

# Deep Mailing - XGBoost Model - Reduzindo Dimensões

O objetivo desse notebook é demonstrar a utilização do XGBoost para a criação de árvores de decisão para a predição de CUPS em mailings.

Em primeiro lugar, definimos os imports que iremos usar...

```
In [1]: import xgboost
import numpy as np
import os
import sys
import logging
import gc
import pickle as pickle
import pandas as pd
import dateutil.parser as parser
import os.path
import math
from sklearn.metrics import accuracy_score
from datetime import datetime
import xgboost as xgb
from xgboost import XGBClassifier
from xgboost import plot_tree
import matplotlib.pyplot as plt
from matplotlib.pylab import rcParams
```

```
/home/ubuntu/git/DeepMailing/env/lib/python3.5/site-packages/sklearn/cross_validation.py:41: DeprecationWarning: This module was deprecated in version 0.18 in favor of the model_selection module into which all the refactored classes and functions are moved. Also note that the interface of the new CV iterators are different from that of this module. This module will be removed in 0.20.
```

```
"This module will be removed in 0.20.", DeprecationWarning)
```

Abaixo definimos os diretórios e nomes dos arquivos intermediários.

```
In [2]: log_location = "../logs/"
arquivo_df_pickled_norm = "../intermediate/df.norm.pickle"
arquivo_df_pickled_norm_train = "../intermediate/df.norm.train.pickle"
arquivo_df_pickled_norm_test = "../intermediate/df.norm.test.pickle"
arquivo_df_pickled_norm_train_x = "../intermediate/df.norm.train.x.pickle.npy"
arquivo_df_pickled_norm_train_y = "../intermediate/df.norm.train.y.pickle.npy"
arquivo_df_pickled_norm_test_x = "../intermediate/df.norm.test.x.pickle.npy"
arquivo_df_pickled_norm_test_y = "../intermediate/df.norm.test.y.pickle.npy"
```

Redefinimos o logger que iremos usar

```
In [3]: logger = logging.getLogger()
logging.basicConfig(format="%(asctime)-15s %(message)s",
                    level=logging.DEBUG,
                    filename=os.path.join(log_location, 'xgboost.log.' + datetime.now().strftime("%Y%m%d%H%M%S.%f") + '.log'))
```

Criamos uma função para imprimir tanto no log quanto no notebook...

```
In [4]: def print_log(msg):
        logging.debug(msg)
        print(msg)
```

Carregamos para a memória o arquivo normalizado e pickled que foi gerado no notebook de "Preparação de Dados".

```
In [5]: print_log("Carregando Pickling normalizado:{}".format(arquivo_df_pickled_norm))
        chamadas = pd.read_pickle(arquivo_df_pickled_norm)
```

Carregando Pickling normalizado:../intermediate/df.norm.pickle

Verificamos as dimensões do dataframe carregado.... E imprimimos uma amostra do dado que precisamos com apenas as colunas relevantes... Como podemos perceber o nosso modelo considera apenas colunas com valores booleanos (0 ou 1)

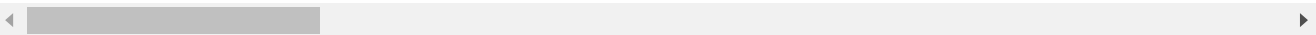
```
In [6]: print_log(chamadas.shape)
        chamadas.loc[:, 'NORM_CARTEIRA_A01':'NORM_DDD_87'].head(10)
```

(1435423, 103)

Out[6]:

	NORM_CARTEIRA_A01	NORM_CARTEIRA_W01	NORM_CARTEIRA_W02	NORM_CARTEIRA_
0	0	0	0	0
1	1	0	0	0
2	0	1	0	0
3	0	0	0	1
4	0	0	0	1
5	0	0	1	0
6	0	0	0	0
7	0	1	0	0
8	0	0	0	1
9	0	0	0	1

10 rows × 98 columns



Para tratarmos as colunas de forma correta, precisamos eliminar os caracteres especiais das colunas, principalmente o espaço e o sinal de - para isso iremos rodar o código abaixo para tratar.

```
In [7]: cols = chamadas.columns
        cols = cols.map(lambda x: x.replace(' ', '_').replace('-', 'menos'))
        chamadas.columns = cols
```

Um problema que temos é q temos muito poucos CUPS sendo acionados, então temos q remover linhas não relevantes no arquivo.

As linhas que devemos remover são aquelas em que não houve CUP e cujas features apesar de ativadas (terem valor 1), raramente são determinantes para um CUP.

No exemplo abaixo, eu vou determinar que se em menos de 15% das vezes em q determinada feature foi ativada, foi ativado um CUP também, eu vou determinar que essa coluna é irrelevante e que todos os valores com ativação dessa coluna e que não haja CUP seja apagados.

Dessa forma, iremos remover linhas ativações de features que não geraram `NORM_CUP = 1` e q sabemos q não são relevantes ao modelo.

No final, calculamos o threshold, que é o numero minimo de CUPs ativados para determinarmos se aquela feature é relevante ou não.

```
In [8]: minima_expressividade = 0.15
dimensionalidade_cup = minima_expressividade / (len(chamadas.query("NORM_CUP==1").index) / len(chamadas.index))
dimensionalidade_cup
```

```
Out[8]: 76.05561638996821
```

Como já calculamos o numero minimo de CUPs gerados para que uma determinada feature seja gerada, podemos agora determinar quais as colunas que devem ser limpas do dataframe

```
In [9]: def limpar_coluna_df(dataframe, coluna):  
        return dataframe.query("not ({}) == 1 and NORM_CUP==0").format(coluna))  
  
        def coluna_tem_cup(dataframe, coluna):  
            return len(dataframe.query("{} == 1 and NORM_CUP==1".format(coluna)).index) > dim  
ensionalidade_cup  
  
        colunas_para_limpar = []  
        lista_colunas = chamadas.loc[:, 'NORM_CARTEIRA_A01:'].columns.values  
        for coluna in list(lista_colunas):  
            if not coluna_tem_cup(chamadas,coluna):  
                colunas_para_limpar.append(coluna)  
                print_log("Coluna {} foi limpa do dataframe por ser irrelevante...".format(co  
luna))
```

[illegible]

Coluna NORM\_DDD\_45 foi limpa do dataframe por ser irrelevante...  
Coluna NORM\_DDD\_87 foi limpa do dataframe por ser irrelevante...

Tendo as colunas irrelevantes, podemos agora limpar elas do nosso dataframe, oque vai fazer com ele fique com menos tamanho e com melhor precisão.

```
In [10]: for coluna in colunas_para_limpar:
          chamadas = limpar_coluna_df(chamadas,coluna)
          gc.collect()
```

Podemos agora analisar o tamanho final do nosso dataframe:

```
In [11]: len(chamadas.index)
```

```
Out[11]: 449681
```

Perfeito, então vamos embaralhar os dados e gerar os nossos dados de treinamento e teste. Vamos considerar 70% para treinamento e 30 % para teste. No final, apagamos o dataframe lido para economizar memória

```
In [12]: print_log("Criando Pickling de train e teste...")
          chamadas = chamadas.sample(int(len(chamadas.index)))
          chamadas_train = chamadas.tail(int(len(chamadas.index) * 0.7))
          chamadas_test = chamadas.head(int(len(chamadas.index) * 0.3))
          del chamadas
```

Criando Pickling de train e teste...

Criamos uma função para gerar um arquivo de referencia de colunas a serem usadas q que vai ser importante na hora de gerar a arvore de decisao...

```
In [13]: def create_column_reference(header_chamadas_x,arquivo_df_pickled_norm_train_x):
          print_log("Criando Arquivo de referencia de colunas...")
          with open(arquivo_df_pickled_norm_train_x+".txt","w") as f:
              counter = 0
              lista_header = list(header_chamadas_x.columns.values)
              for header in lista_header:
                  f.write("{}-{}\n".format(counter,header))
                  counter=counter+1
```

Criamos agora os nossos dataframes de X que são as features e as Y que são os alvos de predição. Também removemos qualquer linha em tentativas seja igual a zero, além de criar um arquivo de referencia com as colunas X que serão usadas no modelo. Esses dataframes serão convertidos para matrizes no formato numpy

```
In [14]: print_log("Separando colunas em X e Y...")
chamadas_train = chamadas_train[(chamadas_train.NORM_TENTATIVAS > 0)]
create_column_reference(chamadas_train.loc[:, chamadas_train.columns.values[2]:'NORM_
DDD_87'].head(1), arquivo_df_pickled_norm_train_x)
chamadas_train_x = chamadas_train.loc[:, chamadas_train.columns.values[2]:'NORM_DDD_8
7'].as_matrix()
chamadas_train_y = chamadas_train.NORM_CUP.as_matrix()
chamadas_test = chamadas_test[(chamadas_test.NORM_TENTATIVAS > 0)]
chamadas_test_x = chamadas_test.loc[:, chamadas_test.columns.values[2]:'NORM_DDD_87'
].as_matrix()
chamadas_test_y = chamadas_test.NORM_CUP.as_matrix()
```

Separando colunas em X e Y...  
Criando Arquivo de referencia de colunas...

Após a criação das matrizes numpy, gravamos elas em arquivos.

```
In [15]: print_log("Criando arquivos finais em formato NUMPY para consumo pelo algoritmo...")

np.save(arquivo_df_pickled_norm_train_x,chamadas_train_x)
np.save(arquivo_df_pickled_norm_train_y,chamadas_train_y)
np.save(arquivo_df_pickled_norm_test_x,chamadas_test_x)
np.save(arquivo_df_pickled_norm_test_y,chamadas_test_y)
```

Criando arquivos finais em formato NUMPY para consumo pelo algoritmo...

Apagamos todos os dados intermediários e rodamos o garbage collector para economizar memória.

```
In [16]: print_log("Removendo objetos desnecessarios")
del chamadas_train_x
del chamadas_train_y
del chamadas_train
del chamadas_test
del chamadas_test_x
del chamadas_test_y
gc.collect()
```

Removendo objetos desnecessarios

Out[16]: 70

Carregamos os objetos numpy em memoria

```
In [17]: print_log("Carregando objetos numpy")
train_x = np.load(arquivo_df_pickled_norm_train_x)
train_y = np.load(arquivo_df_pickled_norm_train_y)
test_x = np.load(arquivo_df_pickled_norm_test_x)
test_y = np.load(arquivo_df_pickled_norm_test_y)
```

Carregando objetos numpy

Contamos quantos CUPS existem em treinamento e teste...

```
In [18]: msg1 = "Train - CUPS Detectados {} num universo de {}".format(len([y for y in train_y
    if y >0]),len(train_y))
msg2 = "Test - CUPS Detectados {} num universo de {}".format(len([y for y in test_y i
    f y >0]),len(test_y))
print_log(msg1)
print_log(msg2)
```

```
Train - CUPS Detectados 1989 num universo de 314776
Test - CUPS Detectados 842 num universo de 134903
```

Configuramos os parametros para o XGBoost, especificando que queremos que ele seja o mais exato possivel, que a medida de erro é erro simples e que queremos apenas uma classificacao binaria com um maximo de 1000 interações

```
In [19]: param = {}
param['eta'] = 0.2
param['objective'] = 'binary:logistic'
param['eval_metric'] = 'error'
param['tree_method'] = 'exact'
param['silent'] = 0
num_round = 1000
```

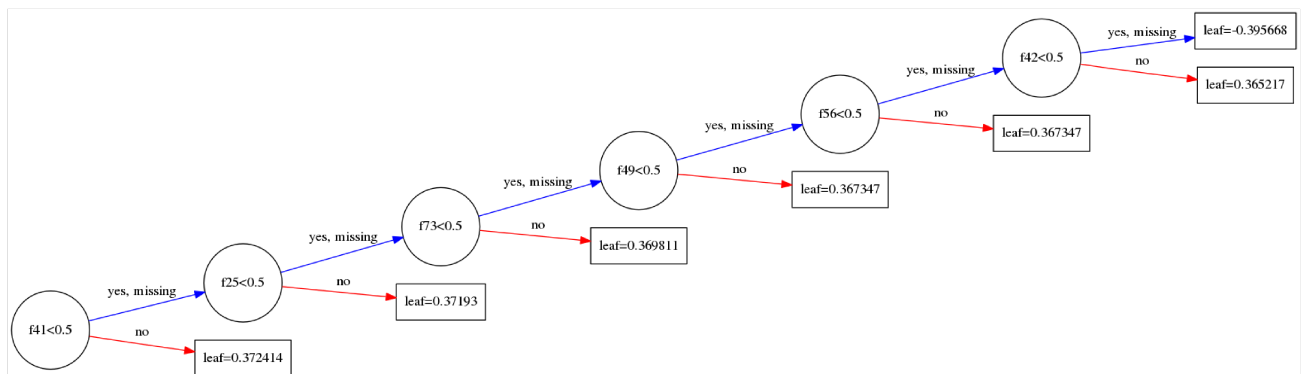
Após a definição dos parâmetros de teste, criamos as matrizes no formato do XGBoost e treinamos o modelo.

```
In [20]: gc.collect()
print_log("Starting model for params:{}".format(param))
dtrain = xgb.DMatrix(train_x, train_y)
dtest = xgb.DMatrix(test_x, test_y)
gpu_res = {}
booster = xgb.train(param, dtrain, num_round, evals=[], evals_result=gpu_res)
```

```
Starting model for params:{'silent': 0, 'eval_metric': 'error', 'eta': 0.2, 'tree_me
thod': 'exact', 'objective': 'binary:logistic'}
```

Após o modelo ser treinado, podemos plotar ele... Para verificar que coluna é cada feature no modelo, por favor ver a lista em anexo no final desse notebook.

```
In [21]: %matplotlib inline
rcParams['figure.figsize'] = 80,50
plot_tree(booster, rankdir='LR')
plt.show()
```



Agora, vamos tentar prever os dados com o nosso modelo treinado...

```
In [22]: test_y_pred = booster.predict(dtest)
test_predictions = np.array([value for value in test_y_pred])
```



E Finalmente medir a precisão da nossa predição... Tanto no total quanto em CUPs detectados.

```
In [23]: accuracy = accuracy_score(test_y, test_predictions.round())
print_log("CUPS Previstos:{}".format(len([x for x in test_predictions if x > 0.5])))
print_log("CUPS na Base Teste:{}".format(len([x for x in test_y if x > 0.5])))
print_log("Accuracy Total:{}".format(accuracy))
print_log("Accuracy em CUPs:{}".format(len([x for x in test_predictions if x > 0.5])
/ len([x for x in test_y if x > 0.5])))
```

```
CUPS Previstos:383
CUPS na Base Teste:842
Accuracy Total:0.9965975552804608
Accuracy em CUPs:0.45486935866983375
```

Após, vamos salvar o modelo gerado em um arquivo para reuso...

```
In [24]: save_file = "../output/{}.model".format(datetime.now().strftime("%Y%m%d.%H%M%S"))
with open(save_file, 'wb') as fp:
    pickle.dump(booster, fp)
print_log("Model saved as {}".format(save_file))
```

```
Model saved as ../output/20171215.180229.model
```

In [25]: %%bash

```
cat ../intermediate/df.norm.train.x.pickle.npy.txt
```

0-NORM\_CARTEIRA\_A01  
1-NORM\_CARTEIRA\_W01  
2-NORM\_CARTEIRA\_W02  
3-NORM\_CARTEIRA\_A02  
4-NORM\_CARTEIRA\_A03  
5-NORM\_SEGMENTO\_CR  
6-NORM\_SEGMENTO\_FA  
7-NORM\_SEGMENTO\_CC  
8-NORM\_SEGMENTO\_LC  
9-NORM\_SEGMENTO\_HC  
10-NORM\_SEGMENTO\_FT  
11-NORM\_SEGMENTO\_FC  
12-NORM\_SEGMENTO\_MA  
13-NORM\_PROPENSAO\_ALTA  
14-NORM\_ORIGEM\_BUREAU  
15-NORM\_ORIGEM\_BASE\_INTERNA  
16-NORM\_STATUS\_BUREAU\_Indefinido  
17-NORM\_STATUS\_BUREAU\_Bom\_menos  
18-NORM\_STATUS\_BUREAU\_Bom  
19-NORM\_STATUS\_BUREAU\_nan  
20-NORM\_STATUS\_INTERNA\_Bom  
21-NORM\_STATUS\_INTERNA\_Hot  
22-NORM\_STATUS\_INTERNA\_Novo  
23-NORM\_STATUS\_INTERNA\_Indefinido  
24-NORM\_STATUS\_INTERNA\_Validado  
25-NORM\_STATUS\_TELEFONE\_Indefinido  
26-NORM\_STATUS\_TELEFONE\_Hot  
27-NORM\_STATUS\_TELEFONE\_Validado  
28-NORM\_STATUS\_TELEFONE\_Bom  
29-NORM\_STATUS\_TELEFONE\_Bom\_menos  
30-NORM\_STATUS\_TELEFONE\_Novo  
31-NORM\_DDD\_63  
32-NORM\_DDD\_24  
33-NORM\_DDD\_96  
34-NORM\_DDD\_22  
35-NORM\_DDD\_38  
36-NORM\_DDD\_79  
37-NORM\_DDD\_69  
38-NORM\_DDD\_53  
39-NORM\_DDD\_34  
40-NORM\_DDD\_94  
41-NORM\_DDD\_84  
42-NORM\_DDD\_75  
43-NORM\_DDD\_88  
44-NORM\_DDD\_68  
45-NORM\_DDD\_14  
46-NORM\_DDD\_32  
47-NORM\_DDD\_85  
48-NORM\_DDD\_37  
49-NORM\_DDD\_91  
50-NORM\_DDD\_92  
51-NORM\_DDD\_99  
52-NORM\_DDD\_86  
53-NORM\_DDD\_97  
54-NORM\_DDD\_49  
55-NORM\_DDD\_33  
56-NORM\_DDD\_61  
57-NORM\_DDD\_65  
58-NORM\_DDD\_67  
59-NORM\_DDD\_44  
60-NORM\_DDD\_42  
61-NORM\_DDD\_89  
62-NORM\_DDD\_19  
63-NORM\_DDD\_77  
64-NORM\_DDD\_46

65-NORM\_DDD\_28  
66-NORM\_DDD\_27  
67-NORM\_DDD\_51  
68-NORM\_DDD\_71  
69-NORM\_DDD\_41  
70-NORM\_DDD\_48  
71-NORM\_DDD\_98  
72-NORM\_DDD\_18  
73-NORM\_DDD\_82  
74-NORM\_DDD\_11  
75-NORM\_DDD\_16  
76-NORM\_DDD\_64  
77-NORM\_DDD\_15  
78-NORM\_DDD\_66  
79-NORM\_DDD\_73  
80-NORM\_DDD\_43  
81-NORM\_DDD\_12  
82-NORM\_DDD\_95  
83-NORM\_DDD\_83  
84-NORM\_DDD\_55  
85-NORM\_DDD\_13  
86-NORM\_DDD\_21  
87-NORM\_DDD\_74  
88-NORM\_DDD\_35  
89-NORM\_DDD\_93  
90-NORM\_DDD\_81  
91-NORM\_DDD\_62  
92-NORM\_DDD\_17  
93-NORM\_DDD\_54  
94-NORM\_DDD\_47  
95-NORM\_DDD\_31  
96-NORM\_DDD\_45  
97-NORM\_DDD\_87

