

# ML\_regression\_and\_classification\_project

August 11, 2024

**Description** A fictional stakeholder shared data containing client transaction sums and some features. All the data is stored in the train\_dataset.csv file, where transaction sums are under the int\_target field. This field is going to be the target to be predicted in this project. The project has a single goal: to predict the target variable. Nevertheless, it will be divided into several subparts for convenience.

```
[48]: # importing necessary libraries
import numpy as np
import pandas as pd

from sklearn.model_selection import train_test_split
from sklearn.model_selection import GridSearchCV

from sklearn.preprocessing import MinMaxScaler
from sklearn.preprocessing import StandardScaler
from sklearn.preprocessing import OneHotEncoder

from sklearn.metrics import accuracy_score
from sklearn.metrics import mean_squared_error as MSE
from sklearn.metrics import classification_report
from sklearn.metrics import confusion_matrix

from sklearn.linear_model import LogisticRegression
from sklearn.linear_model import LinearRegression
from statsmodels.regression.linear_model import OLS

from sklearn.decomposition import PCA
from sklearn.svm import SVC
from sklearn.neighbors import KNeighborsClassifier

from sklearn.tree import DecisionTreeClassifier
from sklearn.ensemble import RandomForestClassifier
from sklearn.ensemble import AdaBoostClassifier
from sklearn.ensemble import GradientBoostingClassifier
from lightgbm import LGBMClassifier
import xgboost as xgb

import plotly.graph_objs as go
```

```
import matplotlib.pyplot as plt
import seaborn as sns
```

1) Importing data

```
[3]: dataset = pd.read_csv('train_dataset.csv')
```

```
[4]: dataset.head()
```

```
[4]:   feature_1  feature_2  feature_3  feature_4  feature_5  feature_6  \
0  -0.043457 -0.027918   1.473594 -10.309556  -1.000000  145.373247
1  -1.535978  0.999851  23.019753  38.430092   1.000000 -110.045207
2  11.567708  0.608310   1.873323   6.412821   0.999995 -280.186852
3   9.220074  0.999864   5.041939 -12.692120  -1.000000  164.773793
4  -3.025434 -0.594711   5.912259 -14.922129  -1.000000   9.546673

      feature_7  feature_8  feature_9  feature_10  ...  feature_12  feature_13  \
0     0.001888   2.555947  -64.208613    0.470329  ...          B    83.317907
1     2.359228   0.571917 -110.318854   -2.620978  ...          D   114.881621
2    133.811870   1.160380 -206.238816   -1.006502  ...          B   145.833459
3     85.009772  -1.401234   92.624296    2.712139  ...          C    98.150579
4     9.153251  -0.157479 -175.895074   -0.832795  ...          C    94.519135

      feature_14  feature_15  feature_16  feature_17  feature_18  feature_19  \
0     85.385030  120.550153   59.736520   64.325558  107.833753   85.853573
1     79.711086   68.505986  116.330797   75.365722  104.229006   90.842232
2    115.159171   64.417992   98.846694   66.381070   82.118547   96.837671
3     78.620349  120.176496  104.297092   89.318362   73.791588   71.870016
4    123.286154   75.458961  111.831234  111.871202   99.733563  100.697358

      feature_20  int_target
0    104.466108        17353
1     95.232654         1214
2    103.887599        16094
3     68.397346        24168
4     97.568328         1105
```

[5 rows x 21 columns]

2. Preliminary data analysis

```
[5]: len(dataset.index)
```

```
[5]: 9000
```

```
[6]: dataset.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 9000 entries, 0 to 8999
```

Data columns (total 21 columns):

| #  | Column     | Non-Null Count | Dtype   |
|----|------------|----------------|---------|
| 0  | feature_1  | 9000 non-null  | float64 |
| 1  | feature_2  | 9000 non-null  | float64 |
| 2  | feature_3  | 9000 non-null  | float64 |
| 3  | feature_4  | 9000 non-null  | float64 |
| 4  | feature_5  | 9000 non-null  | float64 |
| 5  | feature_6  | 9000 non-null  | float64 |
| 6  | feature_7  | 9000 non-null  | float64 |
| 7  | feature_8  | 9000 non-null  | float64 |
| 8  | feature_9  | 9000 non-null  | float64 |
| 9  | feature_10 | 9000 non-null  | float64 |
| 10 | feature_11 | 9000 non-null  | float64 |
| 11 | feature_12 | 9000 non-null  | object  |
| 12 | feature_13 | 9000 non-null  | float64 |
| 13 | feature_14 | 9000 non-null  | float64 |
| 14 | feature_15 | 9000 non-null  | float64 |
| 15 | feature_16 | 9000 non-null  | float64 |
| 16 | feature_17 | 9000 non-null  | float64 |
| 17 | feature_18 | 9000 non-null  | float64 |
| 18 | feature_19 | 9000 non-null  | float64 |
| 19 | feature_20 | 9000 non-null  | float64 |
| 20 | int_target | 9000 non-null  | int64   |

dtypes: float64(19), int64(1), object(1)

memory usage: 1.4+ MB

```
[6]: dataset.describe()
```

```
[6]:
```

|       | feature_1   | feature_2   | feature_3   | feature_4   | feature_5   | \ |
|-------|-------------|-------------|-------------|-------------|-------------|---|
| count | 9000.000000 | 9000.000000 | 9000.000000 | 9000.000000 | 9000.000000 |   |
| mean  | 9.892699    | 0.051539    | 13.691801   | -5.280291   | -0.190470   |   |
| std   | 16.005016   | 0.708241    | 77.058634   | 26.814772   | 0.966878    |   |
| min   | -55.674815  | -1.000000   | 0.002306    | -90.266818  | -1.000000   |   |
| 25%   | -0.590266   | -0.652740   | 0.891789    | -24.131903  | -1.000000   |   |
| 50%   | 10.032527   | 0.108589    | 2.771781    | -6.535842   | -0.999996   |   |
| 75%   | 20.577624   | 0.750147    | 8.313018    | 12.151035   | 1.000000    |   |
| max   | 71.408515   | 1.000000    | 4195.458665 | 96.967006   | 1.000000    |   |

|       | feature_6   | feature_7   | feature_8   | feature_9   | feature_10  | \ |
|-------|-------------|-------------|-------------|-------------|-------------|---|
| count | 9000.000000 | 9000.000000 | 9000.000000 | 9000.000000 | 9000.000000 |   |
| mean  | 33.328407   | 353.997548  | 0.322612    | -69.936887  | -0.334325   |   |
| std   | 193.282240  | 482.976481  | 1.928390    | 130.805512  | 1.802465    |   |
| min   | -782.730110 | 0.000032    | -8.785694   | -562.750841 | -6.678488   |   |
| 25%   | -88.607265  | 38.261808   | -0.867972   | -155.988559 | -1.569338   |   |
| 50%   | 47.154805   | 169.202144  | 0.454511    | -80.249681  | -0.374527   |   |
| 75%   | 165.882796  | 476.800547  | 1.636826    | 10.046369   | 0.874218    |   |

|     |            |             |          |            |          |
|-----|------------|-------------|----------|------------|----------|
| max | 689.183091 | 5099.175987 | 7.563691 | 454.622574 | 6.605742 |
|-----|------------|-------------|----------|------------|----------|

|       | feature_11  | feature_13  | feature_14  | feature_15  | feature_16  | \ |
|-------|-------------|-------------|-------------|-------------|-------------|---|
| count | 9000.000000 | 9000.000000 | 9000.000000 | 9000.000000 | 9000.000000 |   |
| mean  | 0.328602    | 99.738087   | 99.862703   | 99.764452   | 99.958384   |   |
| std   | 1.949699    | 19.970516   | 20.059934   | 19.816391   | 20.073887   |   |
| min   | -6.817115   | 23.296434   | 21.404151   | 20.462943   | 31.695377   |   |
| 25%   | -0.997933   | 86.416262   | 86.377084   | 86.184816   | 86.146061   |   |
| 50%   | 0.361400    | 99.753333   | 99.914720   | 99.938163   | 99.956294   |   |
| 75%   | 1.690059    | 113.192601  | 113.550065  | 113.274773  | 113.503300  |   |
| max   | 7.314190    | 169.972815  | 169.211534  | 176.896232  | 180.756967  |   |

|       | feature_17  | feature_18  | feature_19  | feature_20  | int_target   |
|-------|-------------|-------------|-------------|-------------|--------------|
| count | 9000.000000 | 9000.000000 | 9000.000000 | 9000.000000 | 9000.000000  |
| mean  | 100.053203  | 100.209235  | 100.150711  | 99.819907   | 14962.016000 |
| std   | 20.005191   | 19.839667   | 20.084744   | 20.072153   | 8507.039575  |
| min   | 25.278459   | 6.108243    | 29.523073   | 28.591199   | 501.000000   |
| 25%   | 86.513425   | 86.753705   | 86.684625   | 86.163540   | 7461.750000  |
| 50%   | 99.644334   | 100.166669  | 100.179503  | 99.942090   | 14878.500000 |
| 75%   | 113.486201  | 113.389489  | 113.657307  | 113.275163  | 22404.000000 |
| max   | 182.022092  | 173.917343  | 171.521586  | 178.220128  | 29631.000000 |

3. Dividing the dataset into target and features as well as into train and test datasets

```
[6]: x = dataset.drop(columns=['int_target'])
     y = dataset['int_target']
```

```
[7]: x_tr,x_t,y_tr,y_t = train_test_split(x,y,test_size=0.3,random_state=42)
```

```
[39]: x_tr.head()
```

```
[39]:
```

|      | feature_1 | feature_2 | feature_3  | feature_4  | feature_5 | feature_6   | \ |
|------|-----------|-----------|------------|------------|-----------|-------------|---|
| 7581 | -6.077065 | 0.274977  | 4.922146   | -13.633130 | -1.00000  | 137.750101  |   |
| 8484 | 13.478385 | -0.142569 | 29.020381  | -2.015988  | -0.96514  | -277.913183 |   |
| 6215 | 11.137173 | 0.944967  | 2.699315   | -19.220130 | -1.00000  | 257.743532  |   |
| 6884 | 27.832998 | 0.553589  | 2.006801   | -49.907437 | -1.00000  | 18.190559   |   |
| 3647 | 6.607199  | 0.312878  | 561.074357 | 12.063850  | 1.00000   | -113.262260 |   |

|      | feature_7  | feature_8 | feature_9   | feature_10 | feature_11 | feature_12 | \ |
|------|------------|-----------|-------------|------------|------------|------------|---|
| 7581 | 36.930724  | 0.834688  | -89.663914  | -0.407471  | -0.595614  |            | D |
| 8484 | 181.666863 | -0.022423 | -182.948151 | -2.506720  | -3.411167  |            | B |
| 6215 | 124.036614 | 1.277240  | -86.642872  | -0.461344  | -0.794835  |            | A |
| 6884 | 774.675784 | 0.444103  | -156.934461 | -1.045558  | 1.035041   |            | C |
| 3647 | 43.655079  | -1.820279 | -202.625263 | -3.991031  | -1.689960  |            | C |

|      | feature_13 | feature_14 | feature_15 | feature_16 | feature_17 | feature_18 | \ |
|------|------------|------------|------------|------------|------------|------------|---|
| 7581 | 74.966634  | 108.817839 | 121.177814 | 128.620874 | 77.857225  | 95.340076  |   |
| 8484 | 95.358574  | 139.468850 | 110.479326 | 113.688174 | 108.933680 | 70.552332  |   |

|      |            |            |            |            |            |            |
|------|------------|------------|------------|------------|------------|------------|
| 6215 | 89.014381  | 125.106467 | 93.888226  | 109.688439 | 105.704562 | 90.203266  |
| 6884 | 149.900687 | 113.287872 | 100.622504 | 99.947717  | 81.521003  | 118.928950 |
| 3647 | 75.481512  | 101.213242 | 92.613375  | 105.406862 | 106.978952 | 111.314922 |

|      | feature_19 | feature_20 |
|------|------------|------------|
| 7581 | 99.516410  | 106.408905 |
| 8484 | 102.334319 | 113.123080 |
| 6215 | 109.077749 | 106.867140 |
| 6884 | 121.942998 | 98.779535  |
| 3647 | 89.233944  | 105.230220 |

- There is one categorical feature in the dataset that has a string format. It will be changed to a numeric format.

```
[12]: dataset['feature_12'].value_counts()
```

```
[12]: feature_12
B    2288
C    2265
D    2230
A    2217
Name: count, dtype: int64
```

```
[8]: from sklearn.preprocessing import OneHotEncoder as OHE
ohe = OHE(drop='first', handle_unknown='ignore')
```

```
[9]: etr_12 = ohe.fit_transform(x_tr[['feature_12']]).toarray()
```

```
[10]: et_12 = ohe.transform(x_t[['feature_12']]).toarray()
```

```
[11]: etr_12 = pd.DataFrame(data = etr_12, columns=ohe.categories_[0][1:])
et_12 = pd.DataFrame(data = et_12, columns=ohe.categories_[0][1:])
etr_12.index = x_tr.index
et_12.index = x_t.index
```

```
[13]: print(etr_12.head())
print(et_12.head())
```

|      | B   | C   | D   |
|------|-----|-----|-----|
| 7581 | 0.0 | 0.0 | 1.0 |
| 8484 | 1.0 | 0.0 | 0.0 |
| 6215 | 0.0 | 0.0 | 0.0 |
| 6884 | 0.0 | 1.0 | 0.0 |
| 3647 | 0.0 | 1.0 | 0.0 |

|      | B   | C   | D   |
|------|-----|-----|-----|
| 7940 | 0.0 | 0.0 | 0.0 |
| 1162 | 0.0 | 1.0 | 0.0 |
| 582  | 1.0 | 0.0 | 0.0 |

```
4081  0.0  1.0  0.0
8412  0.0  0.0  1.0
```

```
[12]: x_tr = pd.concat([x_tr,etr_12],axis=1)
x_tr.drop(columns=['feature_12'],inplace=True)
x_t = pd.concat([x_t,et_12],axis=1)
x_t.drop(columns=['feature_12'],inplace=True)
print(x_tr.head())
print(x_t.head())
```

|      | feature_1 | feature_2 | feature_3  | feature_4  | feature_5 | feature_6   | \ |
|------|-----------|-----------|------------|------------|-----------|-------------|---|
| 7581 | -6.077065 | 0.274977  | 4.922146   | -13.633130 | -1.00000  | 137.750101  |   |
| 8484 | 13.478385 | -0.142569 | 29.020381  | -2.015988  | -0.96514  | -277.913183 |   |
| 6215 | 11.137173 | 0.944967  | 2.699315   | -19.220130 | -1.00000  | 257.743532  |   |
| 6884 | 27.832998 | 0.553589  | 2.006801   | -49.907437 | -1.00000  | 18.190559   |   |
| 3647 | 6.607199  | 0.312878  | 561.074357 | 12.063850  | 1.00000   | -113.262260 |   |

|      | feature_7  | feature_8 | feature_9   | feature_10 | ... | feature_14 | \ |
|------|------------|-----------|-------------|------------|-----|------------|---|
| 7581 | 36.930724  | 0.834688  | -89.663914  | -0.407471  | ... | 108.817839 |   |
| 8484 | 181.666863 | -0.022423 | -182.948151 | -2.506720  | ... | 139.468850 |   |
| 6215 | 124.036614 | 1.277240  | -86.642872  | -0.461344  | ... | 125.106467 |   |
| 6884 | 774.675784 | 0.444103  | -156.934461 | -1.045558  | ... | 113.287872 |   |
| 3647 | 43.655079  | -1.820279 | -202.625263 | -3.991031  | ... | 101.213242 |   |

|      | feature_15 | feature_16 | feature_17 | feature_18 | feature_19 | feature_20 | \ |
|------|------------|------------|------------|------------|------------|------------|---|
| 7581 | 121.177814 | 128.620874 | 77.857225  | 95.340076  | 99.516410  | 106.408905 |   |
| 8484 | 110.479326 | 113.688174 | 108.933680 | 70.552332  | 102.334319 | 113.123080 |   |
| 6215 | 93.888226  | 109.688439 | 105.704562 | 90.203266  | 109.077749 | 106.867140 |   |
| 6884 | 100.622504 | 99.947717  | 81.521003  | 118.928950 | 121.942998 | 98.779535  |   |
| 3647 | 92.613375  | 105.406862 | 106.978952 | 111.314922 | 89.233944  | 105.230220 |   |

|      | B   | C   | D   |
|------|-----|-----|-----|
| 7581 | 0.0 | 0.0 | 1.0 |
| 8484 | 1.0 | 0.0 | 0.0 |
| 6215 | 0.0 | 0.0 | 0.0 |
| 6884 | 0.0 | 1.0 | 0.0 |
| 3647 | 0.0 | 1.0 | 0.0 |

[5 rows x 22 columns]

|      | feature_1 | feature_2 | feature_3 | feature_4  | feature_5 | feature_6   | \ |
|------|-----------|-----------|-----------|------------|-----------|-------------|---|
| 7940 | 21.788966 | 0.998660  | 4.177820  | 83.667049  | 1.000000  | -113.342663 |   |
| 1162 | 10.043187 | 0.296065  | 32.259720 | -5.346814  | -0.999955 | 287.185218  |   |
| 582  | 8.701901  | 0.636416  | 7.775502  | -20.108899 | -1.000000 | 78.748458   |   |
| 4081 | 6.182070  | 0.808258  | 62.072918 | 31.335802  | 1.000000  | -59.601221  |   |
| 8412 | 19.381290 | -0.393950 | 3.424870  | -18.088676 | -1.000000 | 475.163917  |   |

|      | feature_7  | feature_8 | feature_9  | feature_10 | ... | feature_14 | \ |
|------|------------|-----------|------------|------------|-----|------------|---|
| 7940 | 474.759022 | 0.249128  | 129.199477 | -5.012149  | ... | 129.395518 |   |

|      |            |           |             |           |     |            |
|------|------------|-----------|-------------|-----------|-----|------------|
| 1162 | 100.865600 | -1.152400 | 2.173456    | 1.658315  | ... | 60.563411  |
| 582  | 75.723079  | 0.740615  | -62.850303  | -0.454813 | ... | 112.716310 |
| 4081 | 38.217993  | -0.306800 | -277.342478 | -3.114036 | ... | 122.854305 |
| 8412 | 375.634383 | -1.346684 | 21.556730   | -0.329047 | ... | 99.175725  |

|      | feature_15 | feature_16 | feature_17 | feature_18 | feature_19 | feature_20 | \ |
|------|------------|------------|------------|------------|------------|------------|---|
| 7940 | 81.421299  | 117.976466 | 107.740527 | 100.519500 | 140.581802 | 118.778455 |   |
| 1162 | 133.962758 | 86.688118  | 91.062096  | 107.407778 | 99.737611  | 99.890741  |   |
| 582  | 129.104165 | 79.461983  | 129.088292 | 115.548850 | 91.866696  | 102.281620 |   |
| 4081 | 117.988280 | 81.842686  | 106.777292 | 95.309154  | 78.172080  | 76.360709  |   |
| 8412 | 76.437763  | 75.404877  | 77.515263  | 121.018053 | 132.856019 | 98.929321  |   |

|      | B   | C   | D   |
|------|-----|-----|-----|
| 7940 | 0.0 | 0.0 | 0.0 |
| 1162 | 0.0 | 1.0 | 0.0 |
| 582  | 1.0 | 0.0 | 0.0 |
| 4081 | 0.0 | 1.0 | 0.0 |
| 8412 | 0.0 | 0.0 | 1.0 |

[5 rows x 22 columns]

5. The scales of the features differ a lot. One general scaling will be utilized for the features.

```
[13]: #Considering the data distribution and exisistence of negative values and
      ↪outliers, StandardScaler was chosen in this case
scaler = StandardScaler()
x_tr_scaled = scaler.fit_transform(x_tr)
x_t_scaled = scaler.transform(x_t)
```

```
[14]: x_tr_scaled = pd.DataFrame(x_tr_scaled,columns=x_tr.columns)
      x_t_scaled = pd.DataFrame(x_t_scaled,columns=x_t.columns)
```

```
[15]: x_tr_scaled.index = x_tr.index
      x_t_scaled.index = x_t.index
```

6. Creating a linear regression model and its assessment using RMSE and other statistical measures.

```
[16]: linreg = LinearRegression()
      linreg.fit(x_tr_scaled,y_tr)
```

```
[16]: LinearRegression()
```

```
[17]: y_pred = linreg.predict(x_t_scaled)
      MSE(y_t,y_pred)**(1/2)
```

```
[17]: 6800.840397189446
```

```
[18]: ols = OLS(y_tr,x_tr_scaled)
```

```
[19]: results = ols.fit()
```

```
[20]: results.summary()
```

[20]:

|                   |                  |                              |           |
|-------------------|------------------|------------------------------|-----------|
| Dep. Variable:    | int_target       | R-squared (uncentered):      | 0.094     |
| Model:            | OLS              | Adj. R-squared (uncentered): | 0.091     |
| Method:           | Least Squares    | F-statistic:                 | 29.74     |
| Date:             | Sun, 11 Aug 2024 | Prob (F-statistic):          | 1.10e-117 |
| Time:             | 15:21:18         | Log-Likelihood:              | -70031.   |
| No. Observations: | 6300             | AIC:                         | 1.401e+05 |
| Df Residuals:     | 6278             | BIC:                         | 1.403e+05 |
| Df Model:         | 22               |                              |           |
| Covariance Type:  | nonrobust        |                              |           |

|            | coef       | std err | t      | P>  t | [0.025    | 0.975]    |
|------------|------------|---------|--------|-------|-----------|-----------|
| feature_1  | -1174.6468 | 284.486 | -4.129 | 0.000 | -1732.337 | -616.956  |
| feature_2  | 865.3265   | 207.325 | 4.174  | 0.000 | 458.899   | 1271.754  |
| feature_3  | 233.2799   | 220.647 | 1.057  | 0.290 | -199.264  | 665.824   |
| feature_4  | -2943.8266 | 428.817 | -6.865 | 0.000 | -3784.454 | -2103.199 |
| feature_5  | -727.0198  | 355.453 | -2.045 | 0.041 | -1423.828 | -30.211   |
| feature_6  | -150.5213  | 302.613 | -0.497 | 0.619 | -743.747  | 442.704   |
| feature_7  | -381.3816  | 268.574 | -1.420 | 0.156 | -907.879  | 145.116   |
| feature_8  | -183.7838  | 295.204 | -0.623 | 0.534 | -762.485  | 394.917   |
| feature_9  | 2907.6143  | 369.462 | 7.870  | 0.000 | 2183.343  | 3631.886  |
| feature_10 | 255.8017   | 268.206 | 0.954  | 0.340 | -269.973  | 781.577   |
| feature_11 | 2196.6574  | 261.298 | 8.407  | 0.000 | 1684.424  | 2708.890  |
| feature_13 | 73.9060    | 205.587 | 0.359  | 0.719 | -329.115  | 476.927   |
| feature_14 | 58.4509    | 205.508 | 0.284  | 0.776 | -344.415  | 461.317   |
| feature_15 | -33.0601   | 205.730 | -0.161 | 0.872 | -436.361  | 370.241   |
| feature_16 | 42.7527    | 205.798 | 0.208  | 0.835 | -360.681  | 446.186   |
| feature_17 | 69.7367    | 205.531 | 0.339  | 0.734 | -333.174  | 472.648   |
| feature_18 | 137.6595   | 205.569 | 0.670  | 0.503 | -265.326  | 540.645   |
| feature_19 | 7.3568     | 205.544 | 0.036  | 0.971 | -395.580  | 410.294   |
| feature_20 | 131.7587   | 205.877 | 0.640  | 0.522 | -271.830  | 535.348   |
| B          | 107.4957   | 253.574 | 0.424  | 0.672 | -389.596  | 604.588   |
| C          | 28.0881    | 252.900 | 0.111  | 0.912 | -467.683  | 523.859   |
| D          | 35.1902    | 253.501 | 0.139  | 0.890 | -461.759  | 532.139   |

|                |        |                   |          |
|----------------|--------|-------------------|----------|
| Omnibus:       | 60.221 | Durbin-Watson:    | 0.347    |
| Prob(Omnibus): | 0.000  | Jarque-Bera (JB): | 62.620   |
| Skew:          | 0.225  | Prob(JB):         | 2.52e-14 |
| Kurtosis:      | 3.189  | Cond. No.         | 5.08     |

Notes:

[1]  $R^2$  is computed without centering (uncentered) since the model does not contain a constant.

[2] Standard Errors assume that the covariance matrix of the errors is correctly specified.

**Update of the Task** The fictional stakeholder saw the results and thought the RMSE measure was too large: almost half the range of the target variable. The linear regression model might not



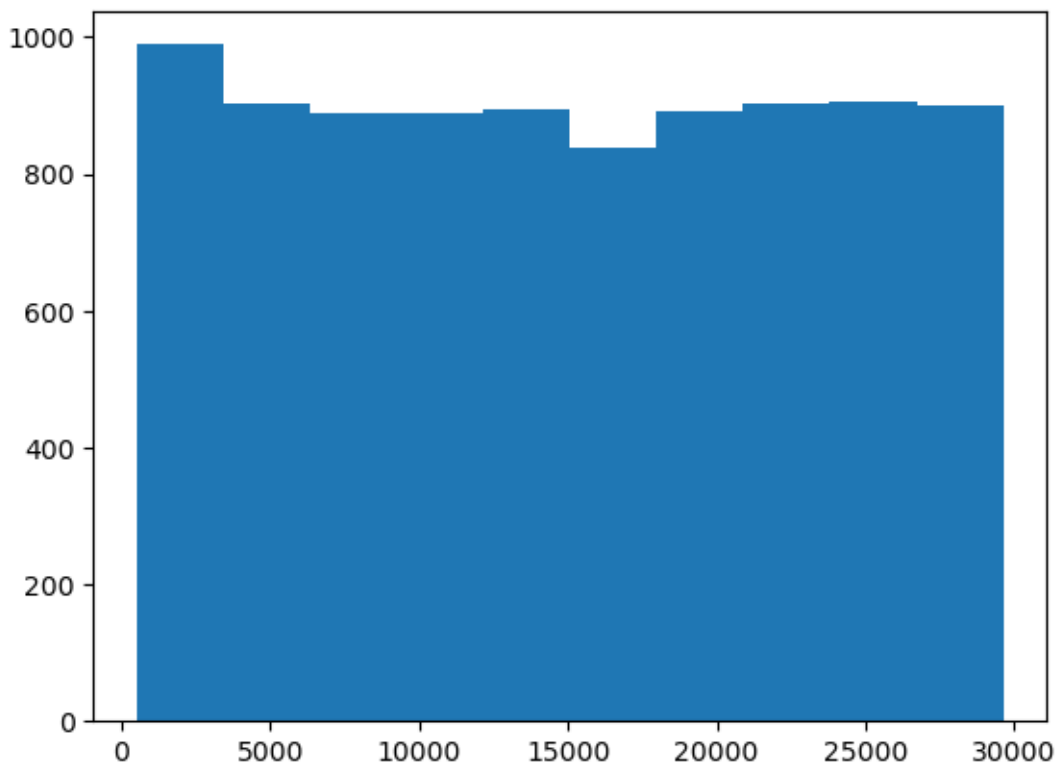
be of good fit for this task. The target variable could be divided into classes. That is to say we will not consider the absolute sums of transaction, but will look at which class the customer belongs to instead (high revenue, medium revenue, or low revenue). The stakeholder asked to show them how the division is going to be done.

7. Looking at the distribution of the target variable to divide the target into classes.

```
[60]: y.describe()
```

```
[60]: count      9000.000000
      mean      14962.016000
      std       8507.039575
      min       501.000000
      25%       7461.750000
      50%      14878.500000
      75%      22404.000000
      max      29631.000000
      Name: int_target, dtype: float64
```

```
[21]: plt.hist(dataset.int_target);
```



8. The distribution seems to be uniform: the target variable could be divided into three equal classes. For that respective percentiles have to be calculated (33% and 67%).

```
[23]: y_33 = np.quantile(y,0.33)
      y_67 = np.quantile(y,0.67)
      print(y_33,y_67)
```

9871.67 20120.66

9. The stakeholder agreed to divide the target variable into classes using these percentiles, but they suggested to round the number. That way the division into classes in the target variable will use the following thresholds: 10000 and below will mean the low revenue class, 20000 and above will mean the high revenue class, and everything else will be the medium revenue class.

```
[23]: # Low revenue class
      class_0 = [0 for i in y if i <= 10000]
      class_0 = pd.Series(class_0)
      class_0.index = y[y<=10000].index
```

```
[24]: # Medium revenue class
      class_1 = [1 for i in y if (i > 10000 and i < 20000)]
      class_1 = pd.Series(class_1)
      class_1.index = y[y[10000 < y] & y[20000 > y]].index
```

C:\Users\azima\AppData\Local\Temp\ipykernel\_1360\2887226351.py:4: FutureWarning: Operation between non boolean Series with different indexes will no longer return a boolean result in a future version. Cast both Series to object type to maintain the prior behavior.

```
class_1.index = y[y[10000 < y] & y[20000 > y]].index
```

```
[26]: # High revenue class
      class_2 = [2 for i in y if i >= 20000]
      class_2 = pd.Series(class_2)
      class_2.index = y[y>=20000].index
```

```
[27]: y_class = pd.concat([class_0,class_1,class_2]).sort_index()
```

```
[28]: y_class.head()
```

```
[28]: 0    1
      1    0
      2    1
      3    2
      4    0
      dtype: int64
```

10. Next steps will involve creating a logistic regression model. This model will serve as the baseline and a filter for the features. The created model will have the Lasso regularization with the regularization coefficient of 0.01. The model's performance will be measured using the accuracy metric.

```
[29]: logreg = LogisticRegression(penalty='l1',C=0.01, solver='liblinear')
```

```
[30]: y_class_tr = pd.Series([y_class[i] for i in y_tr.index],index=y_tr.index)
      y_class_t = pd.Series([y_class[i] for i in y_t.index],index=y_t.index)
```

```
[31]: y_class_tr.head()
```

```
[31]: 7581    1
      8484    0
      6215    1
      6884    1
      3647    0
      dtype: int64
```

```
[32]: y_tr
```

```
[32]: 7581    12868
      8484    1895
      6215    11398
      6884    13635
      3647     3505
      ...
      5734    22641
      5191    24348
      5390    10498
      860     19158
      7270    16376
      Name: int_target, Length: 6300, dtype: int64
```

```
[33]: logreg.fit(x_tr_scaled,y_class_tr)
      y_class_pred = logreg.predict(x_t_scaled)
      accuracy_score(y_class_t,y_class_pred)
```

```
[33]: 0.7355555555555555
```

11. Looking at the resulting feature coefficients of the model

```
[34]: logreg.coef_, logreg.intercept_
```

```
[34]: (array([[ 0.3796715 , -0.1871515 ,  0.          ,  1.07461921,  0.08020088,
           -0.22876785,  0.          , -0.01357207, -0.49872147,  0.01267205,
           -0.73448159,  0.          ,  0.          ,  0.          ,  0.          ,
           0.          ,  0.          ,  0.          ,  0.          ,  0.          ,
           0.          ,  0.          ],
          [-0.23148403,  0.          ,  0.          , -0.84576659,  0.          ,
           0.3754246 ,  0.          ,  0.76255836,  0.          , -0.74824977,
           0.48124861,  0.          ,  0.          ,  0.          ,  0.          ,
           0.          ,  0.          ,  0.          ,  0.          ,  0.          ,
           0.          ,  0.          ]])
```

```

0.          , 0.          ],
[ 0.          , 0.12522877, 0.          , -0.00989079, -0.17520227,
 0.          , -0.15513951, -0.64849311, 0.37998327, 0.67332255,
 0.09582758, 0.          , 0.          , 0.          , 0.          ,
 0.          , 0.          , 0.          , 0.          , 0.          ,
 0.          , 0.          ]]),
array([-0.85347889, -0.84132347, -0.7702121 ]))

```

12. The features that have coefficients of zero for all the classes should not be as useful. Only the non-zero coefficients will be chosen for future analysis with a new dataset created.

```

[35]: coef_dict = {logreg.feature_names_in_[i]: 'y' for i in range(len(logreg.
    ↪feature_names_in_)) if (logreg.coef_[0][i] != 0
        or logreg.coef_[1][i] != 0 or logreg.coef_[2][i] != 0)}
coef_dict

```

```

[35]: {'feature_1': 'y',
      'feature_2': 'y',
      'feature_4': 'y',
      'feature_5': 'y',
      'feature_6': 'y',
      'feature_7': 'y',
      'feature_8': 'y',
      'feature_9': 'y',
      'feature_10': 'y',
      'feature_11': 'y'}

```

```

[36]: list(coef_dict)

```

```

[36]: ['feature_1',
      'feature_2',
      'feature_4',
      'feature_5',
      'feature_6',
      'feature_7',
      'feature_8',
      'feature_9',
      'feature_10',
      'feature_11']

```

```

[37]: x_tr_scaled = x_tr_scaled[list(coef_dict)]
      x_t_scaled = x_t_scaled[list(coef_dict)]
      x_tr_scaled.head()

```

```

[37]:      feature_1  feature_2  feature_4  feature_5  feature_6  feature_7  \
7581   -0.991374   0.322929  -0.317627  -0.837569   0.536843  -0.654027
8484    0.222688  -0.265324   0.114265  -0.801533  -1.603270  -0.358543
6215    0.077338   1.266835  -0.525335  -0.837569   1.154649  -0.476197

```

|      |           |          |           |           |           |           |
|------|-----------|----------|-----------|-----------|-----------|-----------|
| 6884 | 1.113866  | 0.715448 | -1.666202 | -0.837569 | -0.078730 | 0.852107  |
| 3647 | -0.203896 | 0.376326 | 0.637714  | 1.229929  | -0.755537 | -0.640299 |

|      | feature_8 | feature_9 | feature_10 | feature_11 |
|------|-----------|-----------|------------|------------|
| 7581 | 0.253938  | -0.149507 | -0.036008  | -0.464830  |
| 8484 | -0.190346 | -0.862272 | -1.197150  | -1.903308  |
| 6215 | 0.483335  | -0.126424 | -0.065806  | -0.566612  |
| 6884 | 0.051478  | -0.663507 | -0.388948  | 0.368279   |
| 3647 | -1.122266 | -1.012620 | -2.018156  | -1.023936  |

- Looking at the predictive ability of the remaining features. In order to do that principal component analysis will be utilised with visualizing the components on a graph against the target.

```
[42]: # Taking all of the features to assess how much of variance is influenced by
      ↪ each feature.
pca = PCA(n_components = len(coef_dict))
```

```
[39]: data_pca = pca.fit_transform(x_tr_scaled)
      data_pca.shape
```

```
[39]: (6300, 10)
```

```
[40]: pca.explained_variance_ratio_
```

```
[40]: array([0.28700352, 0.17523067, 0.12405576, 0.10048613, 0.09608016,
           0.08450172, 0.06252189, 0.0355892 , 0.02251423, 0.01201671])
```

```
[43]: # The first two features could be left as the most influential ones, but the
      ↪ first three components will also be checked
pca_new = PCA(n_components=2)
data_new = pca_new.fit_transform(x_tr_scaled)
```

```
[44]: data_new
```

```
[44]: array([[ -0.00776042, -0.52236894],
           [-1.28566659,  1.35763105],
           [ 0.6210337 , -0.2135601 ],
           ...,
           [-1.79237927,  1.23449541],
           [ 1.5703436 , -0.2720664 ],
           [ 0.35846336, -1.79997616]])
```

```
[45]: pca_new.explained_variance_ratio_
```

```
[45]: array([0.28700352, 0.17523067])
```

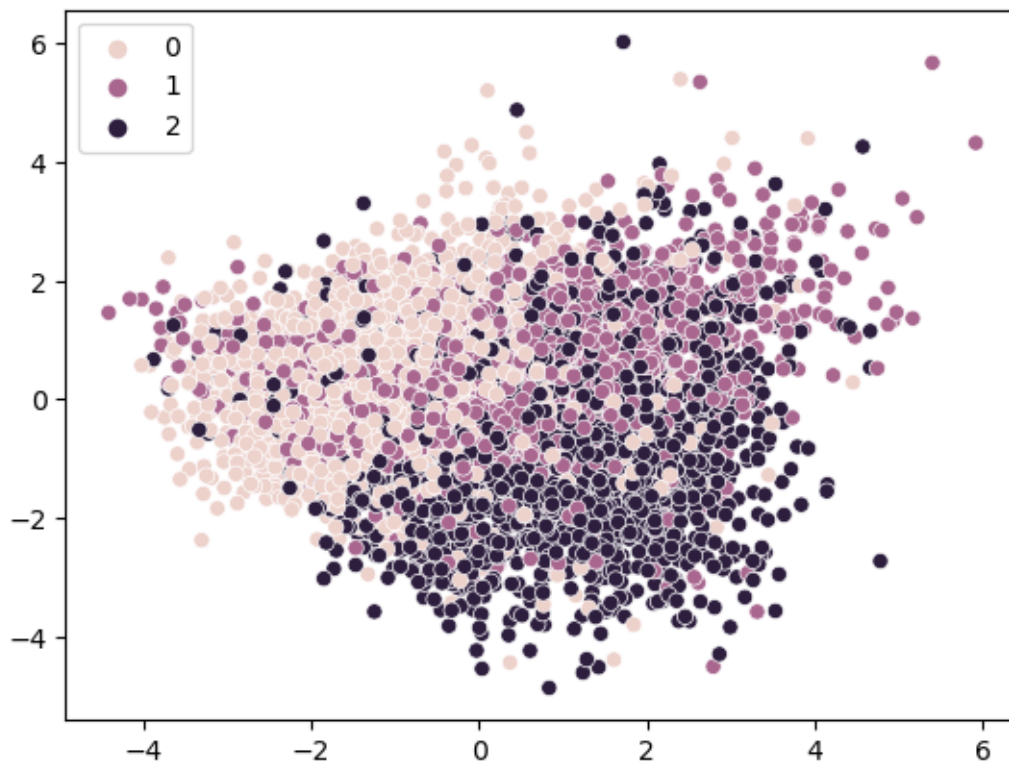
```
[46]: data_new
```

```
[46]: array([[ -0.00776042, -0.52236894],
             [-1.28566659,  1.35763105],
             [ 0.6210337 , -0.2135601 ],
             ...,
             [-1.79237927,  1.23449541],
             [ 1.5703436 , -0.2720664 ],
             [ 0.35846336, -1.79997616]])
```

```
[47]: data_new[:,0]
```

```
[47]: array([-0.00776042, -1.28566659,  0.6210337 , ..., -1.79237927,
            1.5703436 ,  0.35846336])
```

```
[49]: sns.scatterplot(x=data_new[:,0],y=data_new[:,1],hue=y_class_tr);
```



```
[50]: # Additional PCA visualization with the first three components
pca_newer = PCA(n_components=3)
data_newer = pca_newer.fit_transform(x_tr_scaled)
```

```
[51]: trace_pca0 = go.Scatter3d(
    x = data_newer[y_class_tr==0,0],
    y = data_newer[y_class_tr==0,1],
```

```

        z = data_newer[y_class_tr==0,2],
        mode='markers',
        marker = dict(
            size = 5,
            color = 'red',
            opacity = 0.8
        ),
        name='0'
    )
    trace_pca1 = go.Scatter3d(
        x = data_newer[y_class_tr==1,0],
        y = data_newer[y_class_tr==1,1],
        z = data_newer[y_class_tr==1,2],
        mode='markers',
        marker = dict(
            size = 5,
            color = 'green',
            opacity = 0.8
        ),
        name='1'
    )
    trace_pca2 = go.Scatter3d(
        x = data_newer[y_class_tr==2,0],
        y = data_newer[y_class_tr==2,1],
        z = data_newer[y_class_tr==2,2],
        mode='markers',
        marker = dict(
            size = 5,
            color = 'blue',
            opacity = 0.8
        ),
        name='2'
    )

```

```

[52]: layout = go.Layout(
    scene = dict(
        xaxis = dict(title='PC1'),
        yaxis = dict(title='PC2'),
        zaxis = dict(title='PC3'),
    ),
    margin = dict(l=0,r=0,b=0,t=0),
    legend=dict(x=0,y=1)
)

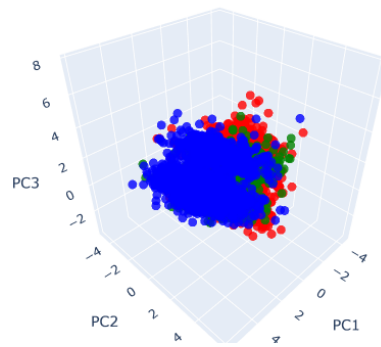
```

```

[53]: fig = go.Figure(data=[trace_pca0,trace_pca1,trace_pca2], layout=layout)
fig.show()

```

• 0  
• 1  
• 2



```
[54]: y_class_tr
```

```
[54]: 7581    1
      8484    0
      6215    1
      6884    1
      3647    0
      ..
      5734    2
      5191    2
      5390    1
      860     1
      7270    1
      Length: 6300, dtype: int64
```

```
[55]: # Preparing data for further analysis
      data_new = pd.DataFrame(data_new, columns=['new1', 'new2'])
      data_new
```

```
[55]:
```

|      | new1      | new2      |
|------|-----------|-----------|
| 0    | -0.007760 | -0.522369 |
| 1    | -1.285667 | 1.357631  |
| 2    | 0.621034  | -0.213560 |
| 3    | 1.625722  | 1.655920  |
| 4    | -2.012947 | 0.609353  |
| ...  | ...       | ...       |
| 6295 | -0.444403 | -2.869545 |
| 6296 | 1.763569  | -0.188656 |
| 6297 | -1.792379 | 1.234495  |
| 6298 | 1.570344  | -0.272066 |
| 6299 | 0.358463  | -1.799976 |

```
[6300 rows x 2 columns]
```



```
[56]: data_new.index = x_tr_scaled.index
```

```
[57]: data_newer = pd.DataFrame(data_newer, columns=['newer1', 'newer2', 'newer_3'])
data_newer.tail(5)
```

```
[57]:      newer1    newer2    newer_3
6295 -0.444403 -2.869545  1.014904
6296  1.763569 -0.188656 -0.132870
6297 -1.792379  1.234495 -2.255040
6298  1.570344 -0.272066 -0.126920
6299  0.358463 -1.799976 -0.661736
```

```
[58]: data_newer.index = x_tr_scaled.index
```

14. Now that the main feature are chosen, a more complex model could be built. For that models based on trees will be utilised. The models will be assessed based on the accuracy metric.

```
[59]: # Trying two types of models with trees
tree_new = DecisionTreeClassifier(random_state = 42)
forest_new = RandomForestClassifier(random_state=42)
```

```
[60]: params_tree = {'max_depth': np.arange(3,7),
                    'max_features': [1,2,3,4]}
search_tree = GridSearchCV(estimator=tree_new, param_grid=params_tree, scoring='accuracy', cv=3)
```

```
[61]: # Looking at the initial dataset (prior to the PCA)
search_tree.fit(x_tr_scaled, y_class_tr)
```

```
[61]: GridSearchCV(cv=3, estimator=DecisionTreeClassifier(random_state=42),
                param_grid={'max_depth': array([3, 4, 5, 6]),
                            'max_features': [1, 2, 3, 4]},
                scoring='accuracy')
```

```
[62]: y_class_pred = search_tree.predict(x_t_scaled)
accuracy_score(y_class_t, y_class_pred)
```

```
[62]: 0.7307407407407407
```

```
[63]: # Now trying the data after the PCA
search_tree.fit(data_new, y_class_tr)
```

```
[63]: GridSearchCV(cv=3, estimator=DecisionTreeClassifier(random_state=42),
                param_grid={'max_depth': array([3, 4, 5, 6]),
                            'max_features': [1, 2, 3, 4]},
                scoring='accuracy')
```

```
[64]: data_new_t = pca_new.transform(x_t_scaled)
data_new_t = pd.DataFrame(data_new_t,columns=['new1','new2'])
data_new_t.index = x_t_scaled.index
```

```
[65]: y_class_pred = search_tree.predict(data_new_t)
accuracy_score(y_class_t,y_class_pred)
```

```
[65]: 0.617037037037037
```

```
[66]: search_tree.fit(data_newer,y_class_tr)
data_newer_t = pca_newer.transform(x_t_scaled)
data_newer_t = pd.DataFrame(data_newer_t,columns=['newer1','newer2','newer_3'])
data_newer_t.index = x_t_scaled.index
y_class_pred = search_tree.predict(data_newer_t)
accuracy_score(y_class_t,y_class_pred)
```

```
[66]: 0.6381481481481481
```

Obviously, the model with a greater number of features gave the higher accuracy of the predictions since the majority of the features have comparable predictive abilities and taking away a large part of the features is not correct. Consecutively, the data that will be used further is the data prior to the PCA so that the results are more accurate.

```
[67]: params_forest = {'n_estimators' : [100,200],
                      'max_depth': np.arange(3,7),
                      'max_features': [1,2,3,4]}
search_forest = GridSearchCV(estimator=forest_new,param_grid=params_forest,scoring='accuracy',cv=3)
```

```
[68]: search_forest.fit(x_tr_scaled,y_class_tr)
```

```
[68]: GridSearchCV(cv=3, estimator=RandomForestClassifier(random_state=42),
                  param_grid={'max_depth': array([3, 4, 5, 6]),
                              'max_features': [1, 2, 3, 4],
                              'n_estimators': [100, 200]},
                  scoring='accuracy')
```

```
[69]: y_class_pred = search_forest.predict(x_t_scaled)
accuracy_score(y_class_t,y_class_pred)
```

```
[69]: 0.7685185185185185
```

15. Trying using other types of model. This time it will be a SVM model and its accuracy will be evaluated.

```
[70]: svc = SVC(kernel='poly',probability=True)
```

```
[71]: params_svc = {'degree': range(2,7), 'C': [0.1,1,10]}
```

```
svc_search =  
↳GridSearchCV(estimator=svc,param_grid=params_svc,scoring='accuracy',cv=3)
```

```
[72]: svc_search.fit(x_tr_scaled,y_class_tr)
```

```
[72]: GridSearchCV(cv=3, estimator=SVC(kernel='poly', probability=True),  
    param_grid={'C': [0.1, 1, 10], 'degree': range(2, 7)},  
    scoring='accuracy')
```

```
[73]: y_class_pred = svc_search.predict(x_t_scaled)  
accuracy_score(y_class_t,y_class_pred)
```

```
[73]: 0.81
```

16. Trying using other types of model. This time it will be a KNN model and its accuracy will be evaluated.

```
[74]: knn = KNeighborsClassifier()  
params = {'n_neighbors':range(2,21)}  
knn_search = GridSearchCV(estimator=knn, param_grid=params, scoring='accuracy',  
    ↳cv=3)  
knn_search.fit(x_tr_scaled,y_class_tr)
```

```
[74]: GridSearchCV(cv=3, estimator=KNeighborsClassifier(),  
    param_grid={'n_neighbors': range(2, 21)}, scoring='accuracy')
```

```
[75]: y_class_pred = knn_search.predict(x_t_scaled)  
accuracy_score(y_class_t,y_class_pred)
```

```
[75]: 0.8107407407407408
```

17. The last two model have showed good results in terms of accurate, but the random forest model has also showed a decent score. Next, boosting could be used where appropriate.

```
[76]: ada_boost0 = AdaBoostClassifier(n_estimators=100,learning_rate=0.  
    ↳1,random_state=42)  
ada_boost0.fit(x_tr_scaled,y_class_tr)  
y_pred = ada_boost0.predict(x_t_scaled)  
accuracy_score(y_class_t,y_pred)
```

```
[76]: 0.7025925925925925
```

```
[77]: ada_boost = AdaBoostClassifier(estimator=search_forest.  
    ↳best_estimator_,n_estimators=100,learning_rate=0.1,random_state=42)
```

```
[78]: ada_boost.fit(x_tr_scaled,y_class_tr)
```

```
[78]: AdaBoostClassifier(estimator=RandomForestClassifier(max_depth=6, max_features=4,  
    n_estimators=200,
```

```
random_state=42),  
learning_rate=0.1, n_estimators=100, random_state=42)
```

```
[79]: y_pred = ada_boost.predict(x_t_scaled)  
accuracy_score(y_class_t,y_pred)
```

```
[79]: 0.8225925925925925
```

```
[80]: gbc = GradientBoostingClassifier(init=search_forest.  
    ↳best_estimator_,n_estimators=100,learning_rate=0.  
    ↳1,max_depth=5,random_state=42)  
gbc.fit(x_tr_scaled, y_class_tr)  
prediction = gbc.predict(x_t_scaled)  
accuracy_score(y_class_t,prediction)
```

```
[80]: 0.8125925925925926
```

```
[81]: gbc2 = GradientBoostingClassifier(init=knn_search.  
    ↳best_estimator_,n_estimators=100,learning_rate=0.  
    ↳1,max_depth=5,random_state=42)  
gbc2.fit(x_tr_scaled, y_class_tr)  
prediction = gbc2.predict(x_t_scaled)  
accuracy_score(y_class_t,prediction)
```

```
[81]: 0.8103703703703704
```

```
[82]: gbc3 = GradientBoostingClassifier(init=svc_search.  
    ↳best_estimator_,n_estimators=100,learning_rate=0.  
    ↳1,max_depth=5,random_state=42)  
gbc3.fit(x_tr_scaled, y_class_tr)  
prediction = gbc3.predict(x_t_scaled)  
accuracy_score(y_class_t,prediction)
```

```
[82]: 0.812962962962963
```

```
[83]: gbc4 = GradientBoostingClassifier(n_estimators=100,learning_rate=0.  
    ↳1,max_depth=5,random_state=42)  
gbc4.fit(x_tr_scaled, y_class_tr)  
prediction = gbc4.predict(x_t_scaled)  
accuracy_score(y_class_t,prediction)
```

```
[83]: 0.8137037037037037
```

```
[84]: xgb = xgb.XGBClassifier(n_estimators=100,learning_rate=0.  
    ↳1,max_depth=5,reg_lambda=1,random_state=42)  
xgb.fit(x_tr_scaled, y_class_tr)  
prediction = xgb.predict(x_t_scaled)  
accuracy_score(y_class_t,prediction)
```

[84]: 0.8107407407407408

```
[85]: lgbm = LGBMClassifier(n_estimators=200, learning_rate=0.
      ↪ 01, reg_lambda=1, max_depth=5, random_state=42)
      lgbm.fit(x_tr_scaled, y_class_tr)
      prediction = lgbm.predict(x_t_scaled)
```

[illegible]

[illegible]

[illegible]

[illegible]



[illegible]

[illegible]

[illegible]

[illegible]

[illegible]

[illegible]

[illegible]

[illegible]



[illegible]

[illegible]



[86] : 0.7907407407407407

**Conclusion** As can be seen the best model overall was one of the AdaBoost models that used a random forest as one of its arguments. Achieved accuracy is slightly above 82%.