

```
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

0. [Submission instructions](#)
1. [Understanding the problem](#)
2. [Data splitting](#)
3. [EDA](#)
4. [Feature engineering](#)
5. [Preprocessing and transformations](#)
6. [Baseline model](#)
7. [Linear models](#)
8. [Different models](#)
9. [Feature selection](#)
10. [Hyperparameter optimization](#)
11. [Interpretation and feature importances](#)
12. [Results on the test set](#)
13. [Summary of the results](#)
14. [Your takeaway from the course](#)

## 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
```

```
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, OrdinalEncod
from sklearn.svm import SVC
from sklearn.tree import DecisionTreeClassifier
from sklearn.model_selection import RandomizedSearchCV
from sklearn.metrics import ConfusionMatrixDisplay
```

## 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 we 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

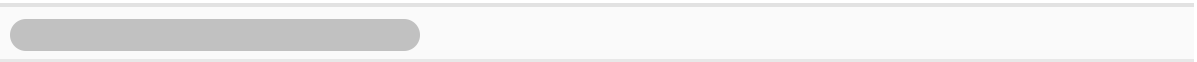
Type your answer here, replacing this text.

```
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 x 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
<b>16395</b>	16396	320000.0	2	1	2	36	0	0	0
<b>21448</b>	21449	440000.0	2	1	2	30	-1	-1	-1
<b>20034</b>	20035	160000.0	2	3	1	44	-2	-2	-2
<b>25755</b>	25756	120000.0	2	2	1	30	0	0	0
<b>1438</b>	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*

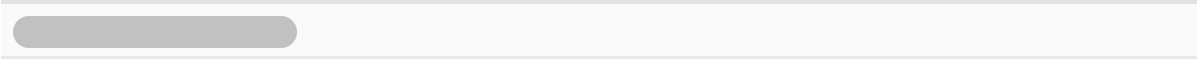
1.1 train\_df.describe() gives us a statistic summary of all columns in the train dataframe.

In [5]: train\_df.describe()

Out [5]:

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

8 rows x 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
```

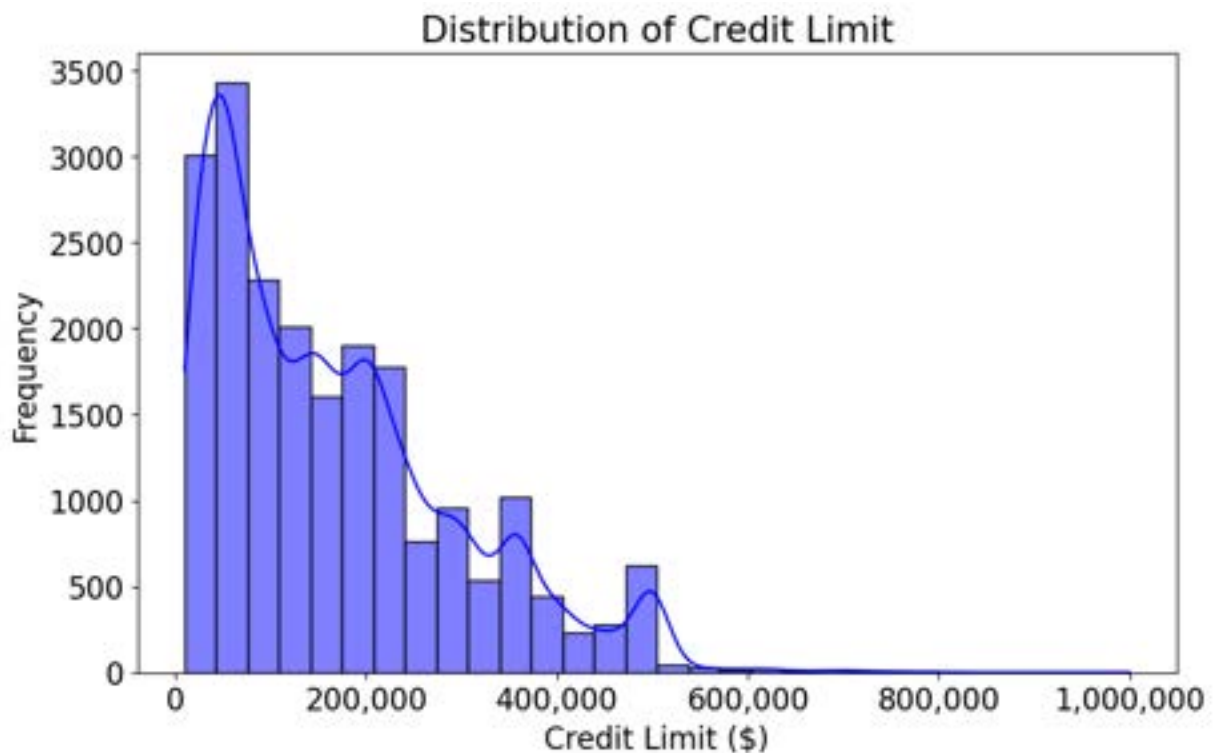
```
In [8]: median_limit_bal = train_df['LIMIT_BAL'].median()
print("Median Credit Limit (LIMIT_BAL):\n", median_limit_bal)
```

```
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.

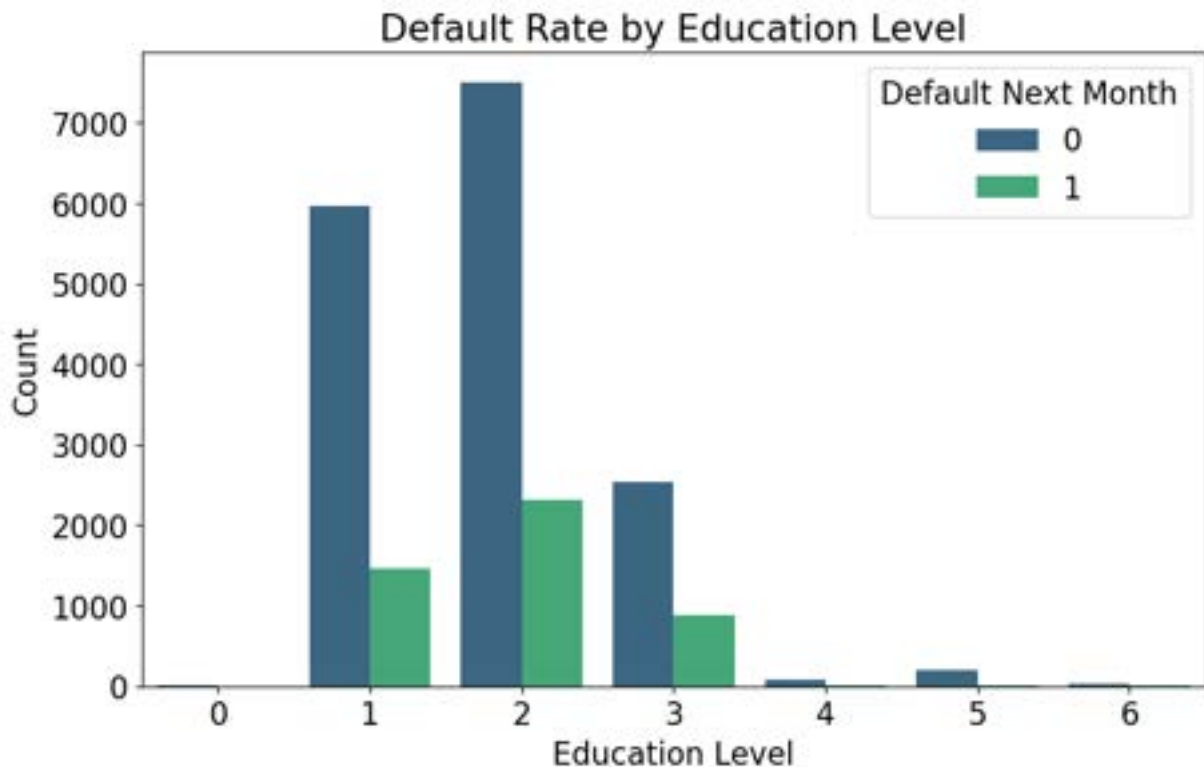
```
In [9]: import seaborn as sns
import matplotlib.ticker as mtick

# 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]: 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()
```



### 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]: 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_AMT6']
# 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, 1)

# 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) / train_df['avg_payment_amt']).mean(axis=1)
```

train\_df

Out[12]:

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

21000 rows × 30 columns

In [13]:

```

# 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_0', 'PAY_2', 'PAY_3', 'PAY_4', 'PAY_5', 'PAY_6']
# 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
<b>25665</b>	25666	40000.0	2	2	2	26	-1	0	0
<b>16464</b>	16465	80000.0	2	3	1	59	0	0	0
<b>22386</b>	22387	170000.0	2	1	2	30	2	2	2
<b>10149</b>	10150	200000.0	2	2	1	41	-2	-2	-2
<b>8729</b>	8730	50000.0	1	2	1	43	0	0	0
...	...	...	...	...	...	...	...	...	...
<b>17548</b>	17549	60000.0	2	2	1	48	0	0	0
<b>11459</b>	11460	310000.0	1	2	1	43	-1	-1	-1
<b>6608</b>	6609	10000.0	2	2	2	22	0	0	0
<b>2414</b>	2415	30000.0	1	2	1	38	1	-1	-1
<b>14757</b>	14758	30000.0	2	3	1	24	2	0	0

9000 rows x 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.

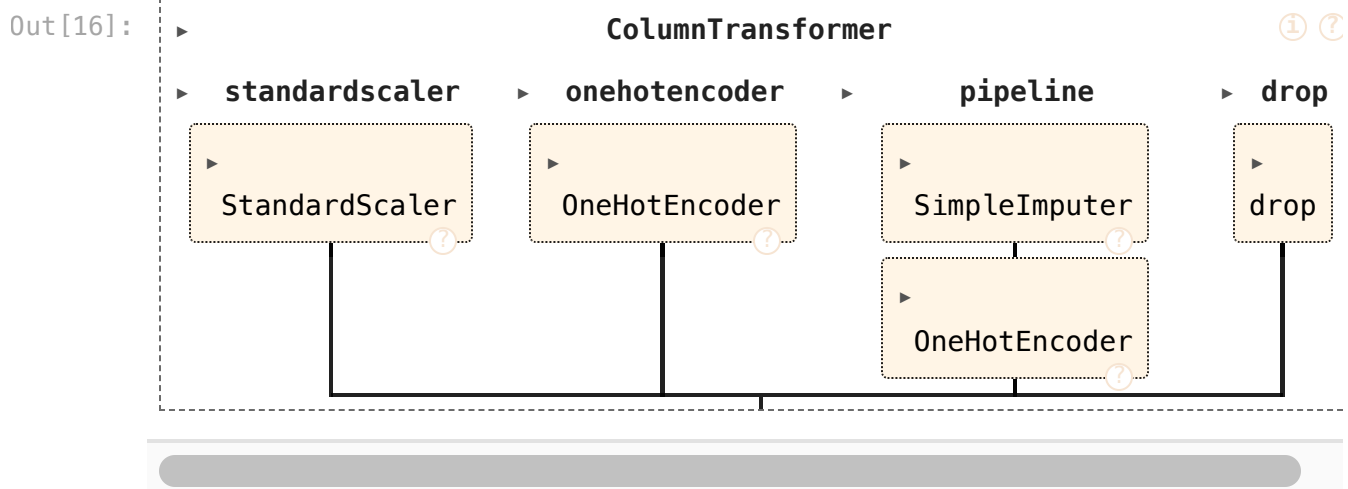
**Solution\_5**

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=False)
categorical_transformer = make_pipeline(
    SimpleImputer(strategy="most_frequent"), OneHotEncoder(handle_unknown="ignore")
)
```

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



## 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_scores.iloc[i])))

        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_train_score=True)
dummy_scores = pd.DataFrame(dummy_scores).T
dummy_scores.index = ['dummy']
dummy_scores
```

```
Out[19]:
```

	fit_time	score_time	test_score	train_score
<b>dummy</b>	0.001 (+/- 0.000)	0.001 (+/- 0.000)	0.221 (+/- 0.022)	0.223 (+/- 0.009)

## 7. Linear models

---



rubric={points:10}

**Your tasks:**

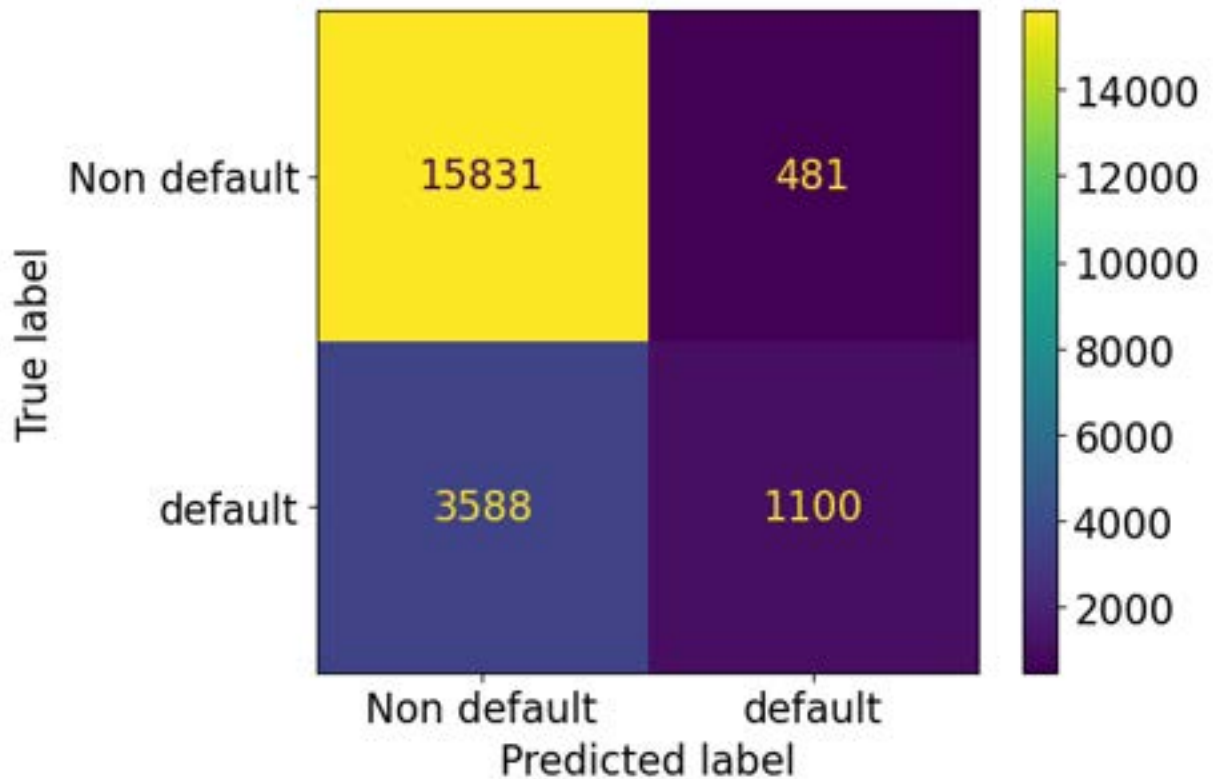
1. Try a linear model as a first real attempt.
2. Carry out hyperparameter tuning to explore different values for the complexity hyperparameter.
3. Report cross-validation scores along with standard deviation.
4. Summarize your results.

**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",  
);
```

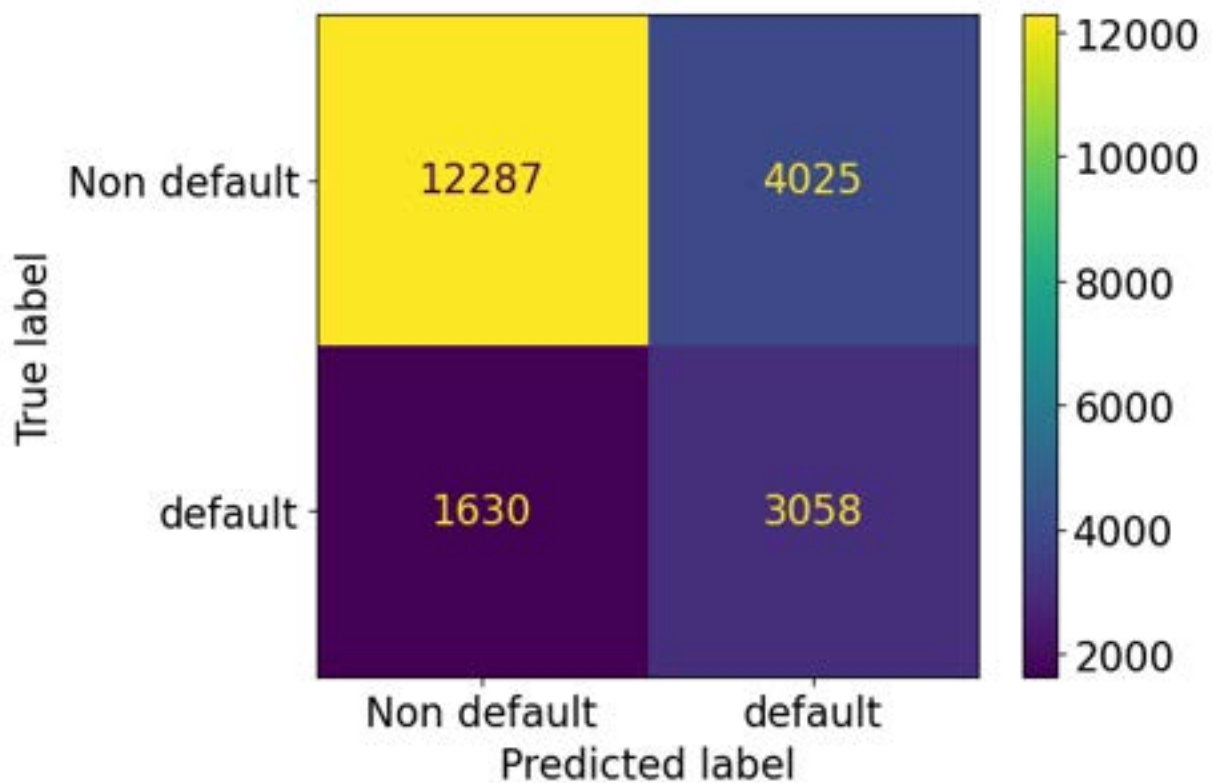


```
In [21]: lr_results = mean_std_cross_val_scores(pipe_lr, X_train, y_train, return_train_score=True)
lr_results = pd.DataFrame(lr_results).T
lr_results.index = ['Logistic regression']
lr_results
```

```
Out[21]:
```

	fit_time	score_time	test_score	train_score
<b>Logistic regression</b>	0.040 (+/- 0.004)	0.004 (+/- 0.000)	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_train)
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
<b>Logistic regression_Balanced</b>	0.046 (+/- 0.003)	0.004 (+/- 0.000)	0.653 (+/- 0.011)	0.653 (+/- 0.003)

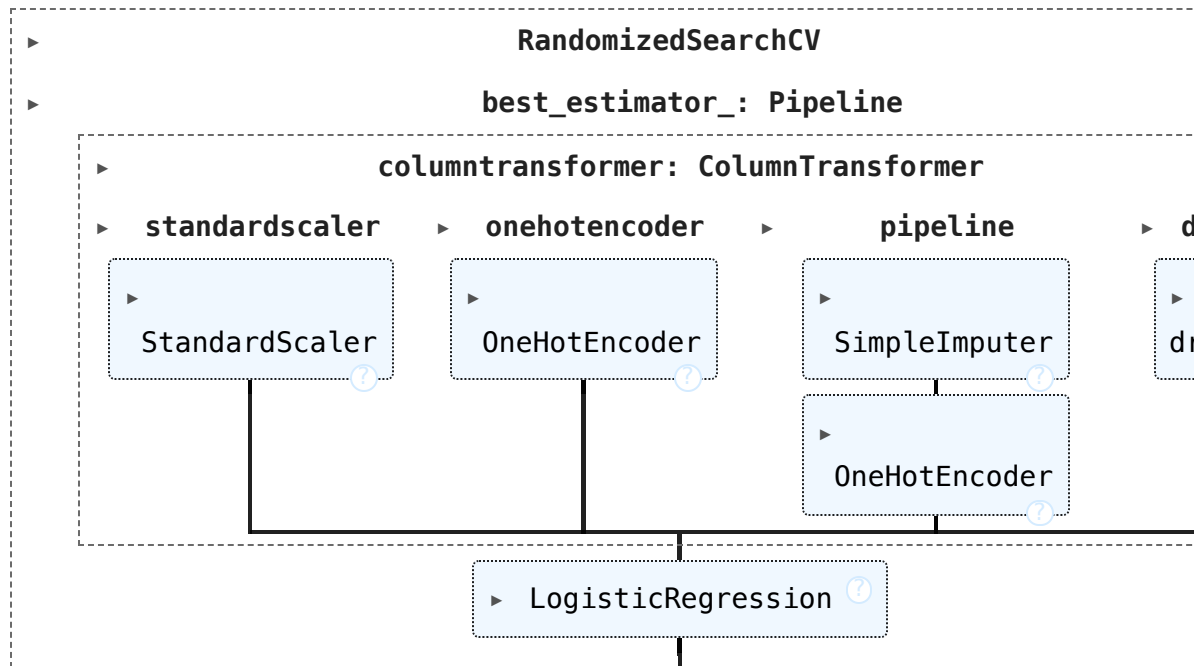
## 2. Hyperparameter Optimization

```
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]:

In [25]: `random_search.best_score_`Out [25]: `0.6537927833010583`In [26]: `random_search.best_params_`Out [26]: `{'logisticregression__C': 0.14240691453096987}`

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
<b>Logistic regression_Balanced</b>	0.045 (+/- 0.005)	0.004 (+/- 0.000)	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, y)

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

```

Out [28]:

	fit_time	score_time	test_score	train_score
<b>LR_Balanced_Tuned</b>	0.040 (+/- 0.000)	0.004 (+/- 0.000)	0.654 (+/- 0.011)	0.653 (+/- 0.002)

#### 4. Summary:

- Because there is a class imbalance so we used logistic regression with equal class weight 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 be 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, HistGradientBoostingClassifier

pipe_rf = make_pipeline(
    preprocessor, RandomForestClassifier(class_weight="balanced", random_state=123)
)
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="balanced")
)

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_score
    )
```

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

Out [32]:

	fit_time	score_time	test_score	train_score
<b>random forest</b>	2.909 (+/- 0.040)	0.045 (+/- 0.001)	0.341 (+/- 0.009)	1.000 (+/- 0.000)
<b>K-Nearest Neighbors</b>	0.010 (+/- 0.001)	0.053 (+/- 0.009)	0.363 (+/- 0.005)	0.472 (+/- 0.005)
<b>LightGBM</b>	0.252 (+/- 0.001)	0.007 (+/- 0.000)	0.375 (+/- 0.013)	0.457 (+/- 0.004)
<b>LightGBM_balanced</b>	0.256 (+/- 0.001)	0.007 (+/- 0.000)	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 improve with feature selection, as the balanced logistic regression model with selected features (0.641) has relatively lower test score than the one without selected features (0.653).

```
In [33]: from sklearn.feature_selection import RFECV
from sklearn.linear_model import LogisticRegression
import numpy as np
import itertools

# Fit the preprocessor on the training data
preprocessor.fit(X_train)

# Transform the training and testing data
X_train_processed = preprocessor.transform(X_train)
X_test_processed = preprocessor.transform(X_test)

# Use the processed data for feature selection
```



```

# Feature selection using RFECV
rfe_cv = RFECV(LogisticRegression(max_iter=2000), cv=10, scoring = recall_sc
rfe_cv.fit(X_train_processed, y_train)

# Get the support mask (selected features)
selected_features_mask = rfe_cv.support_

# Extract feature names from the preprocessor
numeric_feature_names = numeric_features
binary_feature_names = preprocessor.named_transformers_['onehotencoder'].get
categorical_feature_names = preprocessor.named_transformers_['pipeline'].nam
all_feature_names = list(itertools.chain(numeric_feature_names, binary_featu

# Get the names of the selected features
selected_feature_names = [name for name, selected in zip(all_feature_names,

# Create DataFrames with only the selected features
X_train_selected_df = pd.DataFrame(X_train_processed[:, selected_features_ma
X_test_selected_df = pd.DataFrame(X_test_processed[:, selected_features_mask

# Evaluate the logistic regression model on the selected features
logistic_model = LogisticRegression(max_iter=2000, class_weight="balanced")
mean_std_scores = mean_std_cross_val_scores(logistic_model, X_train_selected

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

```

Out[33]:

	fit_time	score_time	test_score	train_score
<b>LR_Balanced_Selected</b>	0.012 (+/- 0.002)	0.002 (+/- 0.000)	0.641 (+/- 0.008)	0.641 (+/- 0.002)

In [34]:

```

combined_scores = pd.concat([dummy_scores, lr_results, lr_balanced_results,
combined_scores

```

Out [34]:

	fit_time	score_time	test_score	train_score
<b>dummy</b>	0.001 (+/- 0.000)	0.001 (+/- 0.000)	0.221 (+/- 0.022)	0.223 (+/- 0.009)
<b>Logistic regression</b>	0.040 (+/- 0.004)	0.004 (+/- 0.000)	0.238 (+/- 0.023)	0.236 (+/- 0.012)
<b>Logistic regression_Balanced</b>	0.045 (+/- 0.005)	0.004 (+/- 0.000)	0.653 (+/- 0.011)	0.653 (+/- 0.003)
<b>LR_Balanced_Tuned</b>	0.040 (+/- 0.000)	0.004 (+/- 0.000)	0.654 (+/- 0.011)	0.653 (+/- 0.002)
<b>random forest</b>	2.909 (+/- 0.040)	0.045 (+/- 0.001)	0.341 (+/- 0.009)	1.000 (+/- 0.000)
<b>K-Nearest Neighbors</b>	0.010 (+/- 0.001)	0.053 (+/- 0.009)	0.363 (+/- 0.005)	0.472 (+/- 0.005)
<b>LightGBM</b>	0.252 (+/- 0.001)	0.007 (+/- 0.000)	0.375 (+/- 0.013)	0.457 (+/- 0.004)
<b>LightGBM_balanced</b>	0.256 (+/- 0.001)	0.007 (+/- 0.000)	0.620 (+/- 0.021)	0.788 (+/- 0.007)
<b>LR_Balanced_Selected</b>	0.012 (+/- 0.002)	0.002 (+/- 0.000)	0.641 (+/- 0.008)	0.641 (+/- 0.002)

## 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