

Bruno Vieira - Projeto de Regressão Logística

September 5, 2024

1 Projeto de Regressão Logística

Neste projeto estaremos trabalhando com um conjunto de dados falso de publicidade, indicando se um usuário de internet específico clicou ou não em uma propaganda. Vamos tentar criar um modelo que preveja se clicará ou não em um anúncio baseado nos recursos desse usuário.

Este conjunto de dados contém os seguintes recursos:

- ‘Daily Time Spent on Site’: tempo no site em minutos.
- ‘Age’: idade do consumidor.
- ‘Area Income’: Média da renda do consumidor na região.
- ‘Daily Internet Usage’: Média em minutos por di que o consumidor está na internet.
- ‘Linha do tópico do anúncio’: Título do anúncio.
- ‘City’: Cidade do consumidor.
- ‘Male’: Se o consumidor era ou não masculino.
- ‘Country’: País do consumidor.
- ‘Timestamp’: hora em que o consumidor clicou no anúncio ou janela fechada.
- ‘Clicked on Ad’: 0 ou 1 indicam se clicou ou não no anúncio.

1.1 Importar bibliotecas

**** Importe algumas bibliotecas que você acha que você precisará ****

```
[ ]: import pandas as pd
import numpy as np
import seaborn as sns
from sklearn.model_selection import train_test_split
from sklearn.linear_model import LogisticRegression
from sklearn.metrics import classification_report
```

1.2 Obter dados

**** Leia o arquivo advertising.csv e grave-o em um DataFrame chamado ad_data. ****

```
[ ]: ad_data = pd.read_csv('advertising.csv')
```

**** Verifique o cabeçalho do ad_data ****

```
[ ]: ad_data.head()
```

```
[ ]:      Daily Time Spent on Site  Age  Area Income  Daily Internet Usage  \
0                68.95    35      61833.90                256.09
1                80.23    31      68441.85                193.77
2                69.47    26      59785.94                236.50
3                74.15    29      54806.18                245.89
4                68.37    35      73889.99                225.58

                Ad Topic Line                City  Male    Country  \
0      Cloned 5thgeneration orchestration  Wrightburgh    0    Tunisia
1      Monitored national standardization    West Jodi    1      Nauru
2      Organic bottom-line service-desk    Davidton    0  San Marino
3  Triple-buffered reciprocal time-frame  West Terrifurt    1      Italy
4      Robust logistical utilization    South Manuel    0      Iceland

                Timestamp  Clicked on Ad
0  2016-03-27 00:53:11                0
1  2016-04-04 01:39:02                0
2  2016-03-13 20:35:42                0
3  2016-01-10 02:31:19                0
4  2016-06-03 03:36:18                0
```

**** Use info() e describe() em ad_data ****

```
[ ]: ad_data.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 1000 entries, 0 to 999
Data columns (total 10 columns):
#   Column                Non-Null Count  Dtype
---  -
0   Daily Time Spent on Site  1000 non-null  float64
1   Age                     1000 non-null  int64
2   Area Income             1000 non-null  float64
3   Daily Internet Usage     1000 non-null  float64
4   Ad Topic Line           1000 non-null  object
5   City                    1000 non-null  object
6   Male                    1000 non-null  int64
7   Country                 1000 non-null  object
8   Timestamp               1000 non-null  object
9   Clicked on Ad           1000 non-null  int64
dtypes: float64(3), int64(3), object(4)
memory usage: 78.2+ KB
```

```
[ ]: ad_data.describe()
```

```
[ ]:      Daily Time Spent on Site      Age  Area Income  \
count      1000.000000  1000.000000  1000.000000
mean         65.000200   36.009000  55000.000080
```

std	15.853615	8.785562	13414.634022
min	32.600000	19.000000	13996.500000
25%	51.360000	29.000000	47031.802500
50%	68.215000	35.000000	57012.300000
75%	78.547500	42.000000	65470.635000
max	91.430000	61.000000	79484.800000

	Daily Internet Usage	Male	Clicked on Ad
count	1000.000000	1000.000000	1000.000000
mean	180.000100	0.481000	0.500000
std	43.902339	0.499889	0.500250
min	104.780000	0.000000	0.000000
25%	138.830000	0.000000	0.000000
50%	183.130000	0.000000	0.500000
75%	218.792500	1.000000	1.000000
max	269.960000	1.000000	1.000000

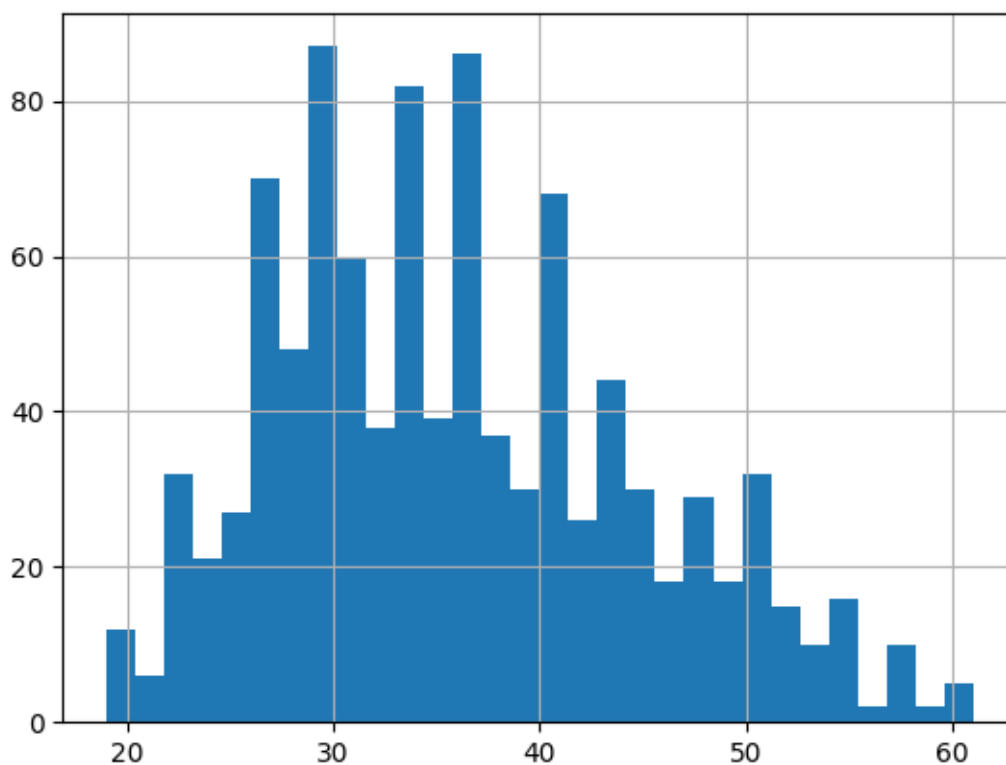
1.3 Análise de dados exploratória

Vamos usar Seaborn para explorar os dados!

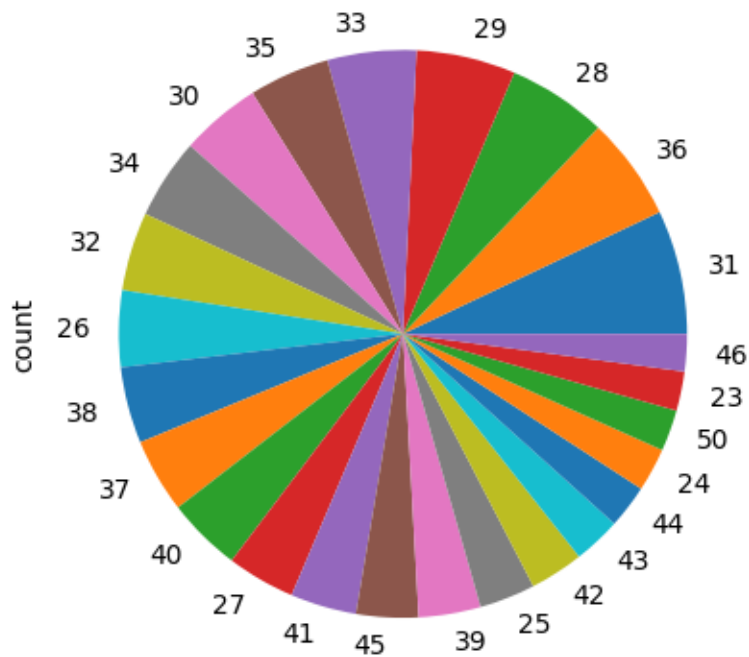
Tente recriar os gráficos abaixo.

**** Crie um histograma de “Age” ****

```
[ ]: histogram = ad_data['Age'].hist(bins=30)
```



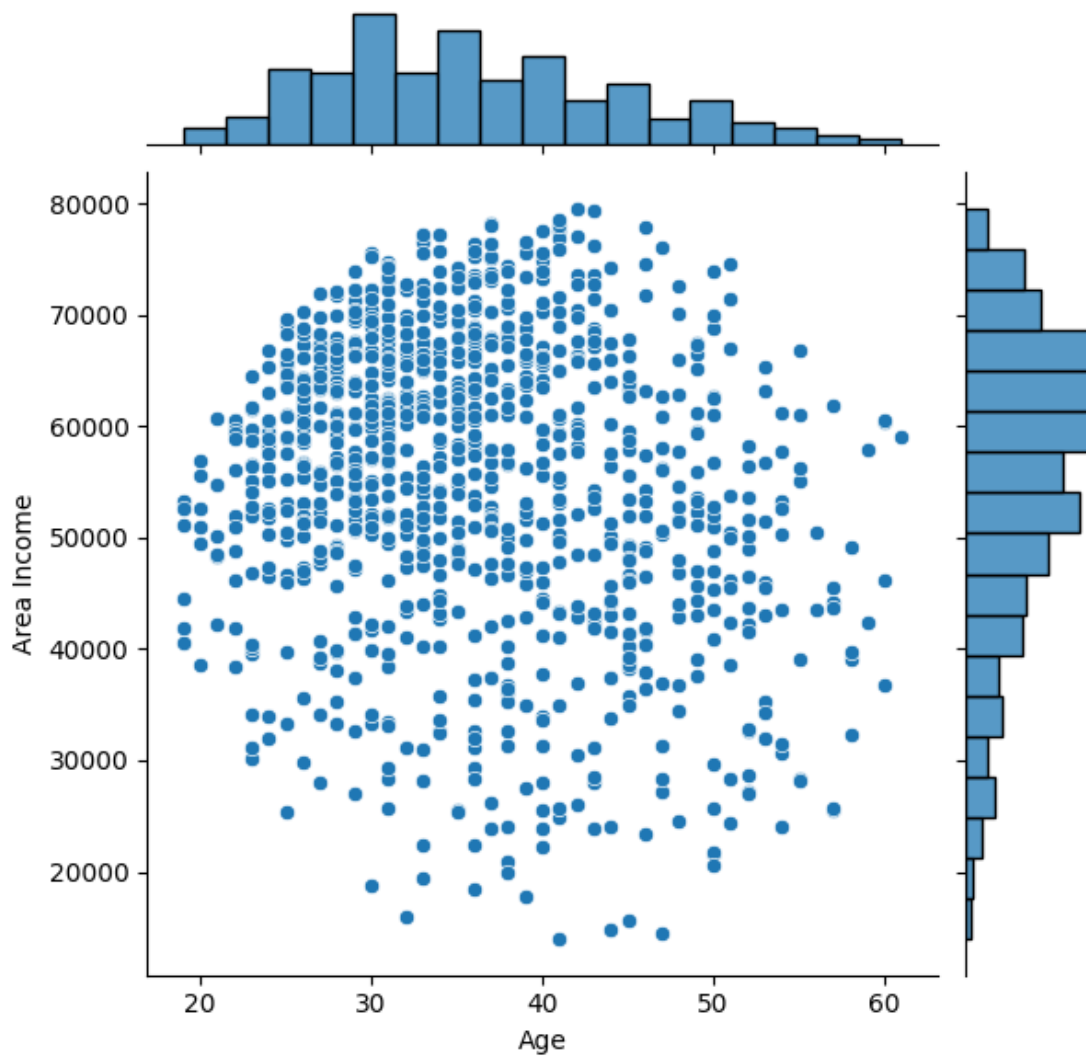
```
[ ]: pie = ad_data['Age'].value_counts().head(25).plot.pie()
```



A maior parte dos usuários estão entre 28 e 35 anos de idade.

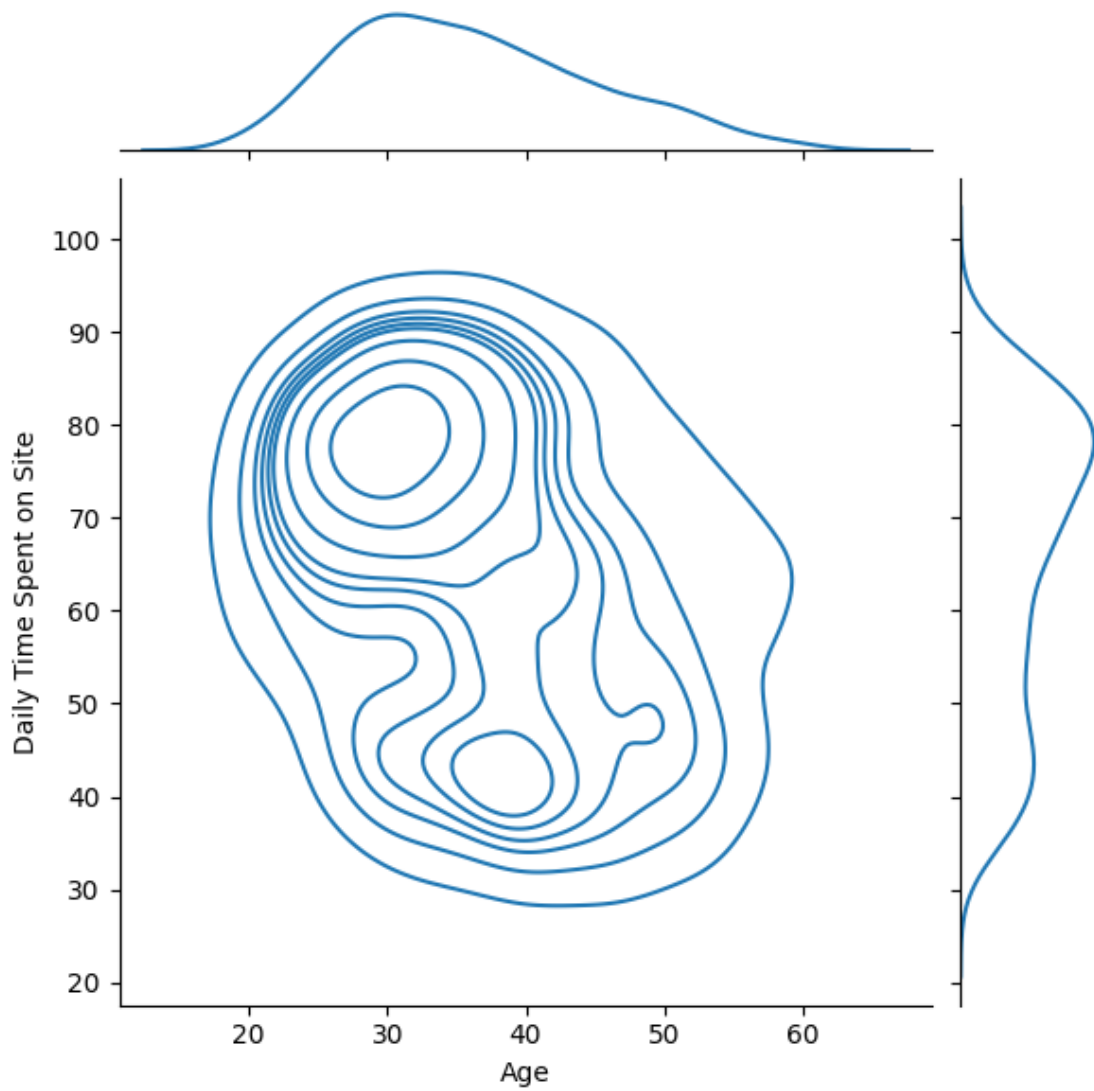
** Crie um joinplot mostrando “Area Income” versus “Age” **

```
[ ]: jointplot_1 = sns.jointplot(x='Age',y='Area Income',data=ad_data)
```



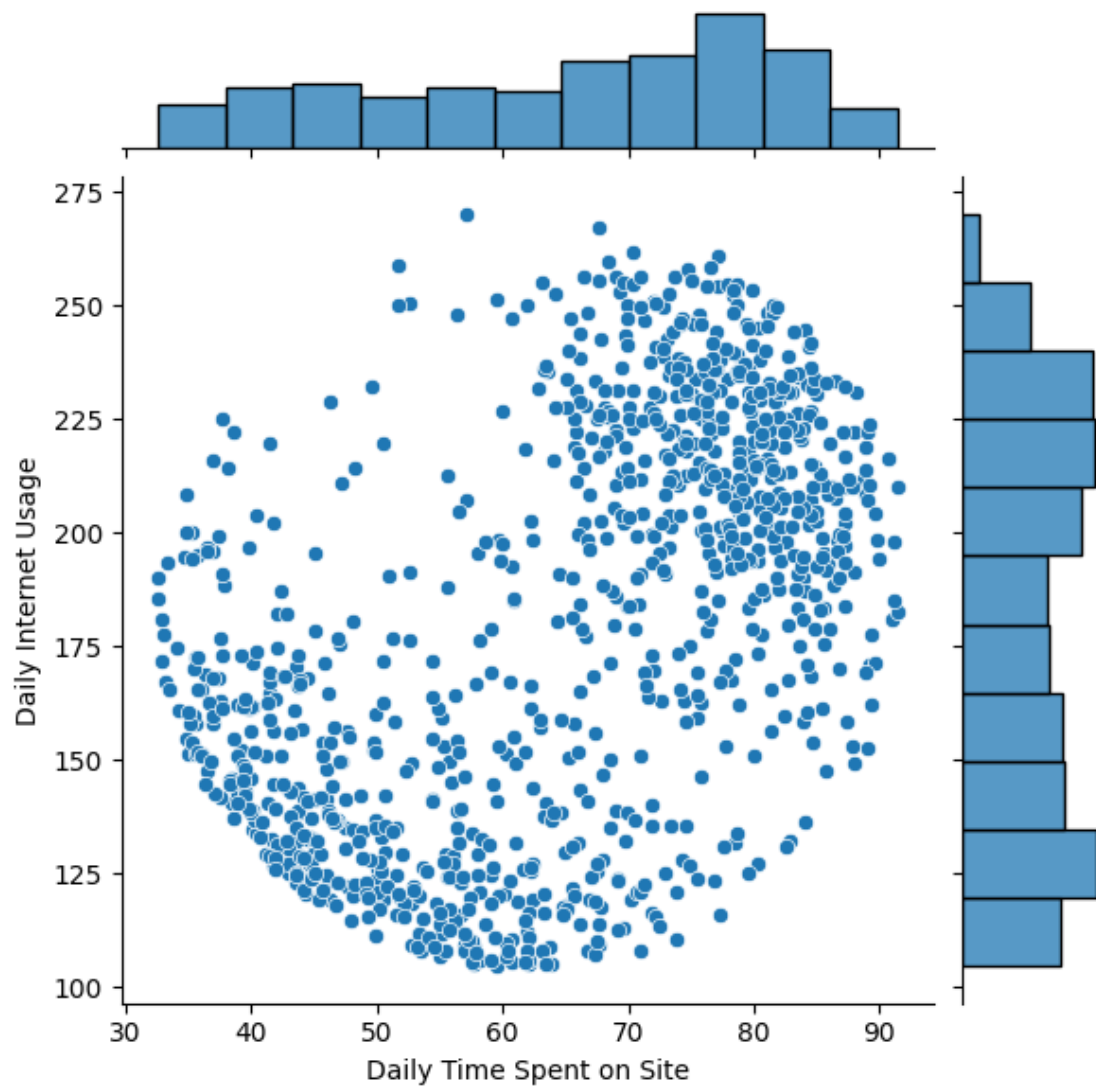
** Crie um jointplot que mostre as distribuições KDE do “Daily Time spent” no site vs “Age”. **

```
[ ]: jointplot_2 = sns.jointplot(x='Age',y='Daily Time Spent on_
↳Site',data=ad_data,kind='kde')
```



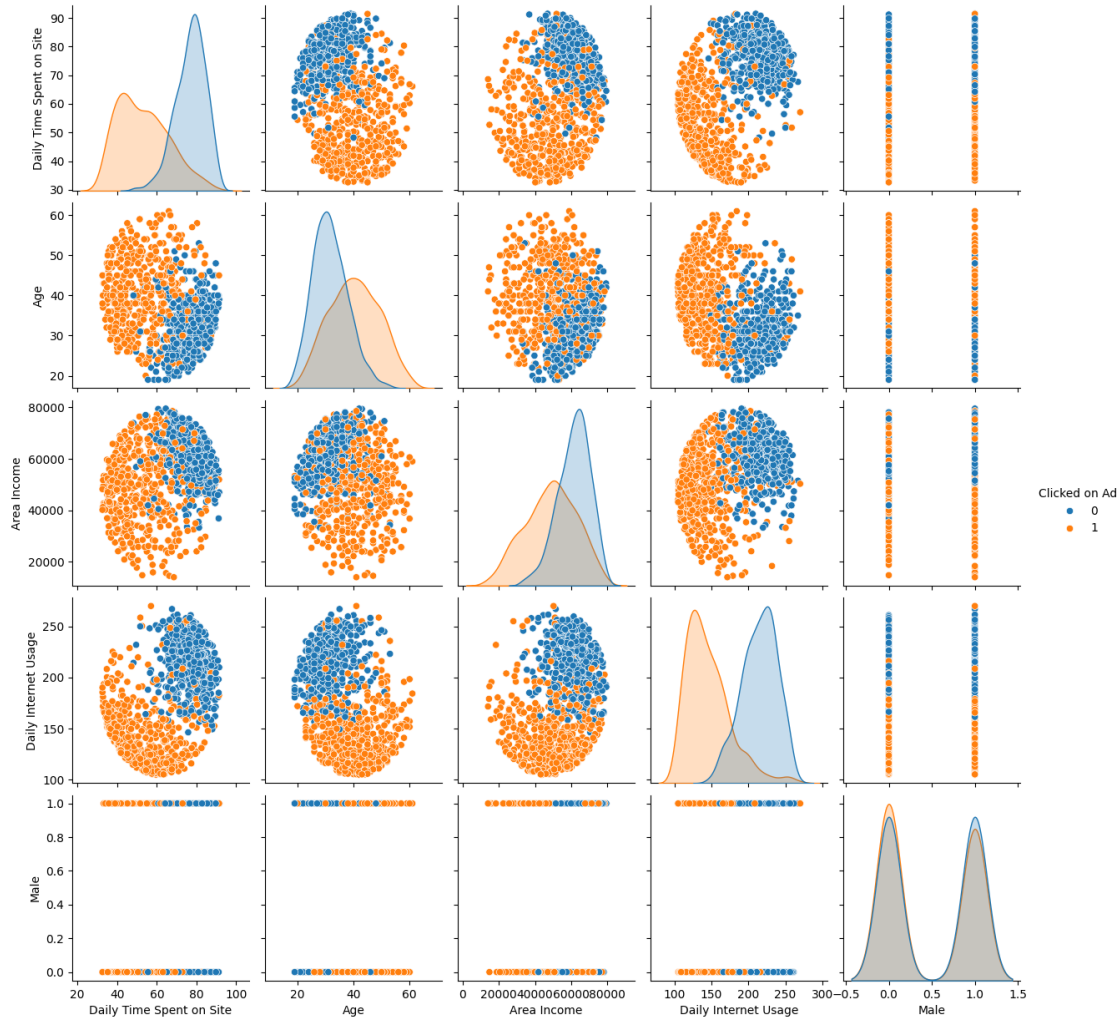
** Crie um jointplot do 'Daily Time Spent on Site' vs. 'Daily Internet Usage'**

```
[ ]: jointplot_3 = sns.jointplot(x='Daily Time Spent on Site',y='Daily Internet_Usage',data=ad_data)
```



** Finalmente, crie um parplot com o matiz definido pelo recurso de coluna 'Clicked on Ad'. **

```
[ ]: parplot = sns.pairplot(ad_data, hue='Clicked on Ad')
```



2 Regressão Logística

Agora é hora de quebrar nossos dados em treino e teste e fitar nosso modelo.

Você terá a liberdade aqui para escolher colunas em que deseja treinar!

**** Divida os dados em conjunto de treinamento e conjunto de testes usando train_test_split ****

```
[ ]: x = ad_data[['Daily Time Spent on Site', 'Age', 'Area Income']]
     y = ad_data['Clicked on Ad']
```

**** Treine e ajuste um modelo de regressão logística no conjunto de treinamento. ****

```
[ ]: x_train, x_test, y_train, y_test = train_test_split(x_test, y_test, test_size=0.
     ↪3, random_state=101)
```



```
[ ]: log = LogisticRegression()
log.fit(x_train,y_train)
print("Os coeficientes do modelo são:", log.coef_)
```

Os coeficientes do modelo são: $\begin{bmatrix} -2.07292261e-01 & 1.48236900e-01 \\ -1.22641002e-04 \end{bmatrix}$

2.1 Previsões e avaliações

**** Agora preveja valores para os dados de teste. ****

```
[ ]: log.predict(x_test)
```

```
[ ]: array([1, 0, 0, 1, 0, 1, 1, 1, 1, 1, 0, 1, 0, 1, 1, 1, 0, 1, 0, 1, 1,
          1, 0, 1, 0, 0, 0, 1, 1, 1, 1, 1, 0, 1, 0, 0, 0, 1, 1, 0, 1, 0,
          0, 0, 1, 0, 0, 0, 1, 0, 1, 1, 1, 0, 1, 1, 0, 1, 1, 1, 0, 0, 1,
          0, 0, 0, 1, 0, 0, 0, 0, 0, 1, 1, 0, 0, 0, 0, 1, 0, 1, 0, 1, 1, 1,
          0, 1, 0, 1, 0, 0, 1, 0, 0, 0, 0, 1, 0, 0, 1, 1, 1, 0, 0, 0, 0, 0,
          0, 0, 1, 1, 1, 0, 1, 1, 0, 0, 0, 1, 0, 0, 1, 0, 0, 1, 1, 0, 1, 0,
          0, 0, 0, 0, 1, 1, 0, 0, 1, 1, 1, 0, 0, 1, 0, 0, 0, 0, 1, 1, 0, 0,
          1, 0, 0, 0, 1, 0, 1, 0, 1, 1, 0, 0, 0, 0, 0, 0, 1, 0, 0, 1, 0,
          0, 0, 0, 0, 1, 0, 0, 0, 0, 1, 1, 1, 0, 1, 0, 1, 0, 0, 0, 1,
          1, 0, 1, 1, 0, 0, 1, 0, 0, 0, 1, 1, 0, 1, 0, 0, 0, 1, 0, 1, 0, 1, 1,
          1, 0, 1, 1, 0, 0, 0, 1, 1, 1, 1, 1, 1, 1, 1, 0, 0, 1, 1, 1, 1,
          1, 1, 1, 1, 0, 0, 1, 0, 0, 1, 0, 1, 0, 1, 0, 1, 1, 0, 0, 0, 1, 1,
          0, 1, 1, 1, 0, 0, 1, 0, 0, 0, 1, 0, 1, 0])
```

**** Crie um relatório de classificação para o modelo. ****

```
[ ]: y_pred = log.predict(x_test)
```

```
[ ]: print(classification_report(y_test, y_pred))
```

	precision	recall	f1-score	support
0	0.93	0.94	0.93	157
1	0.94	0.92	0.93	143
accuracy			0.93	300
macro avg	0.93	0.93	0.93	300
weighted avg	0.93	0.93	0.93	300