

## ▼ Data Science com Python

### ▼ Módulo 5 - Modelagem Clustering

**Professor: Lucas Roberto Correa**

LEMBRETE: Fazer o import dos datasets usados no ambiente do colab antes de executar os comandos.

### ▼ Import dos pacotes

```
# Manipulação dados
import pandas as pd

# Visualização de dados
import seaborn as sns
import matplotlib.pyplot as plt

# Pre processamento
from sklearn.preprocessing import StandardScaler

# Modelos de agrupamento
from sklearn.cluster import KMeans
from sklearn.cluster import AgglomerativeClustering
from sklearn.cluster import DBSCAN

#Métricas
from sklearn.metrics import silhouette_score

# Limpeza de memória
import gc

pd.set_option('display.max_rows', None)
pd.set_option('display.max_columns', None)
pd.set_option('display.width', None)
pd.set_option('display.max_colwidth', None)
```

### ▼ Import dos metadados

link da base: [https://www.kaggle.com/rashmiranu/banking-dataset-classification?select=new\\_train.csv](https://www.kaggle.com/rashmiranu/banking-dataset-classification?select=new_train.csv)

```
meta = pd.read_excel('metadata.xlsx')
```

meta

	Feature	Feature_Type	
0	age	numeric	
1	job	Categorical,nominal	type of job ('admin.','blue-collar','entrepreneur','employed','services','student','unemployed')
2	marital	categorical,nominal	marital status ('divorced','married','single','unknown'; note: 'divorced' means divorced for more than 1 year)
3	education	categorical,nominal	('basic.4y','basic.6y','basic.9y','high.school','illiterate','professional.course','secondary.school','tertiary.degree')
4	default	categorical,nominal	has the client defaulted in any of the previous credit applications?
5	housing	categorical,nominal	Does the client have a home?
6	loan	categorical,nominal	Does the client have a current loan?
7	contact	categorical,nominal	contact channel: telephone, celluar, or other
8	month	categorical,ordinal	last contact month
9	dayofweek	categorical,ordinal	last contact day of the week
10	duration	numeric	last contact duration, in seconds . Important note: this attribute is not defined for telephone calls
11	campaign	numeric	number of contacts performed during this campaign across all channels
12	pdays	numeric	number of days that passed by after the client was last contacted by any channel
13	previous	numeric	number of contacts performed before this campaign

▼ Import da base

```
df = pd.read_csv('new_train.csv', sep=',')
```

```
df.head()
```

	age	job	marital	education	default	housing	loan	contact	month
0	49	blue-collar	married	basic.9y	unknown	no	no	cellular	nov
1	37	entrepreneur	married	university.degree	no	no	no	telephone	nov
2	78	retired	married	basic.4y	no	no	no	cellular	jul
3	36	admin.	married	university.degree	no	yes	no	telephone	may
4	59	retired	divorced	university.degree	no	no	no	cellular	jun

Retirando a target, pois o conjunto de dados será usado para uma análise não supervisionada

```
expl = df.drop(columns=['y'], axis=1)
expl.head()
```

	age	job	marital	education	default	housing	loan	contact	month
0	49	blue-collar	married	basic.9y	unknown	no	no	cellular	nov
1	37	entrepreneur	married	university.degree	no	no	no	telephone	nov
2	78	retired	married	basic.4y	no	no	no	cellular	jul
3	36	admin.	married	university.degree	no	yes	no	telephone	may
4	59	retired	divorced	university.degree	no	no	no	cellular	jun

```
expl_cat = expl[['job', 'marital', 'education', 'default', 'housing', 'loan',
                 'contact', 'month', 'poutcome']]
```

```
expl_num = expl[['duration', 'campaign', 'pdays', 'previous']]
```

Resultado a ser considerado na modelagem

```
expl_num.head()
```

	duration	campaign	pdays	previous
0	227	4	999	0
1	202	2	999	1
2	1148	1	999	0
3	120	2	999	0
4	368	2	999	0

Checagem de nulos

```
expl_num.isnull().sum()
```

```
duration    0
campaign    0
pdays      0
previous    0
dtype: int64
```

## Transformação dos dados com Padronização

```
'''
z = (x - u) / s
onde `u` é a média na amostra de train, `s` é o desvio padrão da amostra.
'''
```

```
scale = StandardScaler()
```

```
expl_num_scale = scale.fit_transform(expl_num)
```

```
expl_num_scale
```

```
array([[ -0.12019627,  0.52298128,  0.19658384, -0.35012691],
       [ -0.2167318 , -0.20368791,  0.19658384,  1.65381294],
       [  3.43617293, -0.56702251,  0.19658384, -0.35012691],
       ...,
       [ -0.49089273,  0.52298128,  0.19658384, -0.35012691],
       [ -0.3596044 , -0.56702251,  0.19658384, -0.35012691],
       [  1.10387435,  0.15964669,  0.19658384, -0.35012691]])
```

```
del df
```

```
gc.collect()
```

```
961
```

### ▼ O algoritmo

```
db = DBSCAN(eps=0.2)
db
```

```
DBSCAN(eps=0.2)
```

### ▼ Aplicando no conjunto de dados

```
expl['DB'] = db.fit_predict(expl_num_scale)
```

```
gc.collect()
```

```
74
```

### ▼ Avaliação de métrica

```
# DBSCAN
eps = [0.2, 0.3, 1]
```

```

for ep in eps:
    clusters = DBSCAN(eps=ep)
    predicao = clusters.fit_predict(expl_num_scale)

    score = silhouette_score(expl_num_scale, predicao)

    print('O valor de silhouette_score é {}, para n_clusters igual a {}'.format(score, ep))

    O valor de silhouette_score é 0.1299856350955824, para n_clusters igual a 0.2
    O valor de silhouette_score é 0.1282766845593569, para n_clusters igual a 0.3
    O valor de silhouette_score é 0.46496986207161534, para n_clusters igual a 1

```

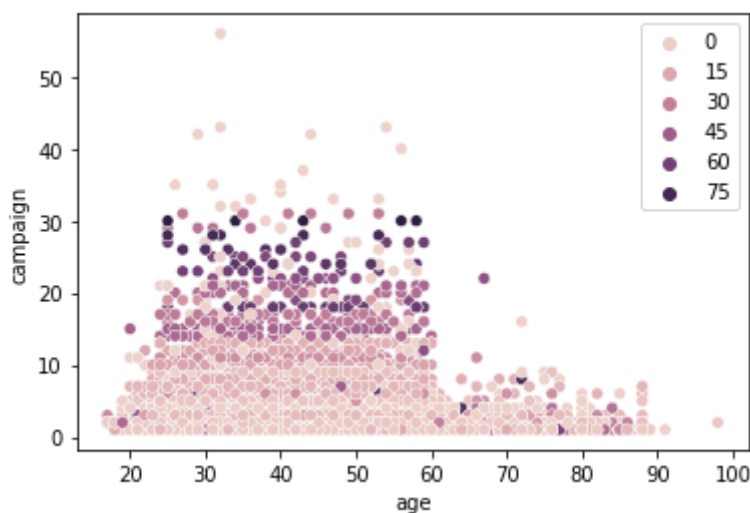
## ▼ Avaliação dos resultados considerando os dados de explicativas categóricas

```
expl.head()
```

	age	job	marital	education	default	housing	loan	contact	month
0	49	blue-collar	married	basic.9y	unknown	no	no	cellular	nov
1	37	entrepreneur	married	university.degree	no	no	no	telephone	nov
2	78	retired	married	basic.4y	no	no	no	cellular	jul
3	36	admin.	married	university.degree	no	yes	no	telephone	may
4	59	retired	divorced	university.degree	no	no	no	cellular	jun

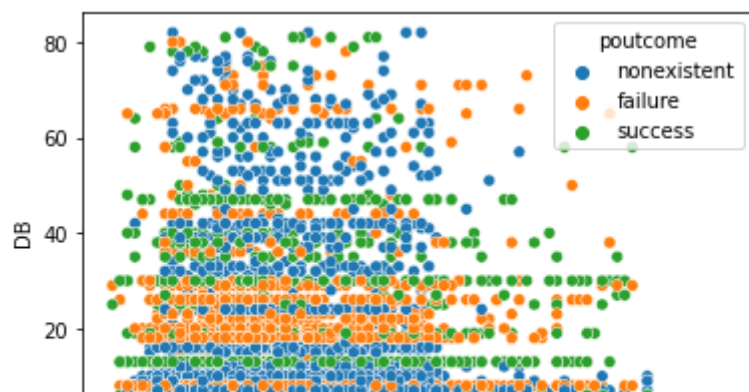
```
sns.scatterplot(data=expl, x="age", y="campaign", hue=db.labels_)
```

<matplotlib.axes.\_subplots.AxesSubplot at 0x7fa37a9bb850>



```
sns.scatterplot(x="age", y="DB", hue="outcome", data=expl)
```

```
<matplotlib.axes._subplots.AxesSubplot at 0x7fa37a8a0f90>
```



```
sns.scatterplot(x="DB", y="marital", hue="outcome", data=expl)
```

```
<matplotlib.axes._subplots.AxesSubplot at 0x7fa37a847910>
```

