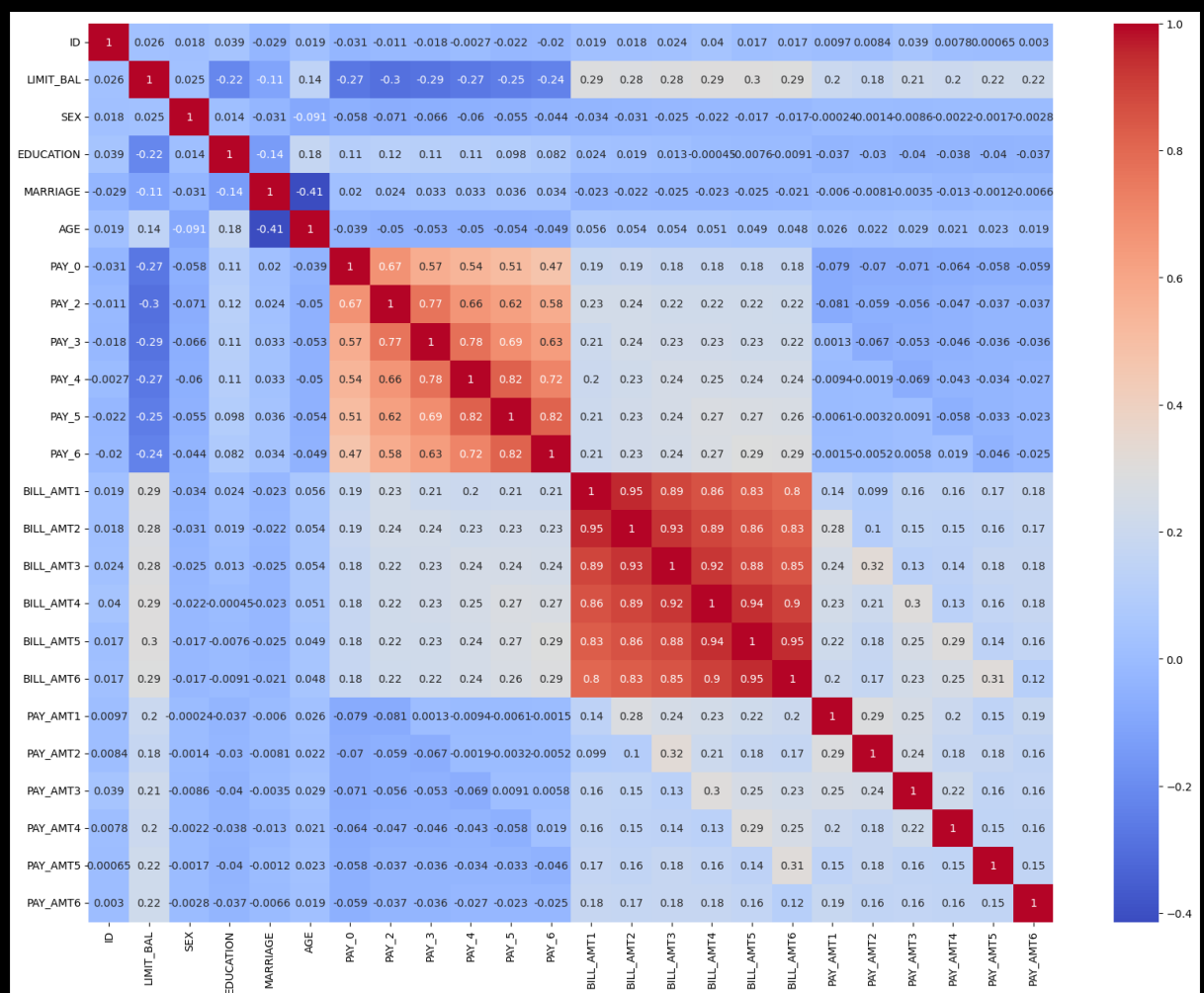
4.3. Correlation Matrix Heatmap

- A heatmap of the correlation matrix shows the relationships between numerical features, indicating how features are related to one another.

python

- `plt.figure(figsize=(20, 15))`
- `sns.heatmap(correlation_matrix, annot=True, cmap='coolwarm')`
- `plt.show()`

Result:



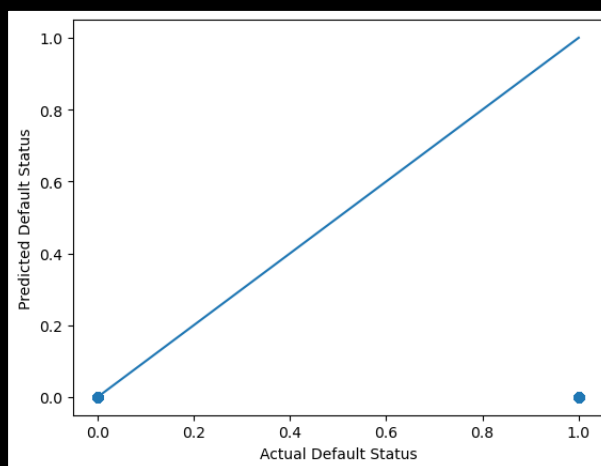
4.4. ROC Curve

- The ROC Curve shows the evaluation of train and test data in scatterplots.

python

- # Evaluation of Training data predcition vs actual training data in a sactterplot
 - o plt.scatter(y_train, y_pred_train)
 - o plt.xlabel('Actual Default Status')
 - o plt.ylabel('Predicted Default Status')
- # add a logistic line
 - o fpr, tpr, thresholds = roc_curve(y_train, y_pred_train)
 - o plt.plot(fpr, tpr)
 - o plt.show()

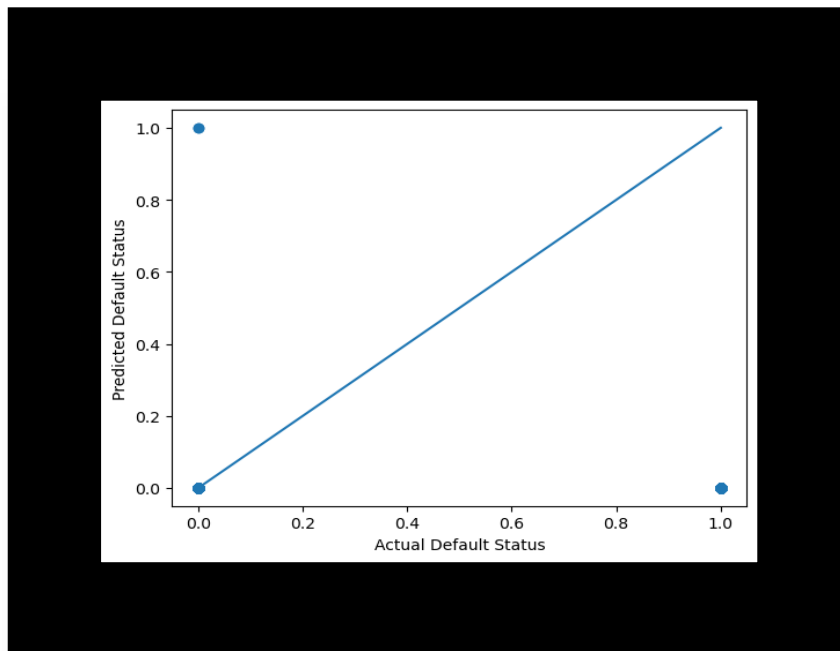
Result:



Scatter plot actual test data vs predicted data

- o plt.scatter(y_test, y_pred_test)
- o plt.xlabel('Actual Default Status')
- o plt.ylabel('Predicted Default Status')
- o # add a suitable line
- o fpr, tpr, thresholds = roc_curve(y_test, y_pred_test)
- o plt.plot(fpr, tpr)
- o plt.show()

Result:



5. Suitable Analytics of the Data

5.1. Descriptive Statistics

python

- `df.describe().T`

Result:

	mean	std	min	25%	50%	75%	max	
count								
ID	3000 0.0	15000.5 00000	8660.39 8374	1.0	7500 .75	1500 0.5	22500 .25	30000 .0
LIMIT_BAL	3000 0.0	167484. 322667	129747. 661567	1000 0.0	5000 0.00	1400 00.0	24000 0.00	10000 00.0
SEX	3000 0.0	1.60373 3	0.48912 9	1.0	1.00	2.0	2.00	2.0
EDUCATION	3000 0.0	1.85313 3	0.79034 9	0.0	1.00	2.0	2.00	6.0
MARRIAGE	3000 0.0	1.55186 7	0.52197 0	0.0	1.00	2.0	2.00	3.0
AGE	3000 0.0	35.4855 00	9.21790 4	21.0	28.0 0	34.0	41.00	79.0
PAY_0	3000 0.0	- 0.01670 0	1.12380 2	-2.0	- 1.00	0.0	0.00	8.0

PAY_2	3000 0.0	- 0.13376 7	1.19718 6	-2.0	- 1.00	0.0	0.00	8.0
PAY_3	3000 0.0	- 0.16620 0	1.19686 8	-2.0	- 1.00	0.0	0.00	8.0
PAY_4	3000 0.0	- 0.22066 7	1.16913 9	-2.0	- 1.00	0.0	0.00	8.0
PAY_5	3000 0.0	- 0.26620 0	1.13318 7	-2.0	- 1.00	0.0	0.00	8.0
PAY_6	3000 0.0	- 0.29110 0	1.14998 8	-2.0	- 1.00	0.0	0.00	8.0
BILL_AMT1	3000 0.0	51223.3 30900	73635.8 60576	- 1655 80.0	3558 .75	2238 1.5	67091 .00	96451 1.0
BILL_AMT2	3000 0.0	49179.0 75167	71173.7 68783	- 6977 7.0	2984 .75	2120 0.0	64006 .25	98393 1.0
BILL_AMT3	3000 0.0	47013.1 54800	69349.3 87427	- 1572 64.0	2666 .25	2008 8.5	60164 .75	16640 89.0
BILL_AMT4	3000 0.0	43262.9 48967	64332.8 56134	- 1700 00.0	2326 .75	1905 2.0	54506 .00	89158 6.0
BILL_AMT5	3000 0.0	40311.4 00967	60797.1 55770	- 8133 4.0	1763 .00	1810 4.5	50190 .50	92717 1.0
BILL_AMT6	3000 0.0	38871.7 60400	59554.1 07537	- 3396 03.0	1256 .00	1707 1.0	49198 .25	96166 4.0
PAY_AMT1	3000 0.0	5663.58 0500	16563.2 80354	0.0	1000 .00	2100 .0	5006. 00	87355 2.0
PAY_AMT2	3000 0.0	5921.16 3500	23040.8 70402	0.0	833. 00	2009 .0	5000. 00	16842 59.0
PAY_AMT3	3000 0.0	5225.68 1500	17606.9 61470	0.0	390. 00	1800 .0	4505. 00	89604 0.0
PAY_AMT4	3000 0.0	4826.07 6867	15666.1 59744	0.0	296. 00	1500 .0	4013. 25	62100 0.0
PAY_AMT5	3000 0.0	4799.38 7633	15278.3 05679	0.0	252. 50	1500 .0	4031. 50	42652 9.0
PAY_AMT6	3000 0.0	5215.50 2567	17777.4 65775	0.0	117. 75	1500 .0	4000. 00	52866 6.0
default_payment_next_month	3000 0.0	0.22120 0	0.41506 2	0.0	0.00	0.0	0.00	1.0

5.2. Feature Correlation

Analyze the correlation between features to identify which features are strongly related to each other.

```
python
```

```
# Correlation Matrix
```

```
o correlation_matrix = X.corr()
```

```

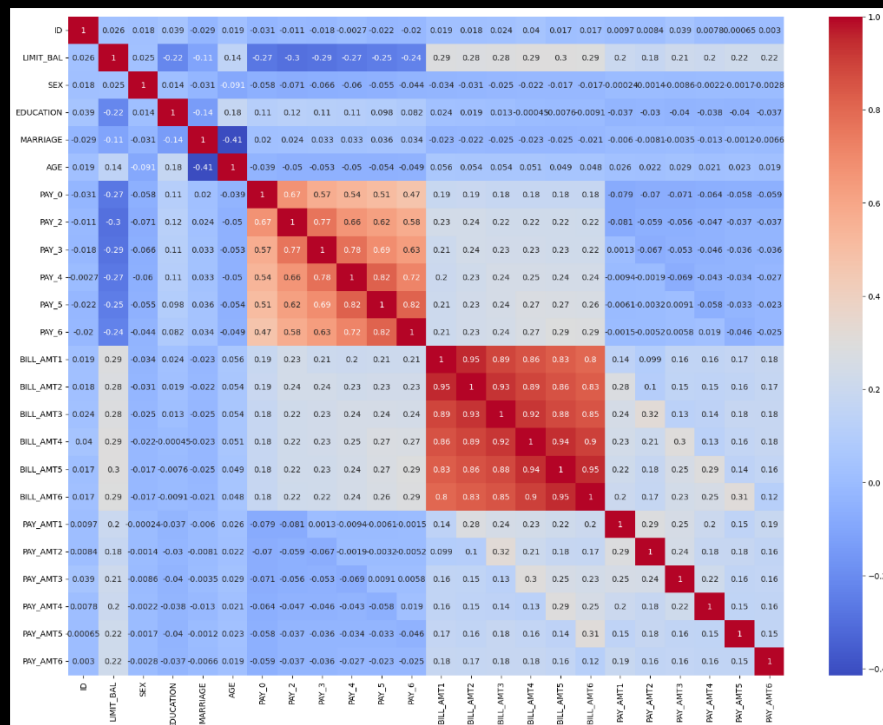
o correlation_matrix

# Heatmap

o plt.figure(figsize=(20, 15))
o sns.heatmap(correlation_matrix, annot=True, cmap='coolwarm')
o plt.show()

```

Result:



6. Evaluation of the Model

6.1. Training Accuracy

Measure the accuracy of the model on the training dataset.

python

```

o accuracy_train = accuracy_score(y_train, y_pred_train)
o accuracy_train

```

Result:

- 0.778125

6.2. Testing Accuracy

Measure the accuracy of the model on the testing dataset.

python

- o `accuracy_test = accuracy_score(y_test, y_pred_test)`
- o `accuracy_test`

Result:

- 0.7811666666666667

6.3 Prediction Example

Provide an example prediction and its interpretation.

python

- o `y_pred_full_case = model.predict([full_x.values.tolist()])`
- o `y_pred_full_case`
- o `# Display the prediction`
- o `if y_pred_full_case[0] == 0:`
- o `print('The customer is not likely to default.')`
- o `else:`
- o `print('The customer is likely to default.')`

Result:

- The customer is not likely to default.

7. Tools used :

- o VS code for code and debugging
- o Data Wrangler for viewing the dataset and gaining key insights from it
- o Chat GPT for the layout used in this report
- o Tabnine AI code completion for assistance in structuring the coding and assisting in debugging
- o Dataset Source
- o Link :
 - o <https://archive.ics.uci.edu/dataset/350/default+of+credit+card+clients>
- o Python Libraries used:
 - o `import pandas as pd`
 - o `import matplotlib.pyplot as plt`

- `import seaborn as sns`
- `from sklearn.preprocessing import LabelEncoder`
- `from sklearn.model_selection import train_test_split`
- `from sklearn.linear_model import LogisticRegression`
- `from sklearn.metrics import accuracy_score`
- `from sklearn.preprocessing import StandardScaler`
- `from sklearn.metrics import roc_curve`
- `from sklearn.metrics import accuracy_score`

8. References :

Books

1. **"Pattern Recognition and Machine Learning"** by Christopher M. Bishop
 - This book provides a thorough introduction to various machine learning algorithms, including Logistic Regression, with mathematical foundations and practical applications.
2. **"Machine Learning: A Probabilistic Perspective"** by Kevin P. Murphy
 - Murphy's book offers detailed explanations of machine learning techniques, including Logistic Regression, with emphasis on probabilistic models.
3. **"Introduction to Statistical Learning: with Applications in R"** by Gareth James, Daniela Witten, Trevor Hastie, and Robert Tibshirani
 - This book is a great resource for understanding statistical learning methods, including Logistic Regression, with practical examples and applications.

Online Courses & Tutorials

1. **Coursera – "Machine Learning" by Andrew Ng**
 - This popular course provides an accessible introduction to various machine learning algorithms, including Logistic Regression. It is available for free and covers both theoretical and practical aspects.
2. **Khan Academy – "Logistic Regression"**
 - Khan Academy offers a video tutorial that covers the basics of Logistic Regression and its application.
3. **edX – "Introduction to Machine Learning with Python"**
 - This course offers a practical introduction to machine learning using Python, including Logistic Regression.

Research Papers

1. **"Logistic Regression: An Overview"** by David W. Hosmer Jr., Stanley Lemeshow, and Rodney X. Sturdivant
 - A foundational paper that provides an in-depth review of Logistic Regression, including its theoretical aspects and applications.
2. **"The Use of Logistic Regression in Epidemiology: A Review"** by T. L. M. B. Croghan
 - This paper discusses the application of Logistic Regression in epidemiology, providing context on its use in different fields.