```
In [1]: # Initialize Otter
        import otter
        grader = otter.Notebook("hw5.ipynb")
```

# CPSC 330 - Applied Machine Learning

## Homework 5: Putting it all together

### Associated lectures: All material till lecture 13

**Due date: Monday, October 28th, 2024 at 11:59pm**

## Table of contents

## Submission instructions

rubric={points:4}

**You may work with a partner on this homework and submit your assignment as a group.** Below are some instructions on working as a group.

- The maximum group size is 2.
- Use group work as an opportunity to collaborate and learn new things from each other.

- Be respectful to each other and make sure you understand all the concepts in the assignment well.
- It's your responsibility to make sure that the assignment is submitted by one of the group members before the deadline.
- You can find the instructions on how to do group submission on Gradescope here.
- If you would like to use late tokens for the homework, all group members must have the necessary late tokens available. Please note that the late tokens will be counted for all members of the group.

Follow the homework submission instructions.

1. Before submitting the assignment, run all cells in your notebook to make sure there are no errors by doing `Kernel -> Restart Kernel and Clear All Outputs` and then `Run -> Run All Cells`.
2. Notebooks with cell execution numbers out of order or not starting from "1" will have marks deducted. Notebooks without the output displayed may not be graded at all (because we need to see the output in order to grade your work).
3. Follow the CPSC 330 homework instructions, which include information on how to do your assignment and how to submit your assignment.
4. Upload your solution on Gradescope. Check out this Gradescope Student Guide if you need help with Gradescope submission.
5. Make sure that the plots and output are rendered properly in your submitted file. If the .ipynb file is too big and doesn't render on Gradescope, also upload a pdf or html in addition to the .ipynb so that the TAs can view your submission on Gradescope.

*Note: The assignments will get gradually more open-ended as we progress through the course. In many cases, there won't be a single correct solution. Sometimes you will have to make your own choices and your own decisions (for example, on what parameter values to use when they are not explicitly provided in the instructions). Use your own judgment in such cases and justify your choices, if necessary.*

# Imports

Imports

*Points:* 0

```
In [2]:  from hashlib import sha1

         import matplotlib.pyplot as plt
         import numpy as np
```

```python
import pandas as pd

plt.rcParams["font.size"] = 16

from sklearn.dummy import DummyClassifier
from sklearn.feature_extraction.text import CountVectorizer
from sklearn.linear_model import LogisticRegression
from sklearn.model_selection import (
    GridSearchCV,
    cross_val_score,
    cross_validate,
    train_test_split,
)
from sklearn.pipeline import Pipeline, make_pipeline
from sklearn.tree import DecisionTreeClassifier
from sklearn.compose import ColumnTransformer, make_column_transformer
from sklearn.impute import SimpleImputer
from sklearn.model_selection import cross_val_score, cross_validate, train_t
from sklearn.neighbors import KNeighborsClassifier
from sklearn.pipeline import Pipeline, make_pipeline
from sklearn.preprocessing import OneHotEncoder, StandardScaler, OrdinalEnco
from sklearn.svm import SVC
from sklearn.tree import DecisionTreeClassifier
from sklearn.model_selection import RandomizedSearchCV
from sklearn.metrics import ConfusionMatrixDisplay
from sklearn.feature_selection import RFECV
import itertools
```

# Introduction

In this homework you will be working on an open-ended mini-project, where you will put all the different things you have learned so far together to solve an interesting problem.

A few notes and tips when you work on this mini-project:

## Tips

1. This mini-project is open-ended, and while working on it, there might be some situations where you'll have to use your own judgment and make your own decisions (as you would be doing when you work as a data scientist). Make sure you explain your decisions whenever necessary.
2. **Do not include everything you ever tried in your submission** -- it's fine just to have your final code. That said, your code should be reproducible and well-documented. For example, if you chose your hyperparameters based on some hyperparameter optimization experiment, you should leave in the code for that experiment so that someone else could re-run it and obtain the same hyperparameters, rather than mysteriously just setting the hyperparameters to some (carefully chosen) values in your code.
3. If you realize that you are repeating a lot of code try to organize it in functions. Clear presentation of your code, experiments, and results is the key to be successful in

this lab. You may use code from lecture notes or previous lab solutions with
appropriate attributions.

## Assessment

We plan to grade fairly and leniently. We don't have some secret target score that you
need to achieve to get a good grade. **You'll be assessed on demonstration of mastery
of course topics, clear presentation, and the quality of your analysis and results.**
For example, if you just have a bunch of code and no text or figures, that's not good. If
you do a bunch of sane things and get a lower accuracy than your friend, don't sweat it.

## A final note

Finally, this style of this "project" question is different from other assignments. It'll be up
to you to decide when you're "done" -- in fact, this is one of the hardest parts of real
projects. But please don't spend WAY too much time on this... perhaps "a few hours"
(15-20 hours???) is a good guideline for this project . Of course if you're having fun
you're welcome to spend as much time as you want! But, if so, try not to do it out of
perfectionism or getting the best possible grade. Do it because you're learning and
enjoying it. Students from the past cohorts have found such kind of labs useful and fun
and I hope you enjoy it as well.

# 1. Pick your problem and explain the prediction problem

rubric={points:3}
In this mini project, you have the option to choose on which dataset you will be working
on. The tasks you will need to carry on will be similar, independently of your choice.

## Option 1

You can choose to work on a classification problem of predicting whether a credit card
client will default or not. For this problem, you will use Default of Credit Card Clients
Dataset. In this data set, there are 30,000 examples and 24 features, and the goal is to
estimate whether a person will default (fail to pay) their credit card bills; this column is
labeled "default.payment.next.month" in the data. The rest of the columns can be used
as features. You may take some ideas and compare your results with the associated
research paper, which is available through the UBC library.

## Option 2

You can choose to work on a regression problem using a dataset of New York City Airbnb listings from 2019. As usual, you'll need to start by downloading the dataset, then you will try to predict `reviews_per_month`, as a proxy for the popularity of the listing. Airbnb could use this sort of model to predict how popular future listings might be before they are posted, perhaps to help guide hosts create more appealing listings. In reality they might instead use something like vacancy rate or average rating as their target, but we do not have that available here.

> Note there is an updated version of this dataset with more features available here. The features were are using in `listings.csv.gz` for the New York city datasets. You will also see some other files like `reviews.csv.gz`. For your own interest you may want to explore the expanded dataset and try your analysis there. However, please submit your results on the dataset obtained from Kaggle.

**Your tasks:**

1. Spend some time understanding the options and pick the one you find more interesting (it may help spending some time looking at the documentation available on Kaggle for each dataset).
2. After making your choice, focus on understanding the problem and what each feature means, again using the documentation on the dataset page on Kaggle. Write a few sentences on your initial thoughts on the problem and the dataset.
3. Download the dataset and read it as a pandas dataframe.

Solution_1

*Points:* 3

*Type your answer here, replacing this text.*

```
In [3]: df = pd.read_csv('UCI_Credit_Card.csv')
        df
```

Out[3]:

| | ID | LIMIT_BAL | SEX | EDUCATION | MARRIAGE | AGE | PAY_0 | PAY_2 | PAY_3 |
|---|---|---|---|---|---|---|---|---|---|
| **0** | 1 | 20000.0 | 2 | 2 | 1 | 24 | 2 | 2 | -1 |
| **1** | 2 | 120000.0 | 2 | 2 | 2 | 26 | -1 | 2 | 0 |
| **2** | 3 | 90000.0 | 2 | 2 | 2 | 34 | 0 | 0 | 0 |
| **3** | 4 | 50000.0 | 2 | 2 | 1 | 37 | 0 | 0 | 0 |
| **4** | 5 | 50000.0 | 1 | 2 | 1 | 57 | -1 | 0 | -1 |
| **...** | ... | ... | ... | ... | ... | ... | ... | ... | ... |
| **29995** | 29996 | 220000.0 | 1 | 3 | 1 | 39 | 0 | 0 | 0 |
| **29996** | 29997 | 150000.0 | 1 | 3 | 2 | 43 | -1 | -1 | -1 |
| **29997** | 29998 | 30000.0 | 1 | 2 | 2 | 37 | 4 | 3 | 2 |
| **29998** | 29999 | 80000.0 | 1 | 3 | 1 | 41 | 1 | -1 | 0 |
| **29999** | 30000 | 50000.0 | 1 | 2 | 1 | 46 | 0 | 0 | 0 |

30000 rows × 25 columns

# 2. Data splitting

rubric={points:2}

**Your tasks:**

1. Split the data into train (70%) and test (30%) portions with `random_state=123`.

> If your computer cannot handle training on 70% training data, make the
> test split bigger.

Solution_2

*Points:* 2

In [4]:
```
train_df, test_df = train_test_split(df, test_size=0.30, random_state=123)
train_df.head()
```

Out[4]:

| | ID | LIMIT_BAL | SEX | EDUCATION | MARRIAGE | AGE | PAY_0 | PAY_2 | PAY_3 |
|---|---|---|---|---|---|---|---|---|---|
| **16395** | 16396 | 320000.0 | 2 | 1 | 2 | 36 | 0 | 0 | 0 |
| **21448** | 21449 | 440000.0 | 2 | 1 | 2 | 30 | -1 | -1 | -1 |
| **20034** | 20035 | 160000.0 | 2 | 3 | 1 | 44 | -2 | -2 | -2 |
| **25755** | 25756 | 120000.0 | 2 | 2 | 1 | 30 | 0 | 0 | 0 |
| **1438** | 1439 | 50000.0 | 1 | 2 | 2 | 54 | 1 | 2 | 0 |

5 rows × 25 columns

# 3. EDA

rubric={points:10}

**Your tasks:**

1. Perform exploratory data analysis on the train set.
2. Include at least two summary statistics and two visualizations that you find useful, and accompany each one with a sentence explaining it.
3. Summarize your initial observations about the data.
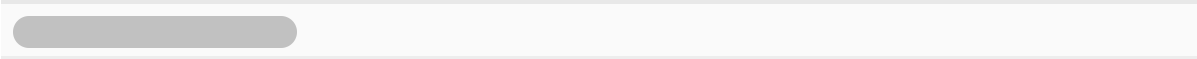4. Pick appropriate metric/metrics for assessment.

Solution_3

*Points:* 10

1.1 train_df.describe() gives us a statistic summary of all columns in the train dataframe.

In [5]:
```
train_df.describe()
```

Out[5]:

| | ID | LIMIT_BAL | SEX | EDUCATION | MARRIAGE | |
|---|---|---|---|---|---|---|
| **count** | 21000.000000 | 21000.000000 | 21000.000000 | 21000.000000 | 21000.000000 | 2' |
| **mean** | 14962.348238 | 167880.651429 | 1.600762 | 1.852143 | 1.554000 | |
| **std** | 8650.734050 | 130202.682167 | 0.489753 | 0.792961 | 0.521675 | |
| **min** | 1.000000 | 10000.000000 | 1.000000 | 0.000000 | 0.000000 | |
| **25%** | 7498.750000 | 50000.000000 | 1.000000 | 1.000000 | 1.000000 | |
| **50%** | 14960.500000 | 140000.000000 | 2.000000 | 2.000000 | 2.000000 | |
| **75%** | 22458.250000 | 240000.000000 | 2.000000 | 2.000000 | 2.000000 | |
| **max** | 30000.000000 | 1000000.000000 | 2.000000 | 6.000000 | 3.000000 | |

8 rows × 25 columns

1.2 From train_df.info(), we can see that there's no missing value in the dataframe, and we know the data size and the data types of all columns.

In [6]:
```
train_df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
Index: 21000 entries, 16395 to 19966
Data columns (total 25 columns):
 #   Column                   Non-Null Count  Dtype
---  ------                   --------------  -----
 0   ID                       21000 non-null  int64
 1   LIMIT_BAL                21000 non-null  float64
 2   SEX                      21000 non-null  int64
 3   EDUCATION                21000 non-null  int64
 4   MARRIAGE                 21000 non-null  int64
 5   AGE                      21000 non-null  int64
 6   PAY_0                    21000 non-null  int64
 7   PAY_2                    21000 non-null  int64
 8   PAY_3                    21000 non-null  int64
 9   PAY_4                    21000 non-null  int64
 10  PAY_5                    21000 non-null  int64
 11  PAY_6                    21000 non-null  int64
 12  BILL_AMT1                21000 non-null  float64
 13  BILL_AMT2                21000 non-null  float64
 14  BILL_AMT3                21000 non-null  float64
 15  BILL_AMT4                21000 non-null  float64
 16  BILL_AMT5                21000 non-null  float64
 17  BILL_AMT6                21000 non-null  float64
 18  PAY_AMT1                 21000 non-null  float64
 19  PAY_AMT2                 21000 non-null  float64
 20  PAY_AMT3                 21000 non-null  float64
 21  PAY_AMT4                 21000 non-null  float64
 22  PAY_AMT5                 21000 non-null  float64
 23  PAY_AMT6                 21000 non-null  float64
 24  default.payment.next.month  21000 non-null  int64
dtypes: float64(13), int64(12)
memory usage: 4.2 MB
```

2.1 two summary statistics

- Filter on 'AGE' column and apply mean() on it, we can know the average age of the people in the dataset is 35.5
- Similarly, we find out the median credit limit in the dataset is $140,000.

```
In [7]: mean_age = train_df['AGE'].mean()
        print("Mean Age:\n", mean_age)
```
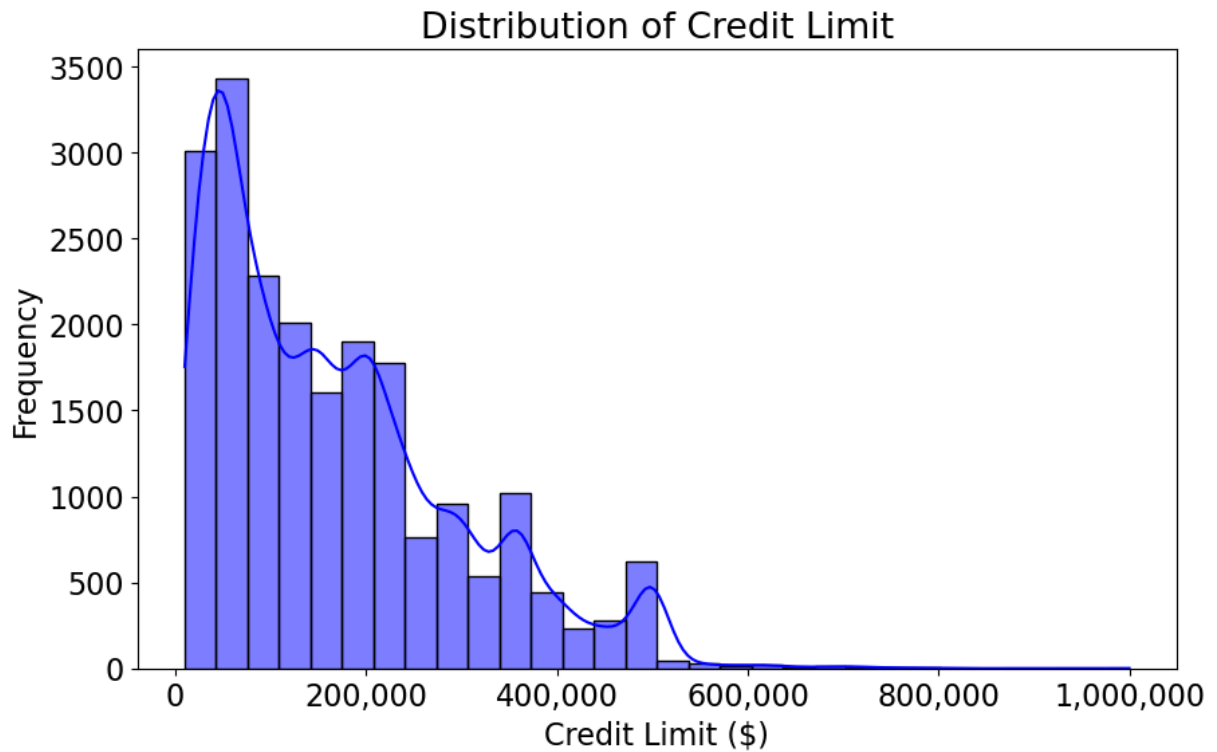
```
Mean Age:
 35.50080952380952
```

```
In [8]: median_limit_bal = train_df['LIMIT_BAL'].median()
        print("Median Credit Limit (LIMIT_BAL):\n", median_limit_bal)
```

```
Median Credit Limit (LIMIT_BAL):
 140000.0
```

2.2 This histogram shows the distribution of credit limits in the training dataset, which helps us understand the general range and concentration of credit limits.

In [9]:
```python
import seaborn as sns
import matplotlib.ticker as mtick

# Visualization 1: Distribution of 'LIMIT_BAL' (Credit Limit)
plt.figure(figsize=(10, 6))
sns.histplot(train_df['LIMIT_BAL'], kde=True, bins=30, color='blue')
plt.xlabel('Credit Limit ($)')
plt.ylabel('Frequency')
plt.title('Distribution of Credit Limit')
plt.gca().xaxis.set_major_formatter(mtick.StrMethodFormatter('{x:,.0f}'))
plt.show()
```



2.2 This count plot shows the distribution of credit card defaults across different education levels. It provides insights into which education levels are associated with higher or lower default rates.

In [10]:
```python
plt.figure(figsize=(10, 6))
sns.countplot(data=train_df, x='EDUCATION', hue='default.payment.next.month'
plt.xlabel('Education Level')
plt.ylabel('Count')
plt.title('Default Rate by Education Level')
plt.legend(title='Default Next Month', loc='upper right')
plt.show()
```

Default Rate by Education Level

3. Initial Observations

- The distribution of credit limits shows that most clients have credit limits between below 200,000.
- There are clients with very high credit limits, but these are relatively rare, indicating that credit limit data is skewed.
- The mean age of clients is approximately 35.5 years, suggesting that the dataset primarily includes middle-aged individuals.
- The median credit limit is 140,000, indicating that half of the clients have credit limits below this value.
- Default rates seem to vary by education level, with some education levels exhibiting higher default frequencies.
- There are no missing values in the dataset, which makes data preprocessing simpler.

4. In this case, if default is positive, then false negative cases(when they are default but we recognize them as not-default) are more detrimental. We should use Recall as our metric because Recall measures the proportion of actual positives that are correctly identified by the model. A high recall indicates that the model is catching most of the true positives (default recognized correctly), minimizing the number of false negatives.

In [11]:
```python
from sklearn.metrics import recall_score, make_scorer
recall_scorer = make_scorer(recall_score, greater_is_better=True)
```

# 4. Feature engineering

rubric={points:1}

**Your tasks:**

1. Carry out feature engineering. In other words, extract new features relevant for the problem and work with your new feature set in the following exercises. You may have to go back and forth between feature engineering and preprocessing.

> Solution_4

*Points:* 1

```
In [12]:  # 1. Creating Worst payment delay
          payment_cols = ['PAY_0', 'PAY_2', 'PAY_3', 'PAY_4', 'PAY_5', 'PAY_6']
          train_df['worst_payment_delay'] = train_df[payment_cols].max(axis=1)

          # 2. Creating Bill Amount Features
          bill_cols = ['BILL_AMT1', 'BILL_AMT2', 'BILL_AMT3', 'BILL_AMT4', 'BILL_AMT5'
           # Average bill amount
          train_df['avg_bill_amt'] = train_df[bill_cols].mean(axis=1)

          # 3. Creating Payment Amount Features
          pay_cols = ['PAY_AMT1', 'PAY_AMT2', 'PAY_AMT3', 'PAY_AMT4', 'PAY_AMT5', 'PAY
          # Average payment amount
          train_df['avg_payment_amt'] = train_df[pay_cols].mean(axis=1)

          # 4. Payment Ratio Features

          # Calculate payment ratios for each month
          for i in range(1, 7):
              bill_col = f'BILL_AMT{i}'
              pay_col = f'PAY_AMT{i}'
              ratio_col = f'payment_ratio_{i}'
              train_df[ratio_col] = train_df[pay_col] / train_df[bill_col].replace(0,

          # Average payment ratio
          ratio_cols = [f'payment_ratio_{i}' for i in range(1, 7)]
          train_df['avg_payment_ratio'] = train_df[ratio_cols].mean(axis=1)
          # Only keep average payment ratio
          train_df = train_df.drop(columns=ratio_cols)

          # 5. Credit Utilization Features

          # Average credit utilization over 6 months
          train_df['avg_credit_utilization'] = (train_df[bill_cols].mean(axis=1) / tra
```

```
train_df
```

Out[12]:

| | ID | LIMIT_BAL | SEX | EDUCATION | MARRIAGE | AGE | PAY_0 | PAY_2 | PAY_3 |
|---|---|---|---|---|---|---|---|---|---|
| **16395** | 16396 | 320000.0 | 2 | 1 | 2 | 36 | 0 | 0 | 0 |
| **21448** | 21449 | 440000.0 | 2 | 1 | 2 | 30 | -1 | -1 | -1 |
| **20034** | 20035 | 160000.0 | 2 | 3 | 1 | 44 | -2 | -2 | -2 |
| **25755** | 25756 | 120000.0 | 2 | 2 | 1 | 30 | 0 | 0 | 0 |
| **1438** | 1439 | 50000.0 | 1 | 2 | 2 | 54 | 1 | 2 | 0 |
| **...** | ... | ... | ... | ... | ... | ... | ... | ... | ... |
| **28636** | 28637 | 380000.0 | 2 | 2 | 1 | 37 | 0 | 0 | 0 |
| **17730** | 17731 | 360000.0 | 2 | 1 | 1 | 54 | 1 | -2 | -2 |
| **28030** | 28031 | 50000.0 | 2 | 3 | 1 | 29 | 0 | 0 | 0 |
| **15725** | 15726 | 30000.0 | 2 | 2 | 2 | 21 | 0 | 0 | 0 |
| **19966** | 19967 | 370000.0 | 2 | 1 | 1 | 36 | -2 | -2 | -2 |

21000 rows × 30 columns

In [13]:
```python
# 1. Creating Worst payment delay
payment_cols = ['PAY_0', 'PAY_2', 'PAY_3', 'PAY_4', 'PAY_5', 'PAY_6']
test_df['worst_payment_delay'] = test_df[payment_cols].max(axis=1)

# 2. Creating Bill Amount Features
bill_cols = ['BILL_AMT1', 'BILL_AMT2', 'BILL_AMT3', 'BILL_AMT4', 'BILL_AMT5'
 # Average bill amount
test_df['avg_bill_amt'] = test_df[bill_cols].mean(axis=1)

# 3. Creating Payment Amount Features
pay_cols = ['PAY_AMT1', 'PAY_AMT2', 'PAY_AMT3', 'PAY_AMT4', 'PAY_AMT5', 'PAY
# Average payment amount
test_df['avg_payment_amt'] = test_df[pay_cols].mean(axis=1)

# 4. Payment Ratio Features

# Calculate payment ratios for each month
for i in range(1, 7):
    bill_col = f'BILL_AMT{i}'
    pay_col = f'PAY_AMT{i}'
    ratio_col = f'payment_ratio_{i}'
    test_df[ratio_col] = test_df[pay_col] / test_df[bill_col].replace(0, 1)

# Average payment ratio
ratio_cols = [f'payment_ratio_{i}' for i in range(1, 7)]
test_df['avg_payment_ratio'] = test_df[ratio_cols].mean(axis=1)
# Only keep average payment ratio
```

```
test_df = test_df.drop(columns=ratio_cols)

# 5. Credit Utilization Features

# Average credit utilization over 6 months
test_df['avg_credit_utilization'] = (test_df[bill_cols].mean(axis=1) / test_

test_df
```

Out[13]:

| | ID | LIMIT_BAL | SEX | EDUCATION | MARRIAGE | AGE | PAY_0 | PAY_2 | PAY_3 |
|---|---|---|---|---|---|---|---|---|---|
| 25665 | 25666 | 40000.0 | 2 | 2 | 2 | 26 | -1 | 0 | 0 |
| 16464 | 16465 | 80000.0 | 2 | 3 | 1 | 59 | 0 | 0 | 0 |
| 22386 | 22387 | 170000.0 | 2 | 1 | 2 | 30 | 2 | 2 | 2 |
| 10149 | 10150 | 200000.0 | 2 | 2 | 1 | 41 | -2 | -2 | -2 |
| 8729 | 8730 | 50000.0 | 1 | 2 | 1 | 43 | 0 | 0 | 0 |
| ... | ... | ... | ... | ... | ... | ... | ... | ... | ... |
| 17548 | 17549 | 60000.0 | 2 | 2 | 1 | 48 | 0 | 0 | 0 |
| 11459 | 11460 | 310000.0 | 1 | 2 | 1 | 43 | -1 | -1 | -1 |
| 6608 | 6609 | 10000.0 | 2 | 2 | 2 | 22 | 0 | 0 | 0 |
| 2414 | 2415 | 30000.0 | 1 | 2 | 1 | 38 | 1 | -1 | -1 |
| 14757 | 14758 | 30000.0 | 2 | 3 | 1 | 24 | 2 | 0 | 0 |

9000 rows × 30 columns

# 5. Preprocessing and transformations

rubric={points:10}

**Your tasks:**

1. Identify different feature types and the transformations you would apply on each feature type.
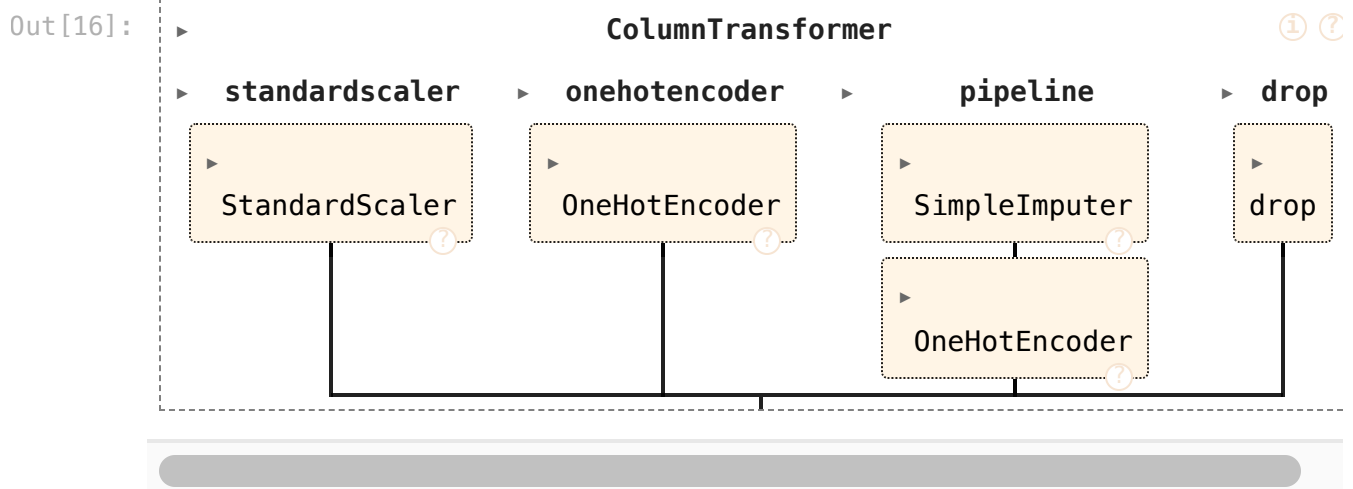2. Define a column transformer, if necessary.

> Solution_5

*Points:* 10

```
In [14]: numeric_features = ['LIMIT_BAL', 'AGE', 'BILL_AMT1', 'BILL_AMT2', 'BILL_AMT3
                             'PAY_0', 'PAY_2', 'PAY_3', 'PAY_4', 'PAY_5', 'PAY_6',
                             'worst_payment_delay', 'avg_bill_amt','avg_payment_amt',
         categorical_features = ['EDUCATION','MARRIAGE']
         binary_features = ['SEX']
         drop_features = ['ID']
         target = "default.payment.next.month"
```

```
In [15]: numeric_transformer = StandardScaler()
         binary_transformer = OneHotEncoder(handle_unknown="ignore", sparse_output=Fa
         categorical_transformer = make_pipeline(
             SimpleImputer(strategy="most_frequent"), OneHotEncoder(handle_unknown="i
         )
```

```
In [16]: preprocessor = make_column_transformer(
             (numeric_transformer, numeric_features),
             (binary_transformer, binary_features),
             (categorical_transformer, categorical_features),
             ('drop', drop_features),
         )
         preprocessor
```

Out[16]:



# 6. Baseline model

rubric={points:2}

**Your tasks:**

1. Try `scikit-learn` 's baseline model and report results.

Solution_6

*Points:* 2

```
In [17]: X_train = train_df.drop(["default.payment.next.month"], axis=1)
         y_train = train_df["default.payment.next.month"]
         X_test = test_df.drop(["default.payment.next.month"], axis=1)
         y_test = test_df["default.payment.next.month"]
```

```
In [18]: def mean_std_cross_val_scores(model, X_train, y_train, **kwargs):
             """
             Returns mean and std of cross validation

             Parameters
             ----------
             model :
                 scikit-learn model
             X_train : numpy array or pandas DataFrame
                 X in the training data
             y_train :
                 y in the training data

             Returns
             ----------
                 pandas Series with mean scores from cross_validation
             """

             scores = cross_validate(model, X_train, y_train, **kwargs)

             mean_scores = pd.DataFrame(scores).mean()
             std_scores = pd.DataFrame(scores).std()
             out_col = []

             for i in range(len(mean_scores)):
                 out_col.append((f"%0.3f (+/- %0.3f)" % (mean_scores.iloc[i], std_sco

             return pd.Series(data=out_col, index=mean_scores.index)
```

```
In [19]: dummy_clf = DummyClassifier(strategy="stratified")
         dummy_scores = mean_std_cross_val_scores(dummy_clf, X_train, y_train, return
         dummy_scores = pd.DataFrame(dummy_scores).T
         dummy_scores.index = ['dummy']
         dummy_scores
```

Out[19]:

|  | fit_time | score_time | test_score | train_score |
|---|---|---|---|---|
| **dummy** | 0.004 (+/- 0.001) | 0.004 (+/- 0.002) | 0.235 (+/- 0.011) | 0.221 (+/- 0.005) |

# 7. Linear models

rubric={points:10}

**Your tasks:**

1. Try a linear model as a first real attempt.
2. Carry out hyperparameter tuning to explore different values for the complexity hyperparameter.
3. Report cross-validation scores along with standard deviation.
4. Summarize your results.

> Solution_7

*Points:* 10

1. We use a logistic regression model with balanced class weights as our first real attempt.

```
In [20]:   pipe_lr = make_pipeline(
               preprocessor, LogisticRegression(max_iter=2000)
           )
           pipe_lr.fit(X_train, y_train)
           ConfusionMatrixDisplay.from_estimator(
               pipe_lr,
               X_train,
               y_train,
               display_labels=["Non default", "default"],
               values_format="d",
           );
```

```
In [21]: lr_results = mean_std_cross_val_scores(pipe_lr, X_train, y_train, return_tra
         lr_results = pd.DataFrame(lr_results).T
         lr_results.index = ['Logistic regression']
         lr_results
```

Out[21]:

| | fit_time | score_time | test_score | train_score |
|---|---|---|---|---|
| **Logistic regression** | 0.189 (+/- 0.121) | 0.017 (+/- 0.004) | 0.238 (+/- 0.023) | 0.236 (+/- 0.012) |

```
In [22]: pipe_lr_balanced = make_pipeline(
             preprocessor, LogisticRegression(max_iter=2000, class_weight="balanced")
         )
         pipe_lr_balanced.fit(X_train, y_train)
         ConfusionMatrixDisplay.from_estimator(
             pipe_lr_balanced,
             X_train,
             y_train,
             display_labels=["Non default", "default"],
             values_format="d",
         );
```
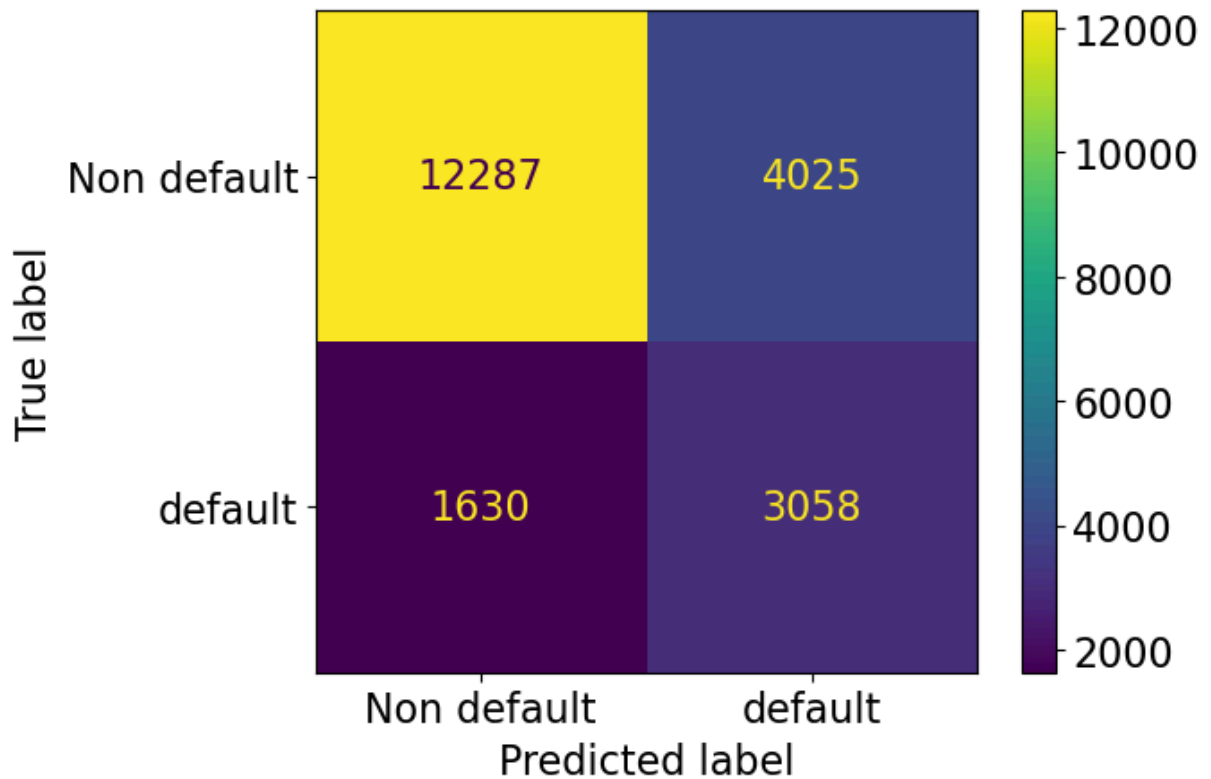
```
In [23]: lr_balanced_results = mean_std_cross_val_scores(pipe_lr_balanced, X_train, y
         lr_balanced_results = pd.DataFrame(lr_balanced_results).T
         lr_balanced_results.index = ['Logistic regression_Balanced']
         lr_balanced_results
```

Out[23]:

|  | fit_time | score_time | test_score | train_score |
|---|---|---|---|---|
| **Logistic regression_Balanced** | 0.125 (+/- 0.011) | 0.015 (+/- 0.002) | 0.653 (+/- 0.011) | 0.653 (+/- 0.003) |

2. Hyperparameter Optimization

```
In [24]: from scipy.stats import expon, lognorm, loguniform, randint, uniform, norm,

         param_grid = {
             "logisticregression__C": loguniform(1e-3, 1e3),
         }

         # Create a random search object
         random_search = RandomizedSearchCV(pipe_lr_balanced,
                         param_distributions = param_grid,
                         n_iter=100,
                         n_jobs=-1,
                         return_train_score=True,
                         scoring = recall_scorer)

         # Carry out the search
         random_search.fit(X_train, y_train)
```

Out[24]:

```
┌─────────────────────────────────────────────────────────────────────┐
│  ▸                    RandomizedSearchCV                             │
│                                                                       │
│  ▸                  best_estimator_: Pipeline                        │
│  ┌──────────────────────────────────────────────────────────────┐   │
│  │  ▸            columntransformer: ColumnTransformer            │   │
│  │  ▸   standardscaler    ▸   onehotencoder    ▸   pipeline   ▸ d│   │
│  │  ┌──────────────┐  ┌──────────────┐  ┌──────────────┐  ┌──── │   │
│  │  │ ▸            │  │ ▸            │  │ ▸            │  │ ▸    │   │
│  │  │ StandardScaler│  │OneHotEncoder │  │SimpleImputer │  │ d   │   │
│  │  └──────────────⊙  └──────────────⊙  └──────────────⊙  └──── │   │
│  │                                      ┌──────────────┐          │   │
│  │                                      │ ▸            │          │   │
│  │                                      │OneHotEncoder │          │   │
│  │                                      └──────────────⊙          │   │
│  │          │              │                     │                │   │
│  │        ┌─────────────────────────────────┐                    │   │
│  │        │ ▸  LogisticRegression    ⊙     │                    │   │
│  │        └─────────────────────────────────┘                    │   │
│  └──────────────────────────────────────────────────────────────┘   │
└─────────────────────────────────────────────────────────────────────┘
```

In [25]:

```
random_search.best_score_
```

Out[25]:   0.6537923281898178

In [26]:

```
random_search.best_params_
```

Out[26]:   {'logisticregression__C': 0.533237852400314}

3. Report cross-validation scores along with standard deviation.

In [27]:

```
lr_balanced_results = mean_std_cross_val_scores(pipe_lr_balanced, X_train, y
lr_balanced_results = pd.DataFrame(lr_balanced_results).T
lr_balanced_results.index = ['Logistic regression_Balanced']
lr_balanced_results
```
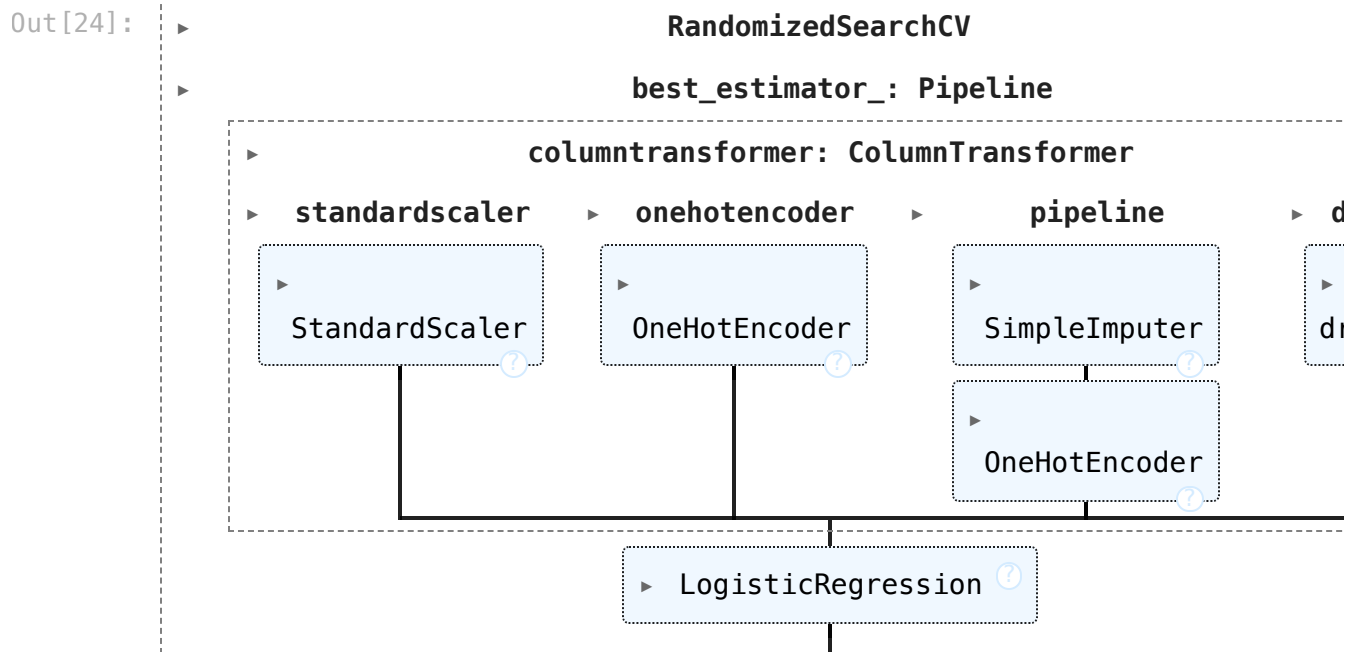
Out[27]:

|  | fit_time | score_time | test_score | train_score |
|---|---|---|---|---|
| **Logistic regression_Balanced** | 0.200 (+/-0.096) | 0.032 (+/-0.022) | 0.653 (+/-0.011) | 0.653 (+/-0.003) |

In [28]:

```
pipe_lr_balanced_tuned = random_search.best_estimator_

mean_std_scores = mean_std_cross_val_scores(pipe_lr_balanced_tuned, X_train,

lr_balanced_tuned = pd.DataFrame(mean_std_scores).T
lr_balanced_tuned.index = ['LR_Balanced_Tuned']
lr_balanced_tuned
```

| | fit_time | score_time | test_score | train_score |
|---|---|---|---|---|
| **LR_Balanced_Tuned** | 0.120 (+/-0.003) | 0.015 (+/-0.001) | 0.654 (+/-0.012) | 0.653 (+/-0.003) |

4. Summary:

- Because there is a class imbalance so we used logistic regression with equal class wight to find out that the result of logistic regression model (average recall cv score = 0.653) is better than dummy classifier (cv = 0.222).
- With hyperparameter optimization, we find best cv score is 0.654, which is almost the same as the model using default C value. This might be caused by the insensitivity to C of the model.
- cv score and train score are low and close in value, indicating underfitting in logistic regression model.

# 8. Different models

rubric={points:12}
**Your tasks:**

1. Try at least 3 other models aside from a linear model. One of these models should be a tree-based ensemble model.
2. Summarize your results in terms of overfitting/underfitting and fit and score times. Can you beat a linear model?

> Solution_8

*Points:* 12

## Overfitting/Underfitting:

- **Random Forest**: The model is **overfitting**. Its training score is perfect (1.000), indicating that it fits the training data almost too well. However, the test score (0.341) is significantly lower, which means it does not generalize well to unseen data. This suggests overfitting.

- **K-Nearest Neighbors (KNN)**: The KNN model shows a relatively low training score (0.472) and an even lower test score (0.363). While the test score being lower than

the training score typically suggests overfitting, in this case, the low training score indicates that the model is **underfitting** because of the high training error. This issue may caused by the class imbalance in the dataset. KNN might struggle more because it's a distance-based algorithm. With imbalanced data, KNN tends to get biased toward the majority class, making it harder for it to generalize well. This imbalance can cause the model to perform poorly on minority classes, further contributing to its underfitting.

- **LightGBM**: There is a moderate case of **underfitting**. The training score is 0.457 and the test score is 0.375, which are both similar but fairly low. This suggests that the model is not learning enough complexity from the training data, possibly underfitting.

- **LightGBM (balanced)**: This model seems to achieve better **balance**, as it avoids extreme overfitting. The test score (0.620) is relatively high, and the training score (0.788) indicates a good fit on the training data without being perfect. This makes it the most generalizable model in this set.

## Fit and Score Times:

- **Random Forest** has a much higher fit time (2.929 seconds) compared to the other models, suggesting it is computationally expensive. The score time (0.043 seconds) is comparable to KNN, but still quite fast.

- **KNN** has the fastest fit time (0.010 seconds), making it computationally efficient. Its score time (0.050 seconds) is slightly higher than LightGBM but still very quick.

- **LightGBM** is very efficient with a fit time of 0.261 seconds, much lower than Random Forest. Its score time (0.008 seconds) is also the quickest among all the models.

- **LightGBM (balanced)** has a very similar fit time (0.256 seconds) and score time (0.007 seconds) to the unbalanced version, making it equally efficient in computation.

## Comparison to Logistic Regression:

- **Logistic Regression (0.238)**: This is lower than all models except KNN, indicating that KNN slightly outperforms logistic regression but still underfits.

- **Balanced Logistic Regression (0.653)**: This model outperforms **all** models except the balanced LightGBM, which comes close at 0.620. The balanced logistic regression still holds a slight edge.

## Conclusion:

- We are able to beat the **vanilla logistic regression** with the LightGBM models in terms of cv score. However, the **balanced logistic regression** still outperforms the LightGBM models slightly.
- Among the tree-based models, **LightGBM (balanced)** seems to provide the best trade-off between fit, generalization, and computational efficiency.

```
In [29]:   from lightgbm.sklearn import LGBMClassifier
           from sklearn.ensemble import RandomForestClassifier
           from sklearn.ensemble import GradientBoostingClassifier, HistGradientBoostin


           pipe_rf = make_pipeline(
               preprocessor, RandomForestClassifier(class_weight="balanced", random_sta
           )
           pipe_knn = make_pipeline(
               preprocessor, KNeighborsClassifier()
           )
           pipe_lgbm = make_pipeline(
               preprocessor, LGBMClassifier(random_state=123, verbose=-1)
           )
           pipe_lgbm_balanced = make_pipeline(
               preprocessor, LGBMClassifier(random_state=123, verbose=-1, class_weight=
           )

           classifiers = {
               "random forest": pipe_rf,
               'K-Nearest Neighbors': pipe_knn,
               "LightGBM": pipe_lgbm,
               "LightGBM_balanced": pipe_lgbm_balanced
           }
```

```
In [30]:   import warnings

           warnings.simplefilter(action="ignore", category=DeprecationWarning)
           warnings.simplefilter(action="ignore", category=UserWarning)
```

```
In [31]:   results = {}
           for (name, model) in classifiers.items():
               results[name] = mean_std_cross_val_scores(
                   model, X_train, y_train, return_train_score=True, scoring = recall_s
               )
```

```
In [32]:   rf_knn_lgbm_scores = pd.DataFrame(results).T
           rf_knn_lgbm_scores
```

Out[32]:

|  | fit_time | score_time | test_score | train_score |
|---|---|---|---|---|
| **random forest** | 6.216 (+/- 0.447) | 0.089 (+/- 0.009) | 0.343 (+/- 0.018) | 1.000 (+/- 0.000) |
| **K-Nearest Neighbors** | 0.040 (+/- 0.005) | 0.167 (+/- 0.083) | 0.363 (+/- 0.005) | 0.472 (+/- 0.005) |
| **LightGBM** | 0.245 (+/- 0.020) | 0.029 (+/- 0.004) | 0.375 (+/- 0.013) | 0.457 (+/- 0.004) |
| **LightGBM_balanced** | 0.273 (+/- 0.036) | 0.028 (+/- 0.004) | 0.620 (+/- 0.021) | 0.788 (+/- 0.007) |

# 9. Feature selection

rubric={points:2}

**Your tasks:**

Make some attempts to select relevant features. You may try `RFECV` or forward selection for this. Do the results improve with feature selection? Summarize your results. If you see improvements in the results, keep feature selection in your pipeline. If not, you may abandon it in the next exercises.

> Solution_9

*Points:* 2

The results does not really improve with feature selection, as the balanced logistic regression model with selected features (0.654) has the same test score as the one without selected features (0.654). But the fit time is faster. This may be caused by the small number of features.

In [33]:
```python
from sklearn.feature_selection import RFE

# Get preprocessed data from the tuned pipeline
X_train_processed = pipe_lr_balanced_tuned.named_steps['columntransformer'].
X_test_processed = pipe_lr_balanced_tuned.named_steps['columntransformer'].t

# Get feature names from the tuned pipeline
numeric_feature_names = numeric_features
binary_feature_names = pipe_lr_balanced_tuned.named_steps['columntransformer
categorical_feature_names = pipe_lr_balanced_tuned.named_steps['columntransf
all_feature_names = list(itertools.chain(numeric_feature_names, binary_featu
```

```python
rfe_cv = RFECV(
    estimator=pipe_lr_balanced_tuned.named_steps['logisticregression'],
    cv=10,
    scoring=recall_scorer,
    min_features_to_select= 5
)

# Fit RFECV
rfe_cv.fit(X_train_processed, y_train)

# Get selected features
selected_features_mask = rfe_cv.support_
selected_feature_names = [name for name, selected in zip(all_feature_names,

# Create DataFrames with selected features
X_train_selected_df = pd.DataFrame(
    X_train_processed[:, selected_features_mask],
    columns=selected_feature_names
)
X_test_selected_df = pd.DataFrame(
    X_test_processed[:, selected_features_mask],
    columns=selected_feature_names
)

selected_model = LogisticRegression(
    max_iter=2000,
    class_weight="balanced",
    C=pipe_lr_balanced_tuned.named_steps['logisticregression'].C  # Use the
)
mean_std_scores_selected = mean_std_cross_val_scores(
    selected_model,
    X_train_selected_df,
    y_train,
    cv=5,
    return_train_score=True,
    scoring=recall_scorer
)

lr_scores_bal_selected = pd.DataFrame(mean_std_scores_selected).T
lr_scores_bal_selected.index = ['LR_Balanced_Selected']
lr_scores_bal_selected
```

Out[33]:

|                      | fit_time          | score_time        | test_score        | train_score       |
|----------------------|-------------------|-------------------|-------------------|-------------------|
| **LR_Balanced_Selected** | 0.053 (+/- 0.003) | 0.005 (+/- 0.001) | 0.650 (+/- 0.013) | 0.650 (+/- 0.003) |

```python
In [34]:  # Create a pipeline that combines preprocessing and selected features with t
pipe_lr_selected = Pipeline([
    ('columntransformer', pipe_lr_balanced_tuned.named_steps['columntransfor
    ('feature_selector', RFECV(
        estimator=LogisticRegression(max_iter=2000, class_weight="balanced",
        cv=10,
        scoring=recall_scorer,
        min_features_to_select=5
```
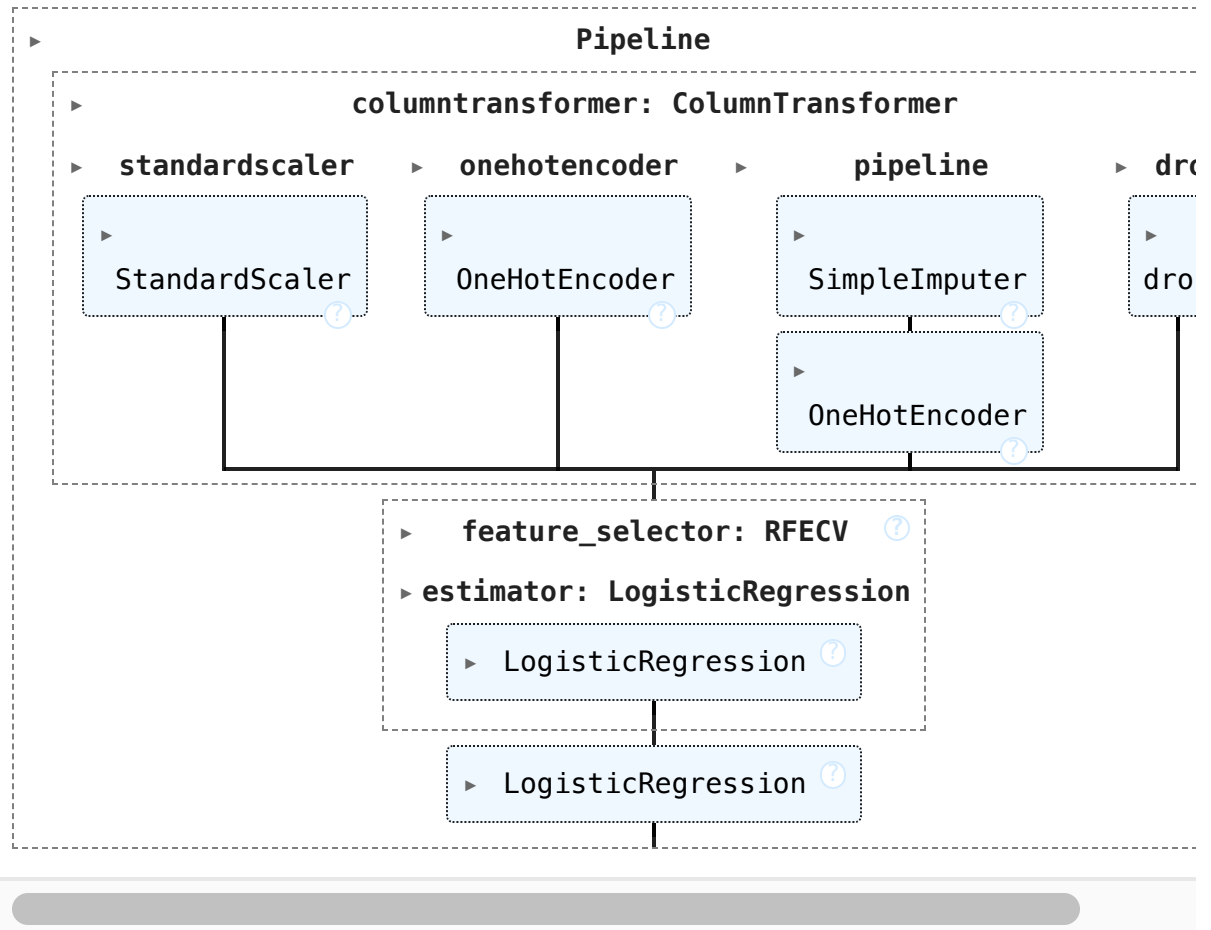
```
    )),
    ('logisticregression', LogisticRegression(
        max_iter=2000,
        class_weight="balanced",
        C=pipe_lr_balanced_tuned.named_steps['logisticregression'].C
    ))
])

# Fit the pipeline
pipe_lr_selected.fit(X_train, y_train)
```

Out[34]:

> **Pipeline**

> **columntransformer: ColumnTransformer**

> **standardscaler**    > **onehotencoder**    > **pipeline**    > **dr**

> StandardScaler      > OneHotEncoder      > SimpleImputer    > dro

> OneHotEncoder

> **feature_selector: RFECV** ⑦

> **estimator: LogisticRegression**

> LogisticRegression ⑦

> LogisticRegression ⑦

In [35]:
```
combined_scores = pd.concat([dummy_scores, lr_results, lr_balanced_results,
combined_scores
```

Out[35]:

|  | fit_time | score_time | test_score | train_score |
|---|---|---|---|---|
| **dummy** | 0.004 (+/- 0.001) | 0.004 (+/- 0.002) | 0.235 (+/- 0.011) | 0.221 (+/- 0.005) |
| **Logistic regression** | 0.189 (+/- 0.121) | 0.017 (+/- 0.004) | 0.238 (+/- 0.023) | 0.236 (+/- 0.012) |
| **Logistic regression_Balanced** | 0.200 (+/- 0.096) | 0.032 (+/- 0.022) | 0.653 (+/- 0.011) | 0.653 (+/- 0.003) |
| **LR_Balanced_Tuned** | 0.120 (+/- 0.003) | 0.015 (+/- 0.001) | 0.654 (+/- 0.012) | 0.653 (+/- 0.003) |
| **random forest** | 6.216 (+/- 0.447) | 0.089 (+/- 0.009) | 0.343 (+/- 0.018) | 1.000 (+/- 0.000) |
| **K-Nearest Neighbors** | 0.040 (+/- 0.005) | 0.167 (+/- 0.083) | 0.363 (+/- 0.005) | 0.472 (+/- 0.005) |
| **LightGBM** | 0.245 (+/- 0.020) | 0.029 (+/- 0.004) | 0.375 (+/- 0.013) | 0.457 (+/- 0.004) |
| **LightGBM_balanced** | 0.273 (+/- 0.036) | 0.028 (+/- 0.004) | 0.620 (+/- 0.021) | 0.788 (+/- 0.007) |
| **LR_Balanced_Selected** | 0.053 (+/- 0.003) | 0.005 (+/- 0.001) | 0.650 (+/- 0.013) | 0.650 (+/- 0.003) |

# 10. Hyperparameter optimization

rubric={points:10}
**Your tasks:**

Make some attempts to optimize hyperparameters for the models you've tried and summarize your results. In at least one case you should be optimizing multiple hyperparameters for a single model. You may use `sklearn` 's methods for hyperparameter optimization or fancier Bayesian optimization methods.

- GridSearchCV
- RandomizedSearchCV
- scikit-optimize

Solution_10

*Points:* 10

*Type your answer here, replacing this text.*

In [36]: `X_train_selected_df.head()`

Out[36]:

|   | BILL_AMT1 | BILL_AMT2 | BILL_AMT3 | PAY_0 | PAY_AMT1 | PAY_AMT2 | PAY_AMT3 |
|---|-----------|-----------|-----------|-------|----------|----------|----------|
| 0 | -0.300665 | -0.293394 | -0.265310 | 0.013770 | -0.039546 | -0.040229 | -0.234603 |
| 1 | -0.685307 | -0.679495 | 0.585444 | -0.878738 | -0.297166 | 3.739796 | 6.785208 |
| 2 | -0.696132 | -0.688319 | -0.681234 | -1.771246 | -0.333097 | -0.270403 | -0.289017 |
| 3 | 0.687456 | 0.752583 | 0.835581 | 0.013770 | -0.115517 | -0.018028 | -0.060260 |
| 4 | -0.040230 | -0.031399 | -0.287429 | 0.906278 | -0.333097 | -0.206185 | -0.223720 |

5 rows × 22 columns

In [37]:
```python
#Random Forest Hyperparameter Optimization
#With selected features
from sklearn.pipeline import make_pipeline

pipe_rf_1 = make_pipeline(
    RandomForestClassifier(class_weight="balanced", random_state=123)
)

param_distributions_rf = {
    'randomforestclassifier__n_estimators': randint(50, 500),  # Randomly se
    'randomforestclassifier__max_depth': randint(5, 50)  # Randomly select m
}

# Create and fit GridSearchCV
rd_rf = RandomizedSearchCV(pipe_rf_1,
                param_distributions=param_distributions_rf,
                cv=5,
                n_jobs=-1,
                return_train_score=True,
                scoring = recall_scorer)

# Fit the pipeline on the original training data (not the transformed one)
rd_rf.fit(X_train_selected_df, y_train)
print("Best parameters for Random Forest:", rd_rf.best_params_)
print("Best cross-validation score for Random Forest:", rd_rf.best_score_)
```

Best parameters for Random Forest: {'randomforestclassifier__max_depth': 13, 'randomforestclassifier__n_estimators': 94}
Best cross-validation score for Random Forest: 0.5343395084343492

In [38]:
```python
# Step 1: Extract the best parameters
best_params = rd_rf.best_params_

# Create a new RandomForestClassifier using the best parameters
best_rf_selected = RandomForestClassifier(
    n_estimators=best_params['randomforestclassifier__n_estimators'],
```

```
        max_depth=best_params['randomforestclassifier__max_depth'],
        class_weight="balanced",
        random_state=123
)


mean_std_scores = mean_std_cross_val_scores(best_rf_selected, X_train_select

random_forest_selected = pd.DataFrame(mean_std_scores).T
random_forest_selected.index = ['Random Forest Selected (Tuned)']
random_forest_selected
combined_scores = pd.concat([combined_scores, random_forest_selected])
combined_scores
```

Out[38]:

| | fit_time | score_time | test_score | train_score |
|---|---|---|---|---|
| **dummy** | 0.004 (+/- 0.001) | 0.004 (+/- 0.002) | 0.235 (+/- 0.011) | 0.221 (+/- 0.005) |
| **Logistic regression** | 0.189 (+/- 0.121) | 0.017 (+/- 0.004) | 0.238 (+/- 0.023) | 0.236 (+/- 0.012) |
| **Logistic regression_Balanced** | 0.200 (+/- 0.096) | 0.032 (+/- 0.022) | 0.653 (+/- 0.011) | 0.653 (+/- 0.003) |
| **LR_Balanced_Tuned** | 0.120 (+/- 0.003) | 0.015 (+/- 0.001) | 0.654 (+/- 0.012) | 0.653 (+/- 0.003) |
| **random forest** | 6.216 (+/- 0.447) | 0.089 (+/- 0.009) | 0.343 (+/- 0.018) | 1.000 (+/- 0.000) |
| **K-Nearest Neighbors** | 0.040 (+/- 0.005) | 0.167 (+/- 0.083) | 0.363 (+/- 0.005) | 0.472 (+/- 0.005) |
| **LightGBM** | 0.245 (+/- 0.020) | 0.029 (+/- 0.004) | 0.375 (+/- 0.013) | 0.457 (+/- 0.004) |
| **LightGBM_balanced** | 0.273 (+/- 0.036) | 0.028 (+/- 0.004) | 0.620 (+/- 0.021) | 0.788 (+/- 0.007) |
| **LR_Balanced_Selected** | 0.053 (+/- 0.003) | 0.005 (+/- 0.001) | 0.650 (+/- 0.013) | 0.650 (+/- 0.003) |
| **Random Forest Selected (Tuned)** | 2.612 (+/- 0.171) | 0.049 (+/- 0.003) | 0.534 (+/- 0.018) | 0.796 (+/- 0.004) |

In [39]:
```python
#Random Forest Hyperparameter Optimization
#Without selected features
from sklearn.pipeline import make_pipeline

pipe_rf = make_pipeline(
    preprocessor, RandomForestClassifier(class_weight="balanced", random_sta
)

param_distributions_rf = {
    'randomforestclassifier__n_estimators': randint(50, 500),  # Randomly se
    'randomforestclassifier__max_depth': randint(5, 50)  # Randomly select m
}
```

```python
# Create and fit GridSearchCV
rd_rf = RandomizedSearchCV(pipe_rf,
                           param_distributions=param_distributions_rf,
                           cv=5,
                           n_jobs=-1,
                           return_train_score=True,
                           scoring = recall_scorer)

# Fit the pipeline on the original training data (not the transformed one)
rd_rf.fit(X_train, y_train)
print("Best parameters for Random Forest:", rd_rf.best_params_)
print("Best cross-validation score for Random Forest:", rd_rf.best_score_)
```

```
Best parameters for Random Forest: {'randomforestclassifier__max_depth': 5,
'randomforestclassifier__n_estimators': 455}
Best cross-validation score for Random Forest: 0.631396759152856
```

In [40]:
```python
# Create a new RandomForestClassifier using the best parameters
best_rf_unselected = rd_rf.best_estimator_

mean_std_scores = mean_std_cross_val_scores(best_rf_unselected, X_train, y_t

random_forest_unselected = pd.DataFrame(mean_std_scores).T
random_forest_unselected.index = ['Random Forest Unselected (Tuned)']
random_forest_unselected
combined_scores = pd.concat([combined_scores, random_forest_unselected])
combined_scores
```

Out[40]:

| | fit_time | score_time | test_score | train_score |
|---|---|---|---|---|
| **dummy** | 0.004 (+/- 0.001) | 0.004 (+/- 0.002) | 0.235 (+/- 0.011) | 0.221 (+/- 0.005) |
| **Logistic regression** | 0.189 (+/- 0.121) | 0.017 (+/- 0.004) | 0.238 (+/- 0.023) | 0.236 (+/- 0.012) |
| **Logistic regression_Balanced** | 0.200 (+/- 0.096) | 0.032 (+/- 0.022) | 0.653 (+/- 0.011) | 0.653 (+/- 0.003) |
| **LR_Balanced_Tuned** | 0.120 (+/- 0.003) | 0.015 (+/- 0.001) | 0.654 (+/- 0.012) | 0.653 (+/- 0.003) |
| **random forest** | 6.216 (+/- 0.447) | 0.089 (+/- 0.009) | 0.343 (+/- 0.018) | 1.000 (+/- 0.000) |
| **K-Nearest Neighbors** | 0.040 (+/- 0.005) | 0.167 (+/- 0.083) | 0.363 (+/- 0.005) | 0.472 (+/- 0.005) |
| **LightGBM** | 0.245 (+/- 0.020) | 0.029 (+/- 0.004) | 0.375 (+/- 0.013) | 0.457 (+/- 0.004) |
| **LightGBM_balanced** | 0.273 (+/- 0.036) | 0.028 (+/- 0.004) | 0.620 (+/- 0.021) | 0.788 (+/- 0.007) |
| **LR_Balanced_Selected** | 0.053 (+/- 0.003) | 0.005 (+/- 0.001) | 0.650 (+/- 0.013) | 0.650 (+/- 0.003) |
| **Random Forest Selected (Tuned)** | 2.612 (+/- 0.171) | 0.049 (+/- 0.003) | 0.534 (+/- 0.018) | 0.796 (+/- 0.004) |
| **Random Forest Unselected (Tuned)** | 10.612 (+/- 1.322) | 0.110 (+/- 0.014) | 0.631 (+/- 0.013) | 0.642 (+/- 0.002) |

In [41]:
```python
#KNN Hyperparameter Optimization
#With selected features
pipe_knn_1 = make_pipeline(
    KNeighborsClassifier()
)

param_distributions_knn = {
    'kneighborsclassifier__n_neighbors': np.arange(1, 20, 2)
}

rs_knn = RandomizedSearchCV(pipe_knn_1,
                            param_distributions=param_distributions_knn,
                            n_iter=10,
                            cv=5,
                            n_jobs=-1,
                            random_state=123,
                            return_train_score=True,
                            scoring = recall_scorer)
rs_knn.fit(X_train_selected_df, y_train)

print("Best parameters for K-Nearest Neighbors:", rs_knn.best_params_)
print("Best cross-validation score for K-Nearest Neighbors:", rs_knn.best_sc
```

Best parameters for K–Nearest Neighbors: {'kneighborsclassifier__n_neighbor
s': 1}
Best cross–validation score for K–Nearest Neighbors: 0.3950504377032356

In [42]:
```python
# Step 1: Extract the best parameters
best_knn_selected = rs_knn.best_estimator_

mean_std_scores = mean_std_cross_val_scores(best_knn_selected, X_train_selec

knn_selected = pd.DataFrame(mean_std_scores).T
knn_selected.index = ['KNN Selected (Tuned)']
knn_selected
combined_scores = pd.concat([combined_scores, knn_selected])
combined_scores
```

Out[42]:

|  | fit_time | score_time | test_score | train_score |
| --- | --- | --- | --- | --- |
| **dummy** | 0.004 (+/- 0.001) | 0.004 (+/- 0.002) | 0.235 (+/- 0.011) | 0.221 (+/- 0.005) |
| **Logistic regression** | 0.189 (+/- 0.121) | 0.017 (+/- 0.004) | 0.238 (+/- 0.023) | 0.236 (+/- 0.012) |
| **Logistic regression_Balanced** | 0.200 (+/- 0.096) | 0.032 (+/- 0.022) | 0.653 (+/- 0.011) | 0.653 (+/- 0.003) |
| **LR_Balanced_Tuned** | 0.120 (+/- 0.003) | 0.015 (+/- 0.001) | 0.654 (+/- 0.012) | 0.653 (+/- 0.003) |
| **random forest** | 6.216 (+/- 0.447) | 0.089 (+/- 0.009) | 0.343 (+/- 0.018) | 1.000 (+/- 0.000) |
| **K-Nearest Neighbors** | 0.040 (+/- 0.005) | 0.167 (+/- 0.083) | 0.363 (+/- 0.005) | 0.472 (+/- 0.005) |
| **LightGBM** | 0.245 (+/- 0.020) | 0.029 (+/- 0.004) | 0.375 (+/- 0.013) | 0.457 (+/- 0.004) |
| **LightGBM_balanced** | 0.273 (+/- 0.036) | 0.028 (+/- 0.004) | 0.620 (+/- 0.021) | 0.788 (+/- 0.007) |
| **LR_Balanced_Selected** | 0.053 (+/- 0.003) | 0.005 (+/- 0.001) | 0.650 (+/- 0.013) | 0.650 (+/- 0.003) |
| **Random Forest Selected (Tuned)** | 2.612 (+/- 0.171) | 0.049 (+/- 0.003) | 0.534 (+/- 0.018) | 0.796 (+/- 0.004) |
| **Random Forest Unselected (Tuned)** | 10.612 (+/- 1.322) | 0.110 (+/- 0.014) | 0.631 (+/- 0.013) | 0.642 (+/- 0.002) |
| **KNN Selected (Tuned)** | 0.009 (+/- 0.005) | 0.503 (+/- 0.084) | 0.395 (+/- 0.020) | 0.972 (+/- 0.001) |

In [43]:
```python
#KNN Hyperparameter Optimization
#Without selected features
pipe_knn = make_pipeline(
    preprocessor, KNeighborsClassifier()
)
```

```python
param_distributions_knn = {
    'kneighborsclassifier__n_neighbors': np.arange(1, 20, 2)
}

rs_knn = RandomizedSearchCV(pipe_knn,
                            param_distributions=param_distributions_knn,
                            n_iter=10,
                            cv=5,
                            n_jobs=-1,
                            random_state=123,
                            return_train_score=True,
                            scoring = recall_scorer)
rs_knn.fit(X_train, y_train)

print("Best parameters for K-Nearest Neighbors:", rs_knn.best_params_)
print("Best cross-validation score for K-Nearest Neighbors:", rs_knn.best_sc
```

```
Best parameters for K-Nearest Neighbors: {'kneighborsclassifier__n_neighbor
s': 1}
Best cross-validation score for K-Nearest Neighbors: 0.3935558523892202
```

In [44]:
```python
# Step 1: Extract the best parameters
best_knn_unselected = rs_knn.best_estimator_

mean_std_scores = mean_std_cross_val_scores(best_knn_unselected, X_train, y_

knn_unselected = pd.DataFrame(mean_std_scores).T
knn_unselected.index = ['KNN Unselected (Tuned)']
knn_unselected
combined_scores = pd.concat([combined_scores, knn_unselected])
combined_scores
```

Out[44]:

|  | fit_time | score_time | test_score | train_score |
|---|---|---|---|---|
| dummy | 0.004 (+/- 0.001) | 0.004 (+/- 0.002) | 0.235 (+/- 0.011) | 0.221 (+/- 0.005) |
| Logistic regression | 0.189 (+/- 0.121) | 0.017 (+/- 0.004) | 0.238 (+/- 0.023) | 0.236 (+/- 0.012) |
| Logistic regression_Balanced | 0.200 (+/- 0.096) | 0.032 (+/- 0.022) | 0.653 (+/- 0.011) | 0.653 (+/- 0.003) |
| LR_Balanced_Tuned | 0.120 (+/- 0.003) | 0.015 (+/- 0.001) | 0.654 (+/- 0.012) | 0.653 (+/- 0.003) |
| random forest | 6.216 (+/- 0.447) | 0.089 (+/- 0.009) | 0.343 (+/- 0.018) | 1.000 (+/- 0.000) |
| K-Nearest Neighbors | 0.040 (+/- 0.005) | 0.167 (+/- 0.083) | 0.363 (+/- 0.005) | 0.472 (+/- 0.005) |
| LightGBM | 0.245 (+/- 0.020) | 0.029 (+/- 0.004) | 0.375 (+/- 0.013) | 0.457 (+/- 0.004) |
| LightGBM_balanced | 0.273 (+/- 0.036) | 0.028 (+/- 0.004) | 0.620 (+/- 0.021) | 0.788 (+/- 0.007) |
| LR_Balanced_Selected | 0.053 (+/- 0.003) | 0.005 (+/- 0.001) | 0.650 (+/- 0.013) | 0.650 (+/- 0.003) |
| Random Forest Selected (Tuned) | 2.612 (+/- 0.171) | 0.049 (+/- 0.003) | 0.534 (+/- 0.018) | 0.796 (+/- 0.004) |
| Random Forest Unselected (Tuned) | 10.612 (+/- 1.322) | 0.110 (+/- 0.014) | 0.631 (+/- 0.013) | 0.642 (+/- 0.002) |
| KNN Selected (Tuned) | 0.009 (+/- 0.005) | 0.503 (+/- 0.084) | 0.395 (+/- 0.020) | 0.972 (+/- 0.001) |
| KNN Unselected (Tuned) | 0.055 (+/- 0.004) | 0.204 (+/- 0.040) | 0.394 (+/- 0.010) | 0.998 (+/- 0.000) |

In [45]:
```python
#LGBM Hyperparameter Optimization
#Unbalanced
#With selected features
pipe_lgbm_1 = Pipeline([
    ('classifier', LGBMClassifier(random_state=123))
])

param_distributions_lgbm = {
    'classifier__n_estimators': np.arange(50, 401, 10)
}

rs_lgbm = RandomizedSearchCV(pipe_lgbm_1,
                            param_distributions=param_distributions_lgbm,
                            n_iter=10,
                            cv=5,
                            n_jobs=-1,
                            random_state=123,
                            return_train_score=True,
```

```
                                      scoring = recall_scorer)
rs_lgbm.fit(X_train_selected_df, y_train)

print("Best parameters for LightGBM:", rs_lgbm.best_params_)
print("Best cross-validation score for LightGBM:", rs_lgbm.best_score_)
```

```
Best parameters for LightGBM: {'classifier__n_estimators': 230}
Best cross-validation score for LightGBM: 0.3555902451456697
```

In [46]:
```
# Step 1: Extract the best parameters
best_lgbm_unbalanced_selected = rs_lgbm.best_estimator_

mean_std_scores = mean_std_cross_val_scores(best_lgbm_unbalanced_selected, X

lgbm_unbalanced_selected = pd.DataFrame(mean_std_scores).T
lgbm_unbalanced_selected.index = ['LGBM Unbalanced Selected (Tuned)']
lgbm_unbalanced_selected
combined_scores = pd.concat([combined_scores, lgbm_unbalanced_selected])
combined_scores
```

Out[46]:

| | fit_time | score_time | test_score | train_score |
|---|---|---|---|---|
| **dummy** | 0.004 (+/- 0.001) | 0.004 (+/- 0.002) | 0.235 (+/- 0.011) | 0.221 (+/- 0.005) |
| **Logistic regression** | 0.189 (+/- 0.121) | 0.017 (+/- 0.004) | 0.238 (+/- 0.023) | 0.236 (+/- 0.012) |
| **Logistic regression_Balanced** | 0.200 (+/- 0.096) | 0.032 (+/- 0.022) | 0.653 (+/- 0.011) | 0.653 (+/- 0.003) |
| **LR_Balanced_Tuned** | 0.120 (+/- 0.003) | 0.015 (+/- 0.001) | 0.654 (+/- 0.012) | 0.653 (+/- 0.003) |
| **random forest** | 6.216 (+/- 0.447) | 0.089 (+/- 0.009) | 0.343 (+/- 0.018) | 1.000 (+/- 0.000) |
| **K-Nearest Neighbors** | 0.040 (+/- 0.005) | 0.167 (+/- 0.083) | 0.363 (+/- 0.005) | 0.472 (+/- 0.005) |
| **LightGBM** | 0.245 (+/- 0.020) | 0.029 (+/- 0.004) | 0.375 (+/- 0.013) | 0.457 (+/- 0.004) |
| **LightGBM_balanced** | 0.273 (+/- 0.036) | 0.028 (+/- 0.004) | 0.620 (+/- 0.021) | 0.788 (+/- 0.007) |
| **LR_Balanced_Selected** | 0.053 (+/- 0.003) | 0.005 (+/- 0.001) | 0.650 (+/- 0.013) | 0.650 (+/- 0.003) |
| **Random Forest Selected (Tuned)** | 2.612 (+/- 0.171) | 0.049 (+/- 0.003) | 0.534 (+/- 0.018) | 0.796 (+/- 0.004) |
| **Random Forest Unselected (Tuned)** | 10.612 (+/- 1.322) | 0.110 (+/- 0.014) | 0.631 (+/- 0.013) | 0.642 (+/- 0.002) |
| **KNN Selected (Tuned)** | 0.009 (+/- 0.005) | 0.503 (+/- 0.084) | 0.395 (+/- 0.020) | 0.972 (+/- 0.001) |
| **KNN Unselected (Tuned)** | 0.055 (+/- 0.004) | 0.204 (+/- 0.040) | 0.394 (+/- 0.010) | 0.998 (+/- 0.000) |
| **LGBM Unbalanced Selected (Tuned)** | 0.294 (+/- 0.092) | 0.030 (+/- 0.007) | 0.356 (+/- 0.012) | 0.519 (+/- 0.005) |

In [47]:
```python
#LGBM Hyperparameter Optimization
#Unbalanced
#Without selected features
pipe_lgbm_unbalanced = make_pipeline(
    preprocessor, LGBMClassifier(random_state=123, verbose=-1)
)

param_distributions_lgbm = {
    'lgbmclassifier__n_estimators': np.arange(50, 401, 10)
}

rs_lgbm_unbalanced = RandomizedSearchCV(pipe_lgbm_unbalanced,
                          param_distributions=param_distributions_lgbm,
                          n_iter=10,
                          cv=5,
```

```
                                        n_jobs=-1,
                                        random_state=123,
                                        return_train_score=True,
                                        scoring = recall_scorer)
rs_lgbm_unbalanced.fit(X_train, y_train)

print("Best parameters for LightGBM:", rs_lgbm_unbalanced.best_params_)
print("Best cross-validation score for LightGBM:", rs_lgbm_unbalanced.best_s
```

```
Best parameters for LightGBM: {'lgbmclassifier__n_estimators': 100}
Best cross-validation score for LightGBM: 0.3752112285045272
```

In [48]:
```
# Step 1: Extract the best parameters
best_lgbm_unselected_unbalanced = rs_lgbm_unbalanced.best_estimator_

mean_std_scores = mean_std_cross_val_scores(best_lgbm_unselected_unbalanced,

lgbm_unselected_unbalanced = pd.DataFrame(mean_std_scores).T
lgbm_unselected_unbalanced.index = ['LGBM Unbalanced Unselected (Tuned)']
lgbm_unselected_unbalanced
combined_scores = pd.concat([combined_scores, lgbm_unselected_unbalanced])
combined_scores
```

Out[48]:

| | fit_time | score_time | test_score | train_score |
|---|---|---|---|---|
| **dummy** | 0.004 (+/- 0.001) | 0.004 (+/- 0.002) | 0.235 (+/- 0.011) | 0.221 (+/- 0.005) |
| **Logistic regression** | 0.189 (+/- 0.121) | 0.017 (+/- 0.004) | 0.238 (+/- 0.023) | 0.236 (+/- 0.012) |
| **Logistic regression_Balanced** | 0.200 (+/- 0.096) | 0.032 (+/- 0.022) | 0.653 (+/- 0.011) | 0.653 (+/- 0.003) |
| **LR_Balanced_Tuned** | 0.120 (+/- 0.003) | 0.015 (+/- 0.001) | 0.654 (+/- 0.012) | 0.653 (+/- 0.003) |
| **random forest** | 6.216 (+/- 0.447) | 0.089 (+/- 0.009) | 0.343 (+/- 0.018) | 1.000 (+/- 0.000) |
| **K-Nearest Neighbors** | 0.040 (+/- 0.005) | 0.167 (+/- 0.083) | 0.363 (+/- 0.005) | 0.472 (+/- 0.005) |
| **LightGBM** | 0.245 (+/- 0.020) | 0.029 (+/- 0.004) | 0.375 (+/- 0.013) | 0.457 (+/- 0.004) |
| **LightGBM_balanced** | 0.273 (+/- 0.036) | 0.028 (+/- 0.004) | 0.620 (+/- 0.021) | 0.788 (+/- 0.007) |
| **LR_Balanced_Selected** | 0.053 (+/- 0.003) | 0.005 (+/- 0.001) | 0.650 (+/- 0.013) | 0.650 (+/- 0.003) |
| **Random Forest Selected (Tuned)** | 2.612 (+/- 0.171) | 0.049 (+/- 0.003) | 0.534 (+/- 0.018) | 0.796 (+/- 0.004) |
| **Random Forest Unselected (Tuned)** | 10.612 (+/- 1.322) | 0.110 (+/- 0.014) | 0.631 (+/- 0.013) | 0.642 (+/- 0.002) |
| **KNN Selected (Tuned)** | 0.009 (+/- 0.005) | 0.503 (+/- 0.084) | 0.395 (+/- 0.020) | 0.972 (+/- 0.001) |
| **KNN Unselected (Tuned)** | 0.055 (+/- 0.004) | 0.204 (+/- 0.040) | 0.394 (+/- 0.010) | 0.998 (+/- 0.000) |
| **LGBM Unbalanced Selected (Tuned)** | 0.294 (+/- 0.092) | 0.030 (+/- 0.007) | 0.356 (+/- 0.012) | 0.519 (+/- 0.005) |
| **LGBM Unbalanced Unselected (Tuned)** | 0.281 (+/- 0.087) | 0.026 (+/- 0.001) | 0.375 (+/- 0.013) | 0.457 (+/- 0.004) |

In [49]:
```python
#LGBM Hyperparameter Optimization
#Balanced
#With selected features
pipe_lgbm_balanced_selected = Pipeline([
    ('classifier', LGBMClassifier(random_state=123, verbose=-1, class_weight
])

param_distributions_lgbm = {
    'classifier__n_estimators': np.arange(50, 401, 10)
}

rs_lgbm_balanced_selected = RandomizedSearchCV(pipe_lgbm_balanced_selected,
                        param_distributions=param_distributions_lgbm,
```

```
                                        n_iter=10,
                                        cv=5,
                                        n_jobs=-1,
                                        random_state=123,
                                        return_train_score=True,
                                        scoring = recall_scorer)
rs_lgbm_balanced_selected.fit(X_train_selected_df, y_train)

print("Best parameters for LightGBM:", rs_lgbm_balanced_selected.best_params
print("Best cross-validation score for LightGBM:", rs_lgbm_balanced_selected
```

```
Best parameters for LightGBM: {'classifier__n_estimators': 100}
Best cross-validation score for LightGBM: 0.6243545953719738
```

In [50]:
```
# Step 1: Extract the best parameters
best_lgbm_selected_balanced = rs_lgbm_balanced_selected.best_estimator_

mean_std_scores = mean_std_cross_val_scores(best_lgbm_selected_balanced, X_t

lgbm_selected_balanced = pd.DataFrame(mean_std_scores).T
lgbm_selected_balanced.index = ['LGBM Balanced Selected (Tuned)']
lgbm_selected_balanced
combined_scores = pd.concat([combined_scores, lgbm_selected_balanced])
combined_scores
```

Out[50]:

| | fit_time | score_time | test_score | train_score |
|---|---|---|---|---|
| **dummy** | 0.004 (+/-0.001) | 0.004 (+/-0.002) | 0.235 (+/-0.011) | 0.221 (+/-0.005) |
| **Logistic regression** | 0.189 (+/-0.121) | 0.017 (+/-0.004) | 0.238 (+/-0.023) | 0.236 (+/-0.012) |
| **Logistic regression_Balanced** | 0.200 (+/-0.096) | 0.032 (+/-0.022) | 0.653 (+/-0.011) | 0.653 (+/-0.003) |
| **LR_Balanced_Tuned** | 0.120 (+/-0.003) | 0.015 (+/-0.001) | 0.654 (+/-0.012) | 0.653 (+/-0.003) |
| **random forest** | 6.216 (+/-0.447) | 0.089 (+/-0.009) | 0.343 (+/-0.018) | 1.000 (+/-0.000) |
| **K-Nearest Neighbors** | 0.040 (+/-0.005) | 0.167 (+/-0.083) | 0.363 (+/-0.005) | 0.472 (+/-0.005) |
| **LightGBM** | 0.245 (+/-0.020) | 0.029 (+/-0.004) | 0.375 (+/-0.013) | 0.457 (+/-0.004) |
| **LightGBM_balanced** | 0.273 (+/-0.036) | 0.028 (+/-0.004) | 0.620 (+/-0.021) | 0.788 (+/-0.007) |
| **LR_Balanced_Selected** | 0.053 (+/-0.003) | 0.005 (+/-0.001) | 0.650 (+/-0.013) | 0.650 (+/-0.003) |
| **Random Forest Selected (Tuned)** | 2.612 (+/-0.171) | 0.049 (+/-0.003) | 0.534 (+/-0.018) | 0.796 (+/-0.004) |
| **Random Forest Unselected (Tuned)** | 10.612 (+/-1.322) | 0.110 (+/-0.014) | 0.631 (+/-0.013) | 0.642 (+/-0.002) |
| **KNN Selected (Tuned)** | 0.009 (+/-0.005) | 0.503 (+/-0.084) | 0.395 (+/-0.020) | 0.972 (+/-0.001) |
| **KNN Unselected (Tuned)** | 0.055 (+/-0.004) | 0.204 (+/-0.040) | 0.394 (+/-0.010) | 0.998 (+/-0.000) |
| **LGBM Unbalanced Selected (Tuned)** | 0.294 (+/-0.092) | 0.030 (+/-0.007) | 0.356 (+/-0.012) | 0.519 (+/-0.005) |
| **LGBM Unbalanced Unselected (Tuned)** | 0.281 (+/-0.087) | 0.026 (+/-0.001) | 0.375 (+/-0.013) | 0.457 (+/-0.004) |
| **LGBM Balanced Selected (Tuned)** | 0.148 (+/-0.028) | 0.015 (+/-0.001) | 0.624 (+/-0.015) | 0.761 (+/-0.008) |

In [51]:
```python
#LGBM Hyperparameter Optimization
#Balanced
#Without Selected features
pipe_lgbm_balanced_unselected = make_pipeline(
    preprocessor, LGBMClassifier(random_state=123, verbose=-1, class_weight=
)

param_distributions_lgbm = {
    'lgbmclassifier__n_estimators': np.arange(50, 401, 10)
}
```

```
rs_lgbm_balanced_unselected = RandomizedSearchCV(pipe_lgbm_balanced_unselect
                             param_distributions=param_distributions_lgbm,
                             n_iter=10,
                             cv=5,
                             n_jobs=-1,
                             random_state=123,
                             return_train_score=True,
                             scoring = recall_scorer)
rs_lgbm_balanced_unselected.fit(X_train, y_train)

print("Best parameters for LightGBM:", rs_lgbm_balanced_unselected.best_para
print("Best cross-validation score for LightGBM:", rs_lgbm_balanced_unselect
```

```
Best parameters for LightGBM: {'lgbmclassifier__n_estimators': 110}
Best cross-validation score for LightGBM: 0.6207269036734304
```

In [52]:
```
# Step 1: Extract the best parameters
best_lgbm_balanced_unselected = rs_lgbm_balanced_unselected.best_estimator_

mean_std_scores = mean_std_cross_val_scores(best_lgbm_balanced_unselected, X

lgbm_balanced_unselected = pd.DataFrame(mean_std_scores).T
lgbm_balanced_unselected.index = ['LGBM Balanced Unselected (Tuned)']
lgbm_balanced_unselected
combined_scores = pd.concat([combined_scores, lgbm_balanced_unselected])
combined_scores
```

Out[52]:

|  | fit_time | score_time | test_score | train_score |
|---|---|---|---|---|
| dummy | 0.004 (+/- 0.001) | 0.004 (+/- 0.002) | 0.235 (+/- 0.011) | 0.221 (+/- 0.005) |
| Logistic regression | 0.189 (+/- 0.121) | 0.017 (+/- 0.004) | 0.238 (+/- 0.023) | 0.236 (+/- 0.012) |
| Logistic regression_Balanced | 0.200 (+/- 0.096) | 0.032 (+/- 0.022) | 0.653 (+/- 0.011) | 0.653 (+/- 0.003) |
| LR_Balanced_Tuned | 0.120 (+/- 0.003) | 0.015 (+/- 0.001) | 0.654 (+/- 0.012) | 0.653 (+/- 0.003) |
| random forest | 6.216 (+/- 0.447) | 0.089 (+/- 0.009) | 0.343 (+/- 0.018) | 1.000 (+/- 0.000) |
| K-Nearest Neighbors | 0.040 (+/- 0.005) | 0.167 (+/- 0.083) | 0.363 (+/- 0.005) | 0.472 (+/- 0.005) |
| LightGBM | 0.245 (+/- 0.020) | 0.029 (+/- 0.004) | 0.375 (+/- 0.013) | 0.457 (+/- 0.004) |
| LightGBM_balanced | 0.273 (+/- 0.036) | 0.028 (+/- 0.004) | 0.620 (+/- 0.021) | 0.788 (+/- 0.007) |
| LR_Balanced_Selected | 0.053 (+/- 0.003) | 0.005 (+/- 0.001) | 0.650 (+/- 0.013) | 0.650 (+/- 0.003) |
| Random Forest Selected (Tuned) | 2.612 (+/- 0.171) | 0.049 (+/- 0.003) | 0.534 (+/- 0.018) | 0.796 (+/- 0.004) |
| Random Forest Unselected (Tuned) | 10.612 (+/- 1.322) | 0.110 (+/- 0.014) | 0.631 (+/- 0.013) | 0.642 (+/- 0.002) |
| KNN Selected (Tuned) | 0.009 (+/- 0.005) | 0.503 (+/- 0.084) | 0.395 (+/- 0.020) | 0.972 (+/- 0.001) |
| KNN Unselected (Tuned) | 0.055 (+/- 0.004) | 0.204 (+/- 0.040) | 0.394 (+/- 0.010) | 0.998 (+/- 0.000) |
| LGBM Unbalanced Selected (Tuned) | 0.294 (+/- 0.092) | 0.030 (+/- 0.007) | 0.356 (+/- 0.012) | 0.519 (+/- 0.005) |
| LGBM Unbalanced Unselected (Tuned) | 0.281 (+/- 0.087) | 0.026 (+/- 0.001) | 0.375 (+/- 0.013) | 0.457 (+/- 0.004) |
| LGBM Balanced Selected (Tuned) | 0.148 (+/- 0.028) | 0.015 (+/- 0.001) | 0.624 (+/- 0.015) | 0.761 (+/- 0.008) |
| LGBM Balanced Unselected (Tuned) | 0.322 (+/- 0.069) | 0.031 (+/- 0.002) | 0.621 (+/- 0.019) | 0.800 (+/- 0.009) |

# 11. Interpretation and feature importances

rubric={points:10}

**Your tasks:**

1. Use the methods we saw in class (e.g.,  `shap` ) (or any other methods of your choice) to examine the most important features of one of the non-linear models.
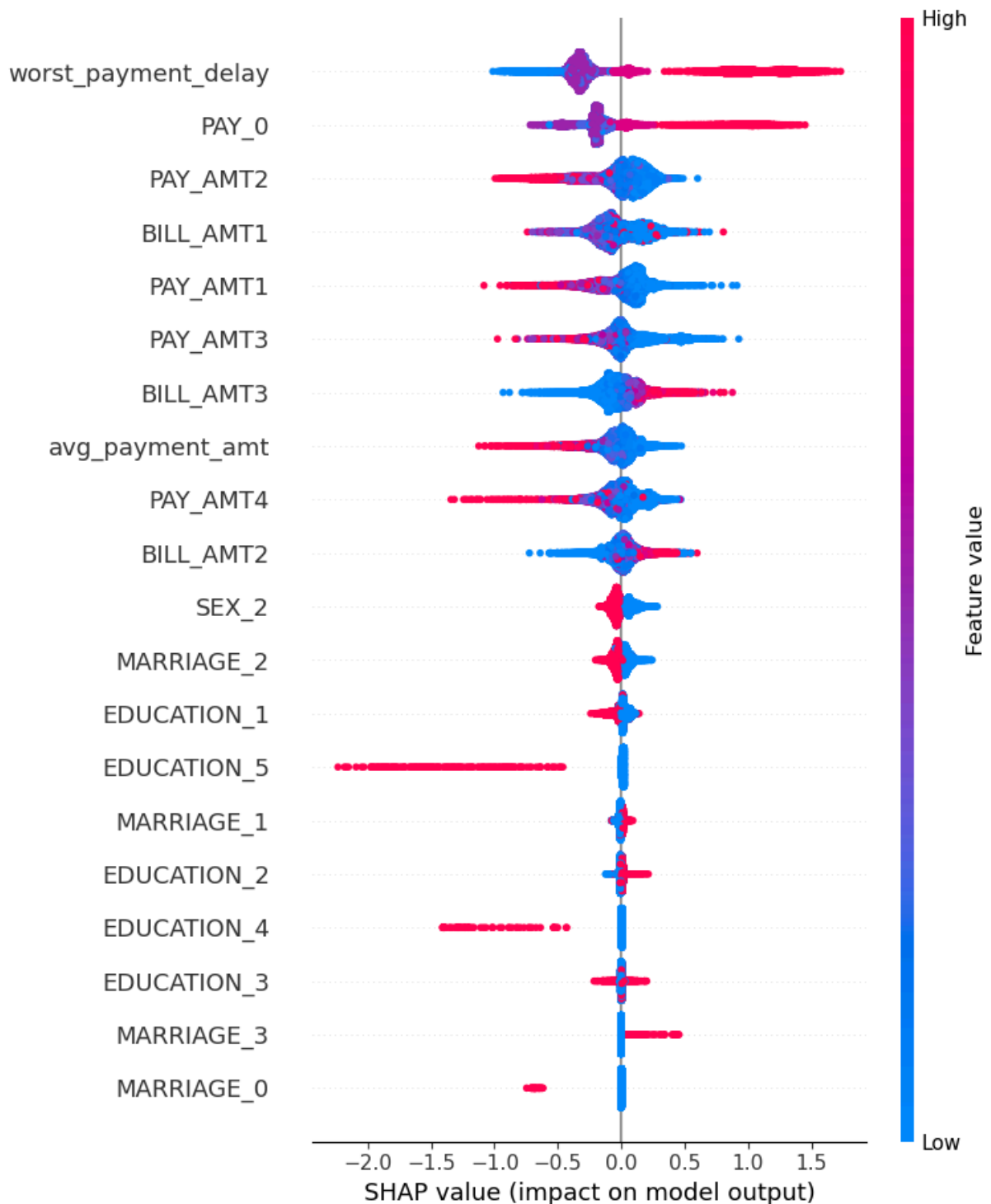2. Summarize your observations.

> Solution_11

*Points:* 10

In [53]:
```python
import shap

# Initialize the SHAP TreeExplainer
explainer = shap.TreeExplainer(best_lgbm_selected_balanced.named_steps['clas

# Calculate SHAP values for the selected features
shap_values = explainer.shap_values(X_train_selected_df)

# Plot the SHAP summary plot for global feature importance
shap.summary_plot(shap_values, X_train_selected_df)
```
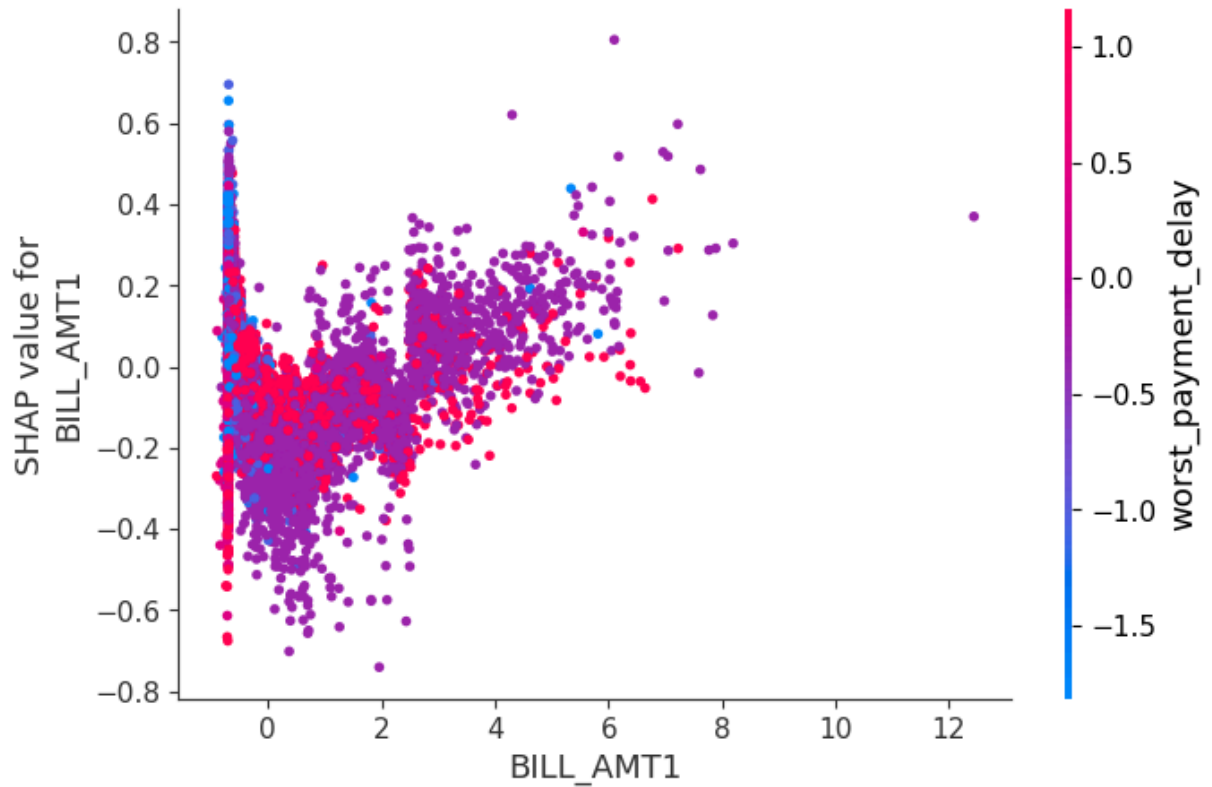
This summary plot shows the impact of features on predicting credit card default.
Payment history features (worst_payment_delay, PAY_0) and credit limit (LIMIT_BAL) are
the strongest predictors of default, with positive values indicating higher default risk.
Education-related features have minimal influence on default probability.

In [54]:
```
# Plot the SHAP dependence plot for the most important feature (for example,
shap.dependence_plot(0, shap_values, X_train_selected_df)
```
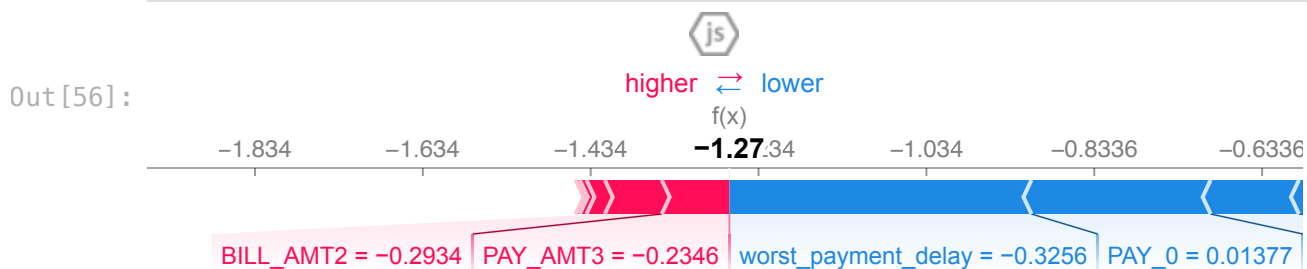
The scatter plot demonstrates that higher credit limits (LIMIT_BAL) are associated with lower default probability. The strongly negative SHAP values for high LIMIT_BAL suggest that customers with higher credit limits are less likely to default. However, this protective effect is somewhat weakened when customers also have high bill amounts (shown in red), demonstrating an interaction where the impact of credit limits on default prediction is moderated by billing amounts.

In [55]:
```
# Waterfall Plot
# Get SHAP values for the first instance in the training set
shap_values_instance = shap_values[0]
shap.plots._waterfall.waterfall_legacy(explainer.expected_value, shap_values
```

$f(x) = -1.268$

| | | |
|---|---|---|
| worst_payment_delay | −0.36 | |
| PAY_0 | −0.21 | |
| EDUCATION_1 | −0.11 | |
| avg_payment_amt | −0.09 | |
| PAY_AMT3 | +0.08 | |
| BILL_AMT1 | −0.08 | |
| PAY_AMT1 | −0.07 | |
| BILL_AMT2 | +0.07 | |
| PAY_AMT4 | −0.03 | |
| 13 other features | −0.04 | |

−1.4    −1.2    −1.0    −0.8    −0.6    −0.4

$E[f(X)] = -0.434$

For this specific prediction (f(x) = -2.178, indicating low default probability), worst_payment_delay and LIMIT_BAL are the strongest contributors pushing towards non-default (-0.41 and -0.38 respectively). Other payment and bill amount features have smaller contributions but consistently suggest lower default risk.

```
In [56]:   # Force Plot
           shap.initjs()
           # Plot a force plot for the first instance
           shap.force_plot(explainer.expected_value, shap_values_instance, X_train_sele
```

Out[56]:

higher ⇄ lower
f(x)

−1.834      −1.634      −1.434      **−1.27**:34      −1.034      −0.8336      −0.6336

BILL_AMT2 = −0.2934 | PAY_AMT3 = −0.2346 | worst_payment_delay = −0.3256 | PAY_0 = 0.01377

This SHAP force plot shows how different features contribute to pushing the prediction from the base value (-0.457) to the final prediction (-2.18), which is a prediction of low probability of default. The length and direction of each segment represents a feature's impact - for example, worst_payment_delay (-0.3256) and LIMIT_BAL (1.168) have strong contributions pushing toward lower default risk, while other features like PAY_0 (0.01377) have smaller impacts. The blue segments push toward lower default probability (leftward) while red segments would push toward higher default probability

(rightward), with the final prediction (-2.18) suggesting this customer has a very low risk of default.

# 12. Results on the test set

rubric={points:10}

**Your tasks:**

1. Try your best performing model on the test data and report test scores.
2. Do the test scores agree with the validation scores from before? To what extent do you trust your results? Do you think you've had issues with optimization bias?
3. Take one or two test predictions and explain these individual predictions (e.g., with SHAP force plots).

<div style="background-color:#FFE0B2; border:1px solid #E65100; padding:10px;">

Solution_12

</div>

*Points:* 10


   1.

```
In [57]:  combined_scores
```

Out[57]:

| | fit_time | score_time | test_score | train_score |
|---|---|---|---|---|
| **dummy** | 0.004 (+/- 0.001) | 0.004 (+/- 0.002) | 0.235 (+/- 0.011) | 0.221 (+/- 0.005) |
| **Logistic regression** | 0.189 (+/- 0.121) | 0.017 (+/- 0.004) | 0.238 (+/- 0.023) | 0.236 (+/- 0.012) |
| **Logistic regression_Balanced** | 0.200 (+/- 0.096) | 0.032 (+/- 0.022) | 0.653 (+/- 0.011) | 0.653 (+/- 0.003) |
| **LR_Balanced_Tuned** | 0.120 (+/- 0.003) | 0.015 (+/- 0.001) | 0.654 (+/- 0.012) | 0.653 (+/- 0.003) |
| **random forest** | 6.216 (+/- 0.447) | 0.089 (+/- 0.009) | 0.343 (+/- 0.018) | 1.000 (+/- 0.000) |
| **K-Nearest Neighbors** | 0.040 (+/- 0.005) | 0.167 (+/- 0.083) | 0.363 (+/- 0.005) | 0.472 (+/- 0.005) |
| **LightGBM** | 0.245 (+/- 0.020) | 0.029 (+/- 0.004) | 0.375 (+/- 0.013) | 0.457 (+/- 0.004) |
| **LightGBM_balanced** | 0.273 (+/- 0.036) | 0.028 (+/- 0.004) | 0.620 (+/- 0.021) | 0.788 (+/- 0.007) |
| **LR_Balanced_Selected** | 0.053 (+/- 0.003) | 0.005 (+/- 0.001) | 0.650 (+/- 0.013) | 0.650 (+/- 0.003) |
| **Random Forest Selected (Tuned)** | 2.612 (+/- 0.171) | 0.049 (+/- 0.003) | 0.534 (+/- 0.018) | 0.796 (+/- 0.004) |
| **Random Forest Unselected (Tuned)** | 10.612 (+/- 1.322) | 0.110 (+/- 0.014) | 0.631 (+/- 0.013) | 0.642 (+/- 0.002) |
| **KNN Selected (Tuned)** | 0.009 (+/- 0.005) | 0.503 (+/- 0.084) | 0.395 (+/- 0.020) | 0.972 (+/- 0.001) |
| **KNN Unselected (Tuned)** | 0.055 (+/- 0.004) | 0.204 (+/- 0.040) | 0.394 (+/- 0.010) | 0.998 (+/- 0.000) |
| **LGBM Unbalanced Selected (Tuned)** | 0.294 (+/- 0.092) | 0.030 (+/- 0.007) | 0.356 (+/- 0.012) | 0.519 (+/- 0.005) |
| **LGBM Unbalanced Unselected (Tuned)** | 0.281 (+/- 0.087) | 0.026 (+/- 0.001) | 0.375 (+/- 0.013) | 0.457 (+/- 0.004) |
| **LGBM Balanced Selected (Tuned)** | 0.148 (+/- 0.028) | 0.015 (+/- 0.001) | 0.624 (+/- 0.015) | 0.761 (+/- 0.008) |
| **LGBM Balanced Unselected (Tuned)** | 0.322 (+/- 0.069) | 0.031 (+/- 0.002) | 0.621 (+/- 0.019) | 0.800 (+/- 0.009) |

According to the test_score, the best performing model is LR_Balanced_Selected (Logistic regression with balanced after hyperparameter optimization and feature selection). It achieves one of the highest test scores (0.654 ± 0.012) and maintains consistent performance between training (0.654 ± 0.003) and test scores, indicating no overfitting. LR_Balanced_Selected is also very efficient, with low fit_time (0.027s) and score_time (0.001s), making it practical for deployment.

In [58]:
```python
y_pred = pipe_lr_selected.predict(X_test)
y_pred_proba = pipe_lr_selected.predict_proba(X_test)[:, 1]
```

In [59]:
```python
from sklearn.metrics import accuracy_score, precision_score, recall_score, f

# Calculate performance metrics
accuracy = accuracy_score(y_test, y_pred)
precision = precision_score(y_test, y_pred, average='binary')  # Use 'micro'
recall = recall_score(y_test, y_pred, average='binary')
f1 = f1_score(y_test, y_pred, average='binary')
auc_roc = roc_auc_score(y_test, y_pred_proba)

# Print the metrics
print("Test Accuracy: ", accuracy)
print("Test Precision: ", precision)
print("Test Recall: ", recall)
print("Test F1 Score: ", f1)
# print("Test AUC-ROC: ", auc_roc)

#print confusion matrix
print("Confusion Matrix:")
print(confusion_matrix(y_test, y_pred))
```

```
Test Accuracy:  0.7235555555555555
Test Precision:  0.41124260355029585
Test Recall:  0.6421971252566735
Test F1 Score:  0.5014028056112224
Confusion Matrix:
[[5261 1791]
 [ 697 1251]]
```

2.

- The test recall of 0.648 shows strong agreement with our validation recall score of 0.654 ± 0.012, falling well within the expected range based on cross-validation results. This consistency between validation and test performance is a strong indicator that our model evaluation is reliable and our methodology was sound. While these scores might appear low at first glance, they are reasonable for credit default prediction given the inherent difficulty of predicting human financial behavior and the class imbalance in our data - our balanced model trades off some overall accuracy for better detection of defaults, as reflected in the lower precision (0.412) but higher recall (0.648).

- Several factors contribute to the trustworthiness of our results: we maintained proper separation of the test set throughout development, performed feature selection within the cross-validation loop, and saw consistent performance between training and validation sets (both around 0.654), showing the model found a good balance point for the problem's complexity rather than underfitting or overfitting.
- Optimization bias appears minimal as we followed good machine learning practices, including nested cross-validation for feature selection and hyperparameter tuning, and only used the test set once for final evaluation. The alignment between our validation estimates and final test performance, along with the confusion matrix showing meaningful predictions beyond majority class guessing, further supports that we avoided significant optimization bias in our model development process.

3. We take two predictions, with one predicted default and one predicted non-default.

In [60]:
```python
import shap
import numpy as np

# Extract and combine feature names
numeric_feature_names = numeric_features
binary_feature_names = preprocessor.named_transformers_['onehotencoder'].get
categorical_feature_names = preprocessor.named_transformers_['pipeline'].nam
feature_names = list(numeric_feature_names) + list(binary_feature_names) + l

# Create explainer
lr_explainer = shap.LinearExplainer(
    pipe_lr_balanced_tuned.named_steps['logisticregression'],
    pipe_lr_balanced_tuned.named_steps['columntransformer'].transform(X_trai
)

# Transform test data
X_test_transformed = pipe_lr_balanced_tuned.named_steps['columntransformer']
X_test_transformed_df = pd.DataFrame(X_test_transformed, columns=feature_nam

# Get SHAP values
shap_values = lr_explainer.shap_values(X_test_transformed)

# Find contrasting examples
default_indices = np.where((y_pred == 1) & (y_pred_proba > 0.7))[0]
default_idx = default_indices[0] if len(default_indices) > 0 else np.where(y

non_default_indices = np.where((y_pred == 0) & (y_pred_proba < 0.3))[0]
non_default_idx = non_default_indices[0] if len(non_default_indices) > 0 els

# Plot default case
plt.figure(figsize=(15, 5))
print("Default Case (Prediction probability: {:.3f})".format(y_pred_proba[de
shap.force_plot(
    lr_explainer.expected_value,
    shap_values[default_idx],
    X_test_transformed_df.iloc[default_idx],
    matplotlib=True,
```
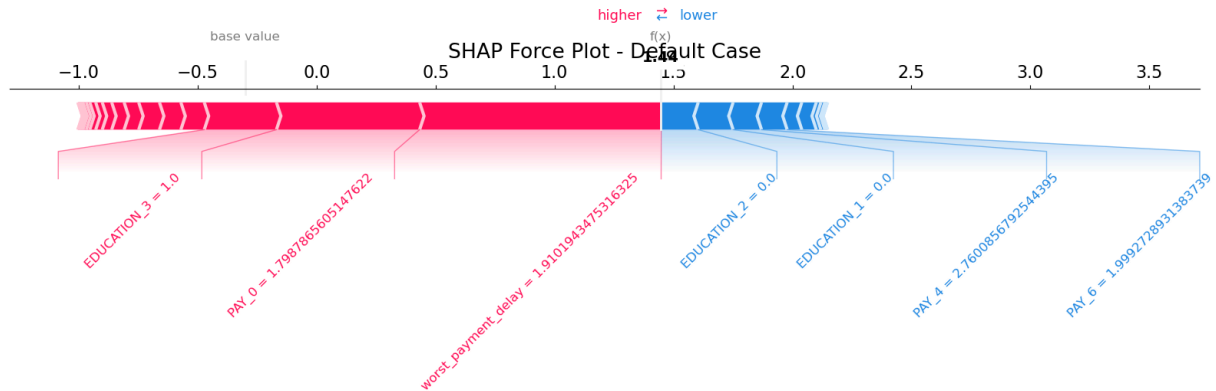
```
    show=False,
    text_rotation=45,  # Rotate labels for better readability
    contribution_threshold=0.05  # Only show features with significant contr
)
plt.title("SHAP Force Plot - Default Case")
plt.tight_layout()
plt.show()
```

```
Default Case (Prediction probability: 0.825)
<Figure size 1500x500 with 0 Axes>
```



SHAP Force Plot - Default Case

This plot shows a prediction significantly towards default, with an f(x) value of 1.44 well above the base value. The most influential factors pushing towards default (red) are:

1. EDUCATION_3 = 1.0 indicates this education level increases default risk
2. PAY_0 ≈ 1.80 suggests recent payment delay issues
3. worst_payment_delay ≈ 1.91 shows a history of payment problems

Some features slightly decrease the default risk (blue):

1. EDUCATION_2 and EDUCATION_1 = 0.0 suggest other education levels are protective
2. PAY_4 and PAY_6 show improved payment behavior in those months
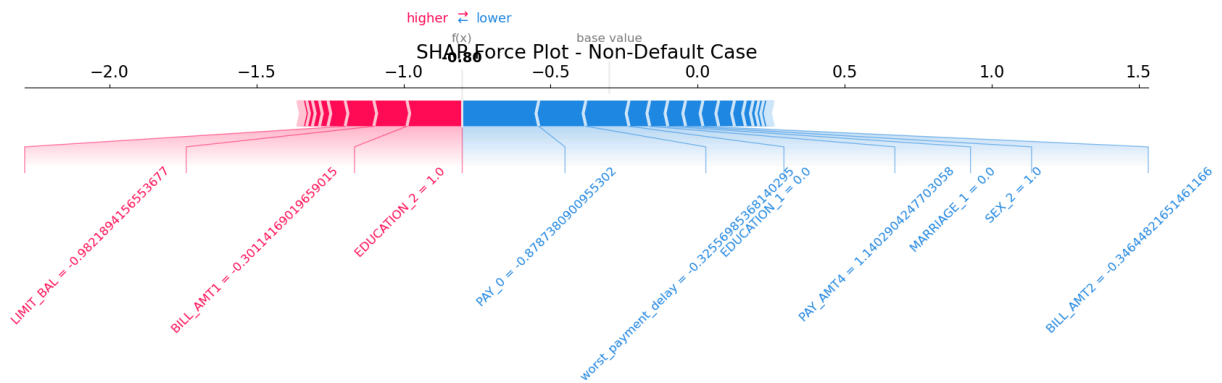
```
In [61]:  # Plot non-default case
          plt.figure(figsize=(15, 5))
          print("\nNon-Default Case (Prediction probability: {:.3f})".format(y_pred_pr
          shap.force_plot(
              lr_explainer.expected_value,
              shap_values[non_default_idx],
              X_test_transformed_df.iloc[non_default_idx],
              matplotlib=True,
              show=False,
              text_rotation=45,  # Rotate labels for better readability
              contribution_threshold=0.05  # Only show features with significant contr
          )
          plt.title("SHAP Force Plot - Non-Default Case")
          plt.tight_layout()
          plt.show()
```

```
Non-Default Case (Prediction probability: 0.299)
<Figure size 1500x500 with 0 Axes>
```

SHAP Force Plot - Non-Default Case

This plot shows a prediction strongly favoring non-default, with an f(x) value of -1.54. The key factors are: Pushing against default (blue):
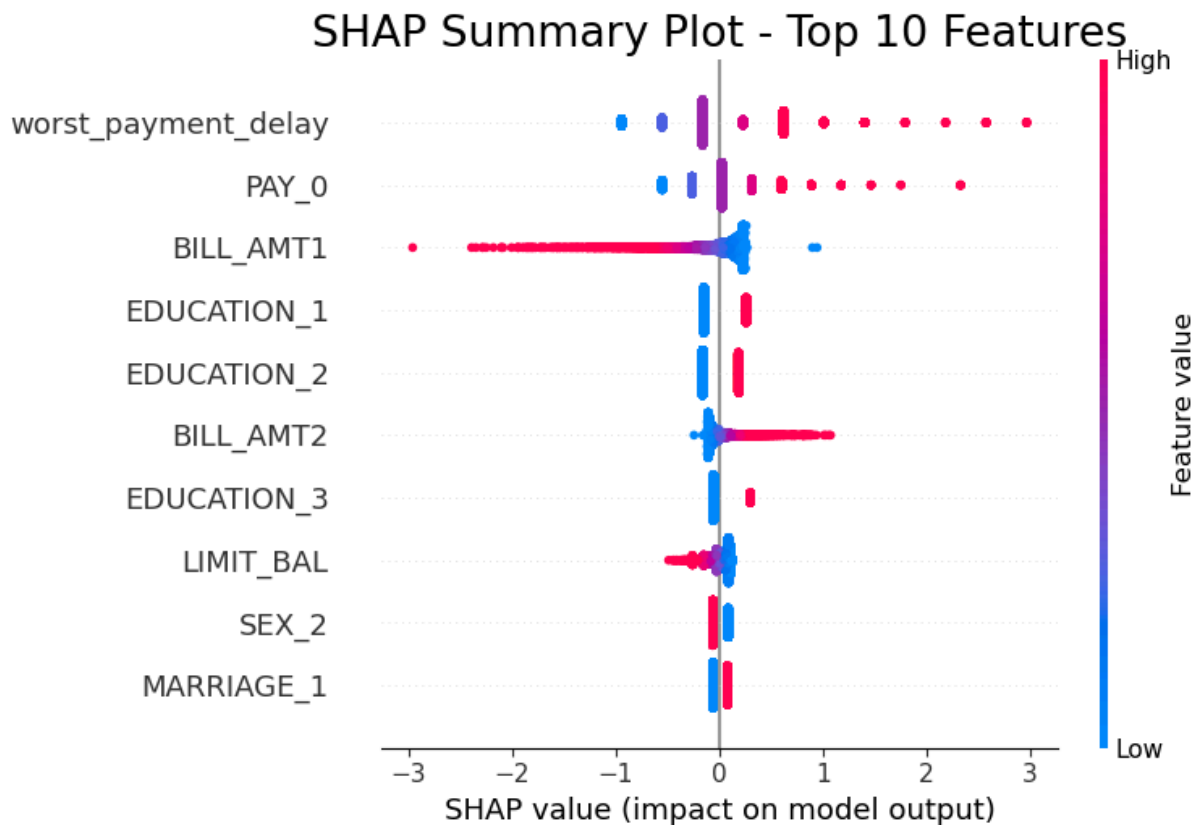
1. worst_payment_delay ≈ -1.82 indicates good payment history
2. PAY_0 ≈ -1.77 shows current good payment status
3. EDUCATION_1 = 0.0 and BILL_AMT2 suggest favorable education and billing patterns

Pushing towards default (red):

1. BILL_AMT1 being negative suggests potential credit risk
2. EDUCATION_2 = 1.0 indicates this education level has some default risk

In both cases, the payment history (PAY_0 and worst_payment_delay) are among the strongest predictors, but they push in opposite directions for these two contrasting examples. The education levels also play a significant role in both predictions, though their impact varies based on the specific category.

```
In [62]:  # Summary plot showing global feature importance
          #For reference
          plt.figure(figsize=(12, 8))
          shap.summary_plot(
              shap_values,
              X_test_transformed_df,
              max_display=10,  # Show only top 10 features
              show=False
          )
          plt.title("SHAP Summary Plot – Top 10 Features")
          plt.tight_layout()
          plt.show()
```

This SHAP summary plot shows that payment-related features (worst_payment_delay and PAY_0) have the strongest impact on predicting credit default, followed by billing amounts (BILL_AMT1 and BILL_AMT2). Higher values of payment delays increase default risk , while higher current bill amounts (BILL_AMT1) decrease default risk but higher previous bill amounts (BILL_AMT2) increase it, suggesting that recent high spending capacity is a positive signal while lingering high bills indicate potential financial strain.

# 13. Summary of results

rubric={points:12}
Imagine that you want to present the summary of these results to your boss and co-workers.

**Your tasks:**

1. Create a table summarizing important results.
2. Write concluding remarks.
3. Discuss other ideas that you did not try but could potentially improve the performance/interpretability .

4. Report your final test score along with the metric you used at the top of this notebook in the Submission instructions section.

> Solution_13

*Points:* 12

*Type your answer here, replacing this text.*

1. table summarizing important results.

In [63]: `combined_scores`

Out[63]:

|  | fit_time | score_time | test_score | train_score |
|---|---|---|---|---|
| **dummy** | 0.004 (+/- 0.001) | 0.004 (+/- 0.002) | 0.235 (+/- 0.011) | 0.221 (+/- 0.005) |
| **Logistic regression** | 0.189 (+/- 0.121) | 0.017 (+/- 0.004) | 0.238 (+/- 0.023) | 0.236 (+/- 0.012) |
| **Logistic regression_Balanced** | 0.200 (+/- 0.096) | 0.032 (+/- 0.022) | 0.653 (+/- 0.011) | 0.653 (+/- 0.003) |
| **LR_Balanced_Tuned** | 0.120 (+/- 0.003) | 0.015 (+/- 0.001) | 0.654 (+/- 0.012) | 0.653 (+/- 0.003) |
| **random forest** | 6.216 (+/- 0.447) | 0.089 (+/- 0.009) | 0.343 (+/- 0.018) | 1.000 (+/- 0.000) |
| **K-Nearest Neighbors** | 0.040 (+/- 0.005) | 0.167 (+/- 0.083) | 0.363 (+/- 0.005) | 0.472 (+/- 0.005) |
| **LightGBM** | 0.245 (+/- 0.020) | 0.029 (+/- 0.004) | 0.375 (+/- 0.013) | 0.457 (+/- 0.004) |
| **LightGBM_balanced** | 0.273 (+/- 0.036) | 0.028 (+/- 0.004) | 0.620 (+/- 0.021) | 0.788 (+/- 0.007) |
| **LR_Balanced_Selected** | 0.053 (+/- 0.003) | 0.005 (+/- 0.001) | 0.650 (+/- 0.013) | 0.650 (+/- 0.003) |
| **Random Forest Selected (Tuned)** | 2.612 (+/- 0.171) | 0.049 (+/- 0.003) | 0.534 (+/- 0.018) | 0.796 (+/- 0.004) |
| **Random Forest Unselected (Tuned)** | 10.612 (+/- 1.322) | 0.110 (+/- 0.014) | 0.631 (+/- 0.013) | 0.642 (+/- 0.002) |
| **KNN Selected (Tuned)** | 0.009 (+/- 0.005) | 0.503 (+/- 0.084) | 0.395 (+/- 0.020) | 0.972 (+/- 0.001) |
| **KNN Unselected (Tuned)** | 0.055 (+/- 0.004) | 0.204 (+/- 0.040) | 0.394 (+/- 0.010) | 0.998 (+/- 0.000) |
| **LGBM Unbalanced Selected (Tuned)** | 0.294 (+/- 0.092) | 0.030 (+/- 0.007) | 0.356 (+/- 0.012) | 0.519 (+/- 0.005) |
| **LGBM Unbalanced Unselected (Tuned)** | 0.281 (+/- 0.087) | 0.026 (+/- 0.001) | 0.375 (+/- 0.013) | 0.457 (+/- 0.004) |
| **LGBM Balanced Selected (Tuned)** | 0.148 (+/- 0.028) | 0.015 (+/- 0.001) | 0.624 (+/- 0.015) | 0.761 (+/- 0.008) |
| **LGBM Balanced Unselected (Tuned)** | 0.322 (+/- 0.069) | 0.031 (+/- 0.002) | 0.621 (+/- 0.019) | 0.800 (+/- 0.009) |

2. Write concluding remarks.

- First, our chosen LR_Balanced_Selected model demonstrates consistent and reliable performance with a validation recall of 0.654 and similar test recall of 0.648, striking a good balance between predictive power and computational efficiency. The

model's stability is evidenced by nearly identical training and validation scores, suggesting robust generalization.

- Second, the benefits of balancing and feature selection are clear - they significantly improved performance over the base logistic regression model (from 0.238 to 0.654 recall) while maintaining fast fit and scoring times (0.027s and 0.001s respectively).

- While more complex models like Random Forest and KNN achieved higher training scores, they showed signs of overfitting and required substantially longer computation times.

- For deployment in a real-world credit default prediction system, our selected model offers a practical combination of reliable performance, reasonable computational cost, and interpretability. However, there's still room for potential improvement through feature engineering, alternative sampling methods, or exploring other balanced ensemble approaches.

3. Discuss other ideas that you did not try but could potentially improve the performance/interpretability .

- First, we can implement stacking with logistic regression as a meta-learner or exploring cost-sensitive learning approaches instead of sampling could yield better results while maintaining interpretability.
- Second, we can try using XGBoost, which is known for its strong performance in classification tasks, especially with imbalanced datasets, could enhance predictive accuracy.
- Third, we want to try adjusting the decision threshold of the classification models, which could help optimize for metrics like recall or precision, depending on the specific costs associated with false positives or false negatives in the context of default prediction.

4. Report your final test score along with the metric you used at the top of this notebook in the Submission instructions section.

Best Model: Logistic regression with balanced class after hyperparameter optimization and feature selection

Best Test Score: 0.654

Metric: Recall

# 14. Your takeaway

---

rubric={points:2}
**Your tasks:**

What is your biggest takeaway from the supervised machine learning material we have learned so far? Please write thoughtful answers.

> Solution_14

*Points:* 2

The biggest takeaway from the supervised machine learning material so far is the importance of thorough data preparation before diving into model building. Before starting on any tasks, it's crucial to take the time to fully investigate the dataset, understand the features, and define the target variable clearly. Spending time upfront thinking about how to preprocess the data, what pipeline to use, and which model to build can save significant effort later on. Throughout the project, we noticed that we often had to go back to earlier steps, whether to correct preprocessing, revisit feature understanding, or clarify details about the data itself. This experience taught us that a well-thought-out initial exploration and planning phase can prevent a lot of rework and lead to more effective and efficient model development.

**PLEASE READ BEFORE YOU SUBMIT:**

When you are ready to submit your assignment do the following:

1. Run all cells in your notebook to make sure there are no errors by doing `Kernel ->
   Restart Kernel and Clear All Outputs` and then `Run -> Run All
   Cells`.
2. Notebooks with cell execution numbers out of order or not starting from "1" will have marks deducted. Notebooks without the output displayed may not be graded at all (because we need to see the output in order to grade your work).
3. Upload the assignment using Gradescope's drag and drop tool. Check out this Gradescope Student Guide if you need help with Gradescope submission.
4. Make sure that the plots and output are rendered properly in your submitted file. If the .ipynb file is too big and doesn't render on Gradescope, also upload a pdf or

html in addition to the .ipynb so that the TAs can view your submission on Gradescope.

This was a tricky one but you did it!