

Anotações:

Verificar se os dados estão balanceados, pois se trata de um problema de classificação

Os dados estão todos codificados, então não há porque fazer uma análise exploratória

Existem muitas variáveis, considerar usar redução de dimensionalidade

Testar os algoritmos com validação cruzada e construir boxplots com os resultados

Testar se feature selection melhora o modelo

1 = insatisfeito, 0 = satisfeito

```
In [28]: import pandas as pd
import numpy as np
from sklearn.preprocessing import MinMaxScaler
from sklearn.model_selection import cross_val_score, KFold
from sklearn.metrics import confusion_matrix, classification_report
from sklearn.naive_bayes import GaussianNB
from sklearn.svm import SVC
from sklearn.metrics import KNeighborsClassifier
from sklearn.linear_model import LogisticRegression
from sklearn.discriminant_analysis import LinearDiscriminantAnalysis
from sklearn.metrics import confusion_matrix
from sklearn.metrics import confusion_matrix
from sklearn.metrics import confusion_matrix

In [71]: df = pd.read_csv('train.csv')

In [41]: df.shape

Out[41]: (10, 10)

In [51]: df.columns

Out[51]: ['ID', 'var1', 'var2', 'var3', 'var4', 'var5', 'var6', 'var7', 'var8', 'var9', 'var10', 'var11', 'var12', 'var13', 'var14', 'var15', 'var16', 'var17', 'var18', 'var19', 'var20', 'var21', 'var22', 'var23', 'var24', 'var25', 'var26', 'var27', 'var28', 'var29', 'var30', 'var31', 'var32', 'var33', 'var34', 'var35', 'var36', 'var37', 'var38', 'var39', 'var40', 'var41', 'var42', 'var43', 'var44', 'var45', 'var46', 'var47', 'var48', 'var49', 'var50', 'var51', 'var52', 'var53', 'var54', 'var55', 'var56', 'var57', 'var58', 'var59', 'var60', 'var61', 'var62', 'var63', 'var64', 'var65', 'var66', 'var67', 'var68', 'var69', 'var70', 'var71', 'var72', 'var73', 'var74', 'var75', 'var76', 'var77', 'var78', 'var79', 'var80', 'var81', 'var82', 'var83', 'var84', 'var85', 'var86', 'var87', 'var88', 'var89', 'var90', 'var91', 'var92', 'var93', 'var94', 'var95', 'var96', 'var97', 'var98', 'var99', 'var100']

In [53]: df.isnull().sum()

Out[53]: 0

In [55]: df.dtypes

Out[55]: ID: int64
var1: float64
var2: float64
var3: float64
var4: float64
var5: float64
var6: float64
var7: float64
var8: float64
var9: float64
var10: float64
var11: float64
var12: float64
var13: float64
var14: float64
var15: float64
var16: float64
var17: float64
var18: float64
var19: float64
var20: float64
var21: float64
var22: float64
var23: float64
var24: float64
var25: float64
var26: float64
var27: float64
var28: float64
var29: float64
var30: float64
var31: float64
var32: float64
var33: float64
var34: float64
var35: float64
var36: float64
var37: float64
var38: float64
var39: float64
var40: float64
var41: float64
var42: float64
var43: float64
var44: float64
var45: float64
var46: float64
var47: float64
var48: float64
var49: float64
var50: float64
var51: float64
var52: float64
var53: float64
var54: float64
var55: float64
var56: float64
var57: float64
var58: float64
var59: float64
var60: float64
var61: float64
var62: float64
var63: float64
var64: float64
var65: float64
var66: float64
var67: float64
var68: float64
var69: float64
var70: float64
var71: float64
var72: float64
var73: float64
var74: float64
var75: float64
var76: float64
var77: float64
var78: float64
var79: float64
var80: float64
var81: float64
var82: float64
var83: float64
var84: float64
var85: float64
var86: float64
var87: float64
var88: float64
var89: float64
var90: float64
var91: float64
var92: float64
var93: float64
var94: float64
var95: float64
var96: float64
var97: float64
var98: float64
var99: float64
var100: float64

In [67]: df.dtypes

Out[67]: ID: int64
var1: float64
var2: float64
var3: float64
var4: float64
var5: float64
var6: float64
var7: float64
var8: float64
var9: float64
var10: float64
var11: float64
var12: float64
var13: float64
var14: float64
var15: float64
var16: float64
var17: float64
var18: float64
var19: float64
var20: float64
var21: float64
var22: float64
var23: float64
var24: float64
var25: float64
var26: float64
var27: float64
var28: float64
var29: float64
var30: float64
var31: float64
var32: float64
var33: float64
var34: float64
var35: float64
var36: float64
var37: float64
var38: float64
var39: float64
var40: float64
var41: float64
var42: float64
var43: float64
var44: float64
var45: float64
var46: float64
var47: float64
var48: float64
var49: float64
var50: float64
var51: float64
var52: float64
var53: float64
var54: float64
var55: float64
var56: float64
var57: float64
var58: float64
var59: float64
var60: float64
var61: float64
var62: float64
var63: float64
var64: float64
var65: float64
var66: float64
var67: float64
var68: float64
var69: float64
var70: float64
var71: float64
var72: float64
var73: float64
var74: float64
var75: float64
var76: float64
var77: float64
var78: float64
var79: float64
var80: float64
var81: float64
var82: float64
var83: float64
var84: float64
var85: float64
var86: float64
var87: float64
var88: float64
var89: float64
var90: float64
var91: float64
var92: float64
var93: float64
var94: float64
var95: float64
var96: float64
var97: float64
var98: float64
var99: float64
var100: float64

In [72]: df = df.iloc[:,1:]

In [73]: df.dtypes

Out[73]: ID: int64
var1: float64
var2: float64
var3: float64
var4: float64
var5: float64
var6: float64
var7: float64
var8: float64
var9: float64
var10: float64
var11: float64
var12: float64
var13: float64
var14: float64
var15: float64
var16: float64
var17: float64
var18: float64
var19: float64
var20: float64
var21: float64
var22: float64
var23: float64
var24: float64
var25: float64
var26: float64
var27: float64
var28: float64
var29: float64
var30: float64
var31: float64
var32: float64
var33: float64
var34: float64
var35: float64
var36: float64
var37: float64
var38: float64
var39: float64
var40: float64
var41: float64
var42: float64
var43: float64
var44: float64
var45: float64
var46: float64
var47: float64
var48: float64
var49: float64
var50: float64
var51: float64
var52: float64
var53: float64
var54: float64
var55: float64
var56: float64
var57: float64
var58: float64
var59: float64
var60: float64
var61: float64
var62: float64
var63: float64
var64: float64
var65: float64
var66: float64
var67: float64
var68: float64
var69: float64
var70: float64
var71: float64
var72: float64
var73: float64
var74: float64
var75: float64
var76: float64
var77: float64
var78: float64
var79: float64
var80: float64
var81: float64
var82: float64
var83: float64
var84: float64
var85: float64
var86: float64
var87: float64
var88: float64
var89: float64
var90: float64
var91: float64
var92: float64
var93: float64
var94: float64
var95: float64
var96: float64
var97: float64
var98: float64
var99: float64
var100: float64

In [75]: df.drop([list(unicos.loc[unicos[0]==1].index),axis=1,inplace=True]
```

Algumas colunas só tem um valor único, e por isso podem ser descartadas



Out [76]: unicos.sort\_values(ascending=False)

|                               |       |
|-------------------------------|-------|
|                               | 0     |
| var38                         | 57736 |
| saldo_medio_var5_uit          | 17330 |
| saldo_var30                   | 16940 |
| saldo_var42                   | 15730 |
| saldo_medio_var5_uit1         | 14778 |
| saldo_medio_var5_hace2        | 14486 |
| saldo_var5                    | 11642 |
| imp_op_var39_comer_uit3       | 9099  |
| imp_op_var41_comer_uit3       | 8961  |
| imp_op_var39_uit1             | 8149  |
| imp_op_var41_uit1             | 8032  |
| saldo_medio_var5_hace3        | 7787  |
| imp_op_var39_comer_uit1       | 7551  |
| imp_op_var41_comer_uit1       | 7421  |
| saldo_var37                   | 4041  |
| imp_trans_var37_uit1          | 3631  |
| saldo_medio_var12_uit3        | 3447  |
| saldo_medio_var12_uit1        | 3405  |
| saldo_var12                   | 3059  |
| saldo_medio_var12_hace2       | 2834  |
| saldo_var24                   | 2614  |
| saldo_medio_var13_corto_uit3  | 2576  |
| imp_var43_emit_uit1           | 2342  |
| saldo_medio_var6_uit3         | 2173  |
| saldo_medio_var6_uit1         | 2116  |
| saldo_var8                    | 1989  |
| saldo_medio_var13_corto_hace2 | 1628  |
| saldo_var26                   | 1592  |
| saldo_var25                   | 1524  |
| saldo_medio_var2_hace2        | 1325  |
| saldo_medio_var12_hace3       | 1152  |
| saldo_medio_var13_corto_hace3 | 968   |
| saldo_medio_var13_corto_uit1  | 943   |
| saldo_medio_var13_corto_uit3  | 859   |
| saldo_var13                   | 730   |
| imp_ent_var13_uit1            | 596   |
| saldo_medio_var13_largo_uit3  | 511   |
| imp_op_var39_effect_uit3      | 462   |
| imp_op_var41_effect_uit3      | 454   |
| saldo_medio_var8_hace3        | 439   |
| imp_aport_var13_hace3         | 425   |
| imp_op_var40_comer_uit3       | 346   |
| saldo_var14                   | 345   |
| imp_op_var39_effect_uit1      | 336   |
| imp_op_var41_effect_uit1      | 331   |
| saldo_medio_var13_largo_hace2 | 295   |
| imp_op_var40_comer_uit1       | 293   |
| saldo_var1                    | 285   |
| saldo_var40                   | 282   |
| saldo_var31                   | 278   |
| saldo_medio_var13_largo_hace3 | 264   |
| saldo_medio_var13_largo_uit1  | 233   |
| saldo_var13_largo             | 229   |
| imp_op_var40_uit1             | 224   |
| var3                          | 208   |
| imp_var7_recib_uit1           | 184   |
| imp_aport_var13_uit1          | 182   |
| saldo_var20                   | 176   |
| num_var45_uit3                | 172   |
| saldo_medio_var44_uit1        | 141   |
| saldo_medio_var44_uit3        | 141   |
| saldo_var44                   | 129   |
| saldo_medio_var17_uit3        | 119   |
| saldo_medio_var17_uit1        | 119   |
| saldo_var17                   | 111   |
| var15                         | 100   |
| saldo_medio_var44_hace2       | 99    |
| num_op_var39_uit3             | 99    |
| num_op_var41_uit3             | 96    |
| num_var45_uit1                | 94    |
| num_op_var39_comer_uit3       | 92    |
| saldo_medio_var17_hace2       | 88    |
| num_op_var41_comer_uit3       | 88    |
| num_compra_var45_hace2        | 85    |
| imp_compra_var44_uit1         | 85    |
| saldo_var32                   | 75    |
| num_op_var39_uit1             | 71    |
| num_med_var45_uit3            | 71    |
| num_op_var41_uit1             | 68    |
| imp_sal_var16_uit1            | 66    |
| num_var45_hace3               | 66    |
| num_op_var39_comer_uit1       | 63    |
| num_op_var41_comer_uit1       | 60    |
| num_op_var41_hace2            | 51    |
| num_op_var39_hace2            | 50    |
| saldo_medio_var33_uit3        | 48    |
| saldo_medio_var33_uit1        | 48    |
| saldo_var33                   | 48    |
| num_op_var40_comer_uit3       | 47    |
| imp_venta_var44_uit1          | 45    |
| saldo_medio_var33_hace2       | 43    |
| imp_aport_var17_uit1          | 42    |
| num_op_var41_effect_uit3      | 40    |
| num_op_var39_effect_uit3      | 40    |
| num_var43_recib_uit1          | 38    |
| num_op_var40_uit3             | 35    |
| num_op_var40_comer_uit1       | 34    |
| num_var22_uit3                | 33    |
| imp_compra_var44_hace3        | 33    |
| imp_reemb_var13_uit1          | 33    |
| saldo_medio_var44_hace3       | 33    |
| num_op_var40_uit1             | 30    |
| imp_op_var40_effect_uit3      | 29    |
| delta_imp_aport_var13_ty3     | 27    |
| num_var43_emit_uit1           | 27    |
| num_op_var39_effect_uit1      | 26    |
| num_op_var41_effect_uit1      | 25    |
| saldo_medio_var33_hace3       | 24    |
| imp_op_var40_effect_uit1      | 23    |
| num_var37_0                   | 22    |
| num_var37                     | 22    |
| num_op_var39_hace2            | 22    |
| num_var22_hace2               | 22    |
| num_op_var41_hace3            | 22    |
| num_var37_med_uit2            | 21    |
| imp_aport_var17_hace3         | 20    |
| imp_reemb_var17_uit1          | 20    |
| num_var22_hace3               | 19    |
| num_op_var40_hace2            | 19    |
| num_var22_uit1                | 18    |
| num_trasp_var11_uit1          | 18    |
| saldo_medio_var17_hace3       | 18    |
| delta_imp_compra_var44_ty3    | 17    |
| imp_aport_var33_hace3         | 16    |
| num_med_var22_uit3            | 15    |
| num_var35                     | 13    |
| num_ent_var16_uit1            | 13    |
| num_var17_0                   | 11    |
| num_var31_0                   | 11    |
| num_var30_0                   | 11    |
| num_compra_var44_uit1         | 10    |
| num_venta_var44_uit1          | 9     |
| delta_imp_aport_var33_ty3     | 9     |
| num_var30                     | 9     |
| delta_num_compra_var44_ty3    | 9     |
| num_var25                     | 9     |
| num_var39_0                   | 9     |
| num_var41_0                   | 9     |
| num_var17                     | 9     |
| num_op_var40_effect_uit3      | 9     |
| num_var26_0                   | 9     |
| num_var28_0                   | 9     |
| num_var42_0                   | 8     |
| num_var4                      | 8     |
| num_med_var13_hace3           | 8     |
| num_aport_var17_uit1          | 8     |
| imp_trasp_var33_in_hace3      | 7     |
| num_med_var33_uit1            | 7     |
| num_op_var40_effect_uit1      | 7     |
| num_var13                     | 7     |
| num_var13_largo               | 7     |
| imp_var37_largo_0             | 7     |
| num_var13_largo_0             | 7     |
| delta_imp_aport_var17_ty3     | 7     |
| imp_trasp_var33_uit1          | 6     |
| num_var7_recib_uit1           | 6     |
| num_sal_var16_uit1            | 6     |
| delta_num_aport_var17_ty3     | 6     |
| delta_num_aport_var13_ty3     | 6     |
| imp_aport_var33_uit1          | 6     |
| num_var12_0                   | 6     |
| num_reemb_var17_uit1          | 5     |
| num_var32_0                   | 5     |
| delta_imp_venta_var44_ty3     | 5     |
| num_var32                     | 5     |
| saldo_medio_var29_hace2       | 5     |
| imp_trasp_var17_out_uit1      | 5     |
| imp_trasp_var17_in_uit1       | 5     |
| num_aport_var17_hace3         | 5     |
| num_var14_0                   | 5     |
| num_var5_0                    | 5     |
| var36                         | 5     |
| num_op_var40_hace3            | 5     |
| delta_num_venta_var44_ty3     | 5     |
| num_var5                      | 5     |
| num_var33_0                   | 4     |
| saldo_var12                   | 4     |
| num_var7_emit_uit1            | 4     |
| imp_var7_emit_uit3            | 4     |
| num_aport_var33_hace3         | 4     |
| imp_venta_var44_hace3         | 4     |
| num_var14                     | 4     |
| imp_var24_0                   | 4     |
| num_meses_var44_uit3          | 4     |
| num_meses_var33_uit3          | 4     |
| num_compra_var44_hace3        | 4     |
| num_meses_var17_uit3          | 4     |
| num_meses_var13_largo_uit3    | 4     |
| num_meses_var12_uit3          | 4     |
| num_meses_var33_ty3           | 4     |
| num_meses_var5_uit3           | 4     |
| saldo_medio_var29_uit3        | 4     |
| saldo_medio_var28_uit3        | 4     |
| saldo_medio_var28_uit1        | 4     |
| num_var13_corto               | 3     |
| num_trasp_var33_in_uit1       | 3     |
| saldo_var18                   | 3     |
| imp_amort_var18_uit1          | 3     |
| num_venta_var44_hace3         | 3     |
| num_var40_0                   | 3     |
| num_var24                     | 3     |
| num_var3                      | 3     |
| delta_num_trasp_var33_ty3     | 3     |
| num_trasp_var17_in_hace3      | 3     |
| num_var13_corto_0             | 3     |
| imp_trasp_var17_in_hace3      | 3     |
| saldo_medio_var13_medio_hace2 | 3     |
| saldo_var34                   | 3     |
| saldo_medio_var13_medio_uit1  | 3     |
| saldo_medio_var13_medio_uit3  | 3     |
| saldo_var6                    | 3     |
| imp_amort_var34_uit1          | 3     |
| num_aport_var33_uit1          | 3     |
| num_var44_0                   | 3     |
| delta_imp_reemb_var17_ty3     | 3     |
| num_var1                      | 3     |
| delta_imp_trasp_var17_in_ty3  | 3     |
| num_var1_0                    | 3     |
| delta_imp_trasp_var33_in_ty3  | 3     |
| delta_num_trasp_var17_ty3     | 3     |
| saldo_var13_medio             | 3     |
| saldo_var29                   | 3     |
| num_var8_0                    | 3     |
| delta_num_reemb_var17_ty3     | 3     |
| num_meses_var29_uit3          | 3     |
| ind_var1                      | 2     |
| num_var6_0                    | 2     |
| ind_var1_0                    | 2     |
| num_var8                      | 2     |
| num_reemb_var13_uit1          | 2     |
| ind_var5                      | 2     |
| ind_var6_0                    | 2     |
| num_var18_0                   | 2     |
| num_meses_var13_medio_uit3    | 2     |
| num_var16                     | 2     |
| num_reemb_var17_hace3         | 2     |
| num_var6                      | 2     |
| ind_var25_cle                 | 2     |
| ind_var26                     | 2     |
| num_var20_0                   | 2     |
| ind_var26_cle                 | 2     |
| num_var20                     | 2     |
| saldo_medio_var29_hace3       | 2     |
| ind_var26_0                   | 2     |
| ind_var13_largo_0             | 2     |
| ind_var6_0                    | 2     |
| ind_var13_0                   | 2     |
| ind_var13_medio_0             | 2     |
| ind_var13_medio               | 2     |
| ind_var13                     | 2     |
| ind_var14_0                   | 2     |
| ind_var14                     | 2     |
| ind_var13_corto               | 2     |
| ind_var17_0                   | 2     |
| ind_var13_corto_0             | 2     |
| ind_var17                     | 2     |
| ind_var18_0                   | 2     |
| ind_var18                     | 2     |
| ind_var19                     | 2     |
| ind_var20_0                   | 2     |
| ind_var20                     | 2     |
| num_reemb_var33_uit1          | 2     |
| ind_var12                     | 2     |
| ind_var24                     | 2     |
| ind_var12_0                   | 2     |
| num_trasp_var33_out_uit1      | 2     |
| ind_var8_0                    | 2     |
| ind_var3                      | 2     |
| num_var33_in_hace3            | 2     |
| num_trasp_var17_out_uit1      | 2     |
| num_trasp_var17_in_uit1       | 2     |
| ind_var6                      | 2     |
| ind_var13_largo               | 2     |
| num_var13_medio_0             | 2     |
| num_var13_medio               | 2     |
| ind_var24_0                   | 2     |
| TARGET                        | 2     |
| ind_var25_0                   | 2     |
| delta_imp_reemb_var33_ty3     | 2     |
| ind_var29                     | 2     |
| delta_imp_amort_var18_ty3     | 2     |
| delta_imp_reemb_var34_ty3     | 2     |
| ind_var39                     | 2     |
| num_var44                     | 2     |
| delta_imp_reemb_var13_ty3     | 2     |
| delta_imp_trasp_var17_out_ty3 | 2     |
| ind_var30                     | 2     |
| delta_imp_trasp_var33_out_ty3 | 2     |
| num_var39                     | 2     |
| delta_num_reemb_var13_ty3     | 2     |
| delta_num_reemb_var33_ty3     | 2     |
| delta_num_trasp_var17_out_ty3 | 2     |
| ind_var29_0                   | 2     |
| ind_var30_0                   | 2     |
| num_var31_0                   | 2     |
| num_var34                     | 2     |
| ind_var37_cle                 | 2     |
| ind_var41_0                   | 2     |
| ind_var40                     | 2     |
| ind_var40_0                   | 2     |
| ind_var39_0                   | 2     |
| ind_var37_0                   | 2     |
| ind_var34                     | 2     |
| ind_var31                     | 2     |
| ind_var34_0                   | 2     |
| ind_var33                     | 2     |
| ind_var32_0                   | 2     |
| ind_var32                     | 2     |
| ind_var32_cle                 | 2     |
| num_var40                     | 2     |
| delta_num_trasp_var33_out_ty3 | 2     |
| ind_var7_emit_uit1            | 2     |
| imp_reemb_var33_uit1          | 2     |
| ind_var7_recib_uit1           | 2     |
| ind_var8_cle_uit1             | 2     |
| ind_var43_emit_uit1           | 2     |
| num_var29                     | 2     |
| imp_trasp_var32_out_uit1      | 2     |
| ind_var43_recib_uit1          | 2     |
| num_var29_0                   | 2     |
| ind_var25                     | 2     |
| ind_var10cte_uit1             | 2     |
| ind_var16_uit1                | 2     |
| num_var7_emit_uit1            | 2     |
| imp_reemb_var17_hace3         | 2     |
| num_var34_0                   | 2     |
| ind_var44                     | 2     |
| num_var44_0                   | 2     |
| num_var27                     | 1     |
| saldo_medio_var13_medio_hace3 | 1     |
| ind_var46_0                   | 1     |
| num_var27_0                   | 1     |
| ind_var2_0                    | 1     |
| num_var28                     | 1     |
| saldo_var27                   | 1     |
| saldo_var28                   | 1     |
| ind_var2                      | 1     |
| num_var2_0_uit1               | 1     |
| num_var2_uit1                 | 1     |
| ind_var28                     | 1     |
| saldo_var41                   | 1     |
| imp_reemb_var33_hace3         | 1     |
| ind_var28_0                   | 1     |
| ind_var27                     | 1     |
| num_trasp_var17_out_hace3     | 1     |
| imp_amort_var34_hace3         | 1     |
| imp_amort_var18_hace3         | 1     |
| ind_var46                     | 1     |
| saldo_var2_uit1               | 1     |
| imp_reemb_var13_hace3         | 1     |
| num_var41                     | 1     |
| num_trasp_var33_out_hace3     | 1     |
| ind_var27_0                   | 1     |
| imp_trasp_var17_out_hace3     | 1     |
| num_reemb_var33_hace3         | 1     |
| imp_trasp_var33_out_hace3     | 1     |
| num_var46_0                   | 1     |
| num_var46                     | 1     |
| num_reemb_var13_hace3         | 1     |
| saldo_var46                   | 1     |
| ind_var41                     | 1     |

A coluna var38 tem 57736 valores únicos, sendo que há 76020 linhas. Avaliar se preciso apagar ela

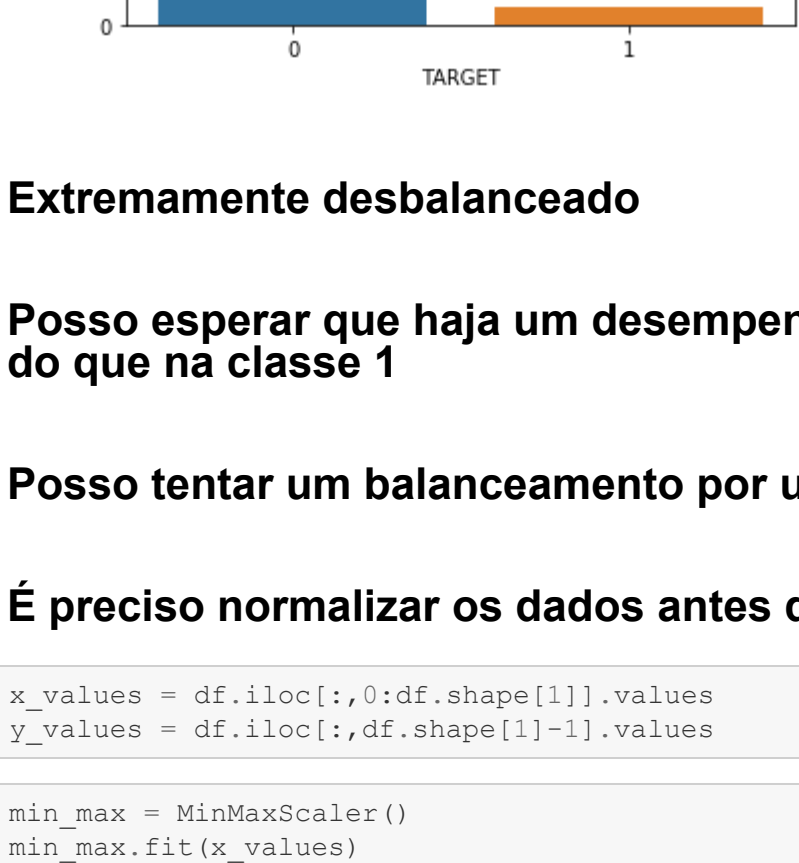
In [77]: df.duplicated().sum()

Out [77]: 4807

Removendo duplicatas

In [78]: df.drop\_duplicates(inplace=True)

In [13]: sns.countplot(data=df,x='TARGET');



Extremamente desbalanceado

Posso esperar que haja um desempenho muito melhor na identificação da classe 0 do que na classe 1

Posso tentar um balanceamento por undersampling (Tomek links)

É preciso normalizar os dados antes de aplicar PCA

In [79]: x\_values = df.iloc[:,0:df.shape[1]].values

Out [80]: y\_values = df.iloc[:,df.shape[1]-1].values

In [80]: min\_max = MinMaxScaler()

Out [81]: min\_max.fit(x\_values)

Out [82]: x\_values\_s = min\_max.transform(x\_values)

In [81]: x\_values\_s

Out [81]: array([[0.99976406, 0.18, ..., ..., 0., ..., 0.00154526, ..., 0.],

Out [82]: [0.99976406, 0.29, ..., ..., 0., ..., 0.0020025, ..., 0.],

Out [83]: [0.99976406, 0.18, ..., ..., 0., ..., 0.00282212, ..., 0.],

Out [84]: [0.99976406, 0.34, ..., ..., 0., ..., 0.00515084, ..., 0.],

Out [85]: [0.99976406, 0.18, ..., ..., 0., ..., 0.003126, ..., 0.],

Out [86]: [0.99976406, 0.2, ..., ..., 0., ..., 0.00359128, ..., 0.],

Out [87]: ...])

In [82]: pca = PCA(n\_components=50)

Out [83]: pca.fit(x\_values\_s)

Out [84]: x\_components = pca.fit.transform(x\_values\_s)

In [83]: pca\_fit.explained\_variance\_ratio\_.sum()

Out [83]: 0.9921693943006986

In [84]: x\_components.shape

Out [84]: (71213, 50)

50 componentes explicam praticamente toda a informação (99,2%)

Fazendo testes sem balanceamento com validação cruzada

RandomForest, Regressão logística, SVM, NaiveBayes, LDA, KNN, XGBClassifier

In [89]: kf = KFold(n\_splits=10, shuffle=True)

In [123]: modelos = (('RandomForest', RandomForestClassifier()), ('KNN', KNeighborsClassifier()), ('Reg. L.', LogisticRegression()), ('SVM', SVC()), ('NB', GaussianNB()), ('LDA', LinearDiscriminantAnalysis()), ('XGBoost', XGBClassifier()))

In [ ]: resultados = []

for \_modelo in modelos:

resultados.append(cross\_val\_score(estimator=\_modelo, x=x\_components, y=y\_values, cv=kf, scoring='balanced\_accuracy'))

In [125]: plt.boxplot(x=resultados, labels=[x[0] for x in modelos]):



Os modelos RandomForest, Regressão Logística, LDA e XGBoost preveram perfeitamente.

Testando os modelos que se saíram melhor individualmente para coletar as acurácias por classe

In [85]: x\_treino, x\_teste, y\_treino, y\_teste = train\_test\_split(x\_components, y\_values, test\_size=0.3, random\_state=2)

In [86]: x\_treino.shape, x\_teste.shape

Out [86]: ((49849, 50), (21364, 50))

In [87]: randf = RandomForestClassifier(bootstrap=True, ccp\_alpha=0.0, class\_weight=None, criterion='gini', max\_depth=None, max\_features='auto', max\_leaf\_nodes=None, max\_samples=None, min\_impurity\_decrease=0.0, min\_impurity\_split=None, min\_samples\_leaf=1, min\_samples\_split=2, min\_weight\_fraction\_leaf=0.0, n\_estimators=100, n\_jobs=None, oob\_score=False, random\_state=None, verbose=0, warm\_start=False)

Out [87]:

In [88]: prev\_randf = randf.predict(x\_teste)

In [89]: print(classification\_report(y\_teste, prev\_randf))

|              | precision | recall | f1-score | support |
|--------------|-----------|--------|----------|---------|
| 0            | 1.00      | 1.00   | 1.00     | 20536   |
| 1            | 1.00      | 1.00   | 1.00     | 828     |
| accuracy     | 1.00      | 1.00   | 1.00     | 21364   |
| macro avg    | 1.00      | 1.00   | 1.00     | 21364   |
| weighted avg | 1.00      | 1.00   | 1.00     | 21364   |

O Random Forest já prevê perfeitamente, não há porque ir em frente