

Regressão Logística

Esse documento tem como principal objetivo mostrar a teoria da regressão logística, sua relação com a álgebra linear e um exemplo de implementação realizado a partir da programação. Assim, esse documento foi dividido nos seguintes subtópicos:

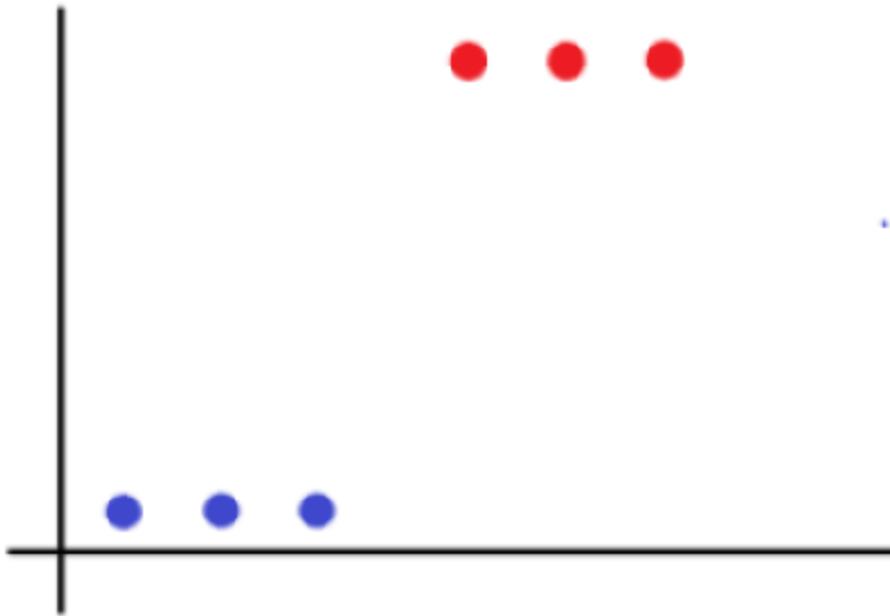
- Teoria
- Regressão Linear
- Implementação
- Conclusão

Teoria

A regressão logística é uma técnica de análise de dados com funcionalidades muito importantes voltadas, por exemplo, para a probabilidade e encontrar a relação entre duas variáveis distintas. Ela é muito utilizada em modelos de Machine Learning e Inteligência artificial, de acordo com a AWS, e possui diversos benefícios, como sua simplicidade de implementação, velocidade de processamento, flexibilidade e visibilidade.

Assim, é possível realizar uma análise preditiva a partir do resultado obtido da probabilidade de determinado cenário acontecer.

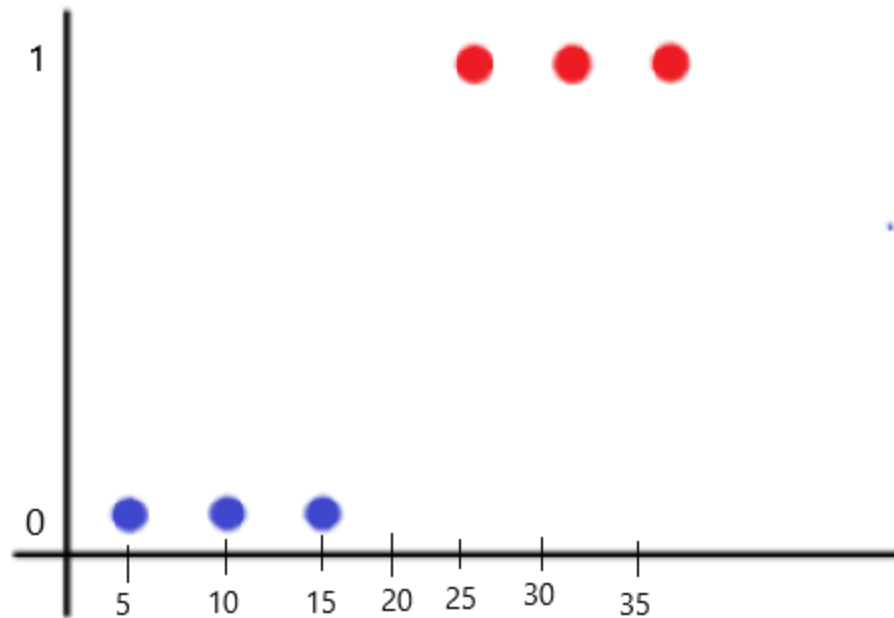
Vamos supor, por exemplo que somos uma agência de seguros e gostaríamos de saber a probabilidade de uma pessoa sofrer um acidente com base no tempo que ela passou na autoescola. Teríamos um gráfico da seguinte forma:



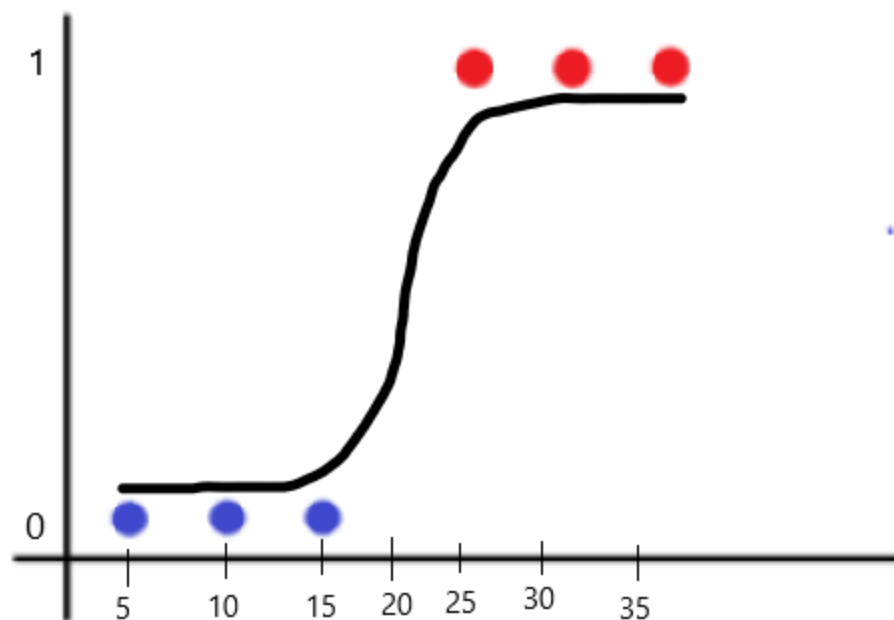
Os círculos vermelhos representam as pessoas que nunca sofreram acidente, enquanto os círculos azuis representam os que já sofreram, enquanto isso, o eixo x representa o número de dias que passaram na autoescola.

Se utilizássemos a regressão linear, ela não iria conseguir definir exatamente a probabilidade por causa da falta de linearidade em relação aos dados e seu limite em relação à variável binária de sofrer acidente. Isso ocorre, porque as pessoas que já sofreram acidente são classificadas em "Sim" ou "Não". Dessa forma, poderíamos representar numericamente esse gráfico da seguinte forma, sendo o eixo y se essas pessoas já sofreram acidente, e o eixo x o número de dias que eles permaneceram na autoescola.

Assim, a regressão logística geralmente é utilizada para realizar modelos de análise de dados de probabilidade de dados categóricos, muitas vezes binários como nesse exemplo.



Desse modo, ao utilizar a regressão logística, estaríamos desenhando uma sigmoide para entender a probabilidade de um dos resultados dessa variável binária acontecer, com isso, podemos perceber um gráfico da seguinte maneira:



Essa sigmoide apresenta a essa equação:

$$f(x) = \frac{1}{1 + e^{-x}}$$

A partir da equação, nota-se que todos os valores para $f(x)$ estão entre 1 e 0, pois, caso x tenha valor igual a infinito, seu limite tenderá a 0, com a equação se igualando a $1/1$. Enquanto isso, se seu expoente for igual a menos infinito, seu valor crescerá infinitamente e, conseqüentemente, tenderá a 0.

Nota-se que, quando os dados apresentados não possuem como resultado exatamente 1 e 0, eles são classificados a partir de sua proximidade, então todo número x que estiver classificado como $0.5 < x$, será

considerado como 1.

Para calcular essa probabilidade de ser um ou outro, podemos representar pela equação dada a seguir:

$$P(Y = 1 | X = x_i) = p_i$$

$$P(Y = 0 | X = x_i) = 1 - p_i$$

A primeira equação representa as chances de darem sucesso (1), enquanto a segunda representa as chances de darem fracasso (0).

Nesse caso não é utilizada a regressão linear, pois ela acaba superando o valor 1 do gráfico, podendo fornecer valores maiores que 1 e menor do que 0, violando o critério de probabilidade, cujos valores precisamos estar entre 0 e 1. Assim, a partir da regressão logística, temos essa condição de pé, além de uma possibilidade de ajuste de acordo com os dados apresentados no conjunto.

Modelos

A regressão logística também é considerada um modelo linear generalizado (MLG), ou seja, é uma generalização flexível da regressão linear ordinária que permite variáveis de resposta que têm modelos de distribuição de erros diferente de uma distribuição normal.

A partir disso, a regressão logística pode ter uma certa semelhança com algumas funções de ligação sendo os dois famosos modelos o **logit** e o **probit**. Pelo fato de serem modelos da regressão logística, eles são classificados modelos de escolha qualitativa e acabam compensando algumas limitações do Modelo de Probabilidade Linear.

1. Logit

O modelo logit é o mais tradicional quando é apresentada uma variável binária como variável dependente. De forma resumida, o logit é uma função que consiste aplicar o logaritmo nas *odds*, sendo possível calcular a razão de possibilidades a partir da diferença entre dois logits, como $\text{logit}(p) - \text{logit}(q)$. Por fim, é importante ressaltar que a função logit é representada pela função inversa da sigmoide, mostrada a seguir:

$$\text{Ln}\left(\frac{P}{1-P}\right) = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_k X_k$$

Nessa equação, P representa a probabilidade, sendo portanto, um número entre 0 e 1.

2. Probit

Essa função de ligação é bem mais complexa do que a logit e é caracterizada por uma distribuição normal cumulativa. Além disso, ela apresenta técnicas semelhantes e com o mesmo objetivo, tentar calcular a probabilidade de determinado resíduo pertencer a uma categoria ou outra.

Sua fórmula pode ser verificada a seguir:

$$PD_i = \Phi \left(\sum_{j=1}^n \beta_j X_{ij} \right)$$

Nota-se que Φ representa a distribuição normal cumulativa da função.

Regressão Linear

Além dessa sua importância e aplicações, podemos perceber que a equação da regressão logística apresenta uma certa semelhança com a equação da regressão linear, outra técnica muito utilizada para análise de dados e análise preditiva.

$$y = \beta_0 + \beta_1 x$$

Essa equação da regressão linear é composta pelos seguintes elementos:

y = variável resposta ou dependente

β_0 = intercepto

β_1 = coeficiente angular

x = variável explicativa ou independente

Acaba que a regressão linear pode ser representada em forma de matriz a partir da seguinte multiplicação:

$$\begin{bmatrix} y_1 \\ y_2 \\ y_3 \\ \vdots \\ y_n \end{bmatrix} = \begin{bmatrix} 1 & x_1 & x_1^2 & \dots & x_1^m \\ 1 & x_2 & x_2^2 & \dots & x_2^m \\ 1 & x_3 & x_3^2 & \dots & x_3^m \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & x_n & x_n^2 & \dots & x_n^m \end{bmatrix} \begin{bmatrix} \beta_0 \\ \beta_1 \\ \beta_2 \\ \vdots \\ \beta_m \end{bmatrix}$$

A equação mostrada anteriormente não coincide totalmente com a multiplicação mostrada, pois ela representa uma regressão linear simples, enquanto a matriz representa as variações de regressão linear a partir do conjunto de resíduos apresentados.

Nota-se que nessa multiplicação de matrizes, a matriz composto pelas diferentes variações de x é composto por equações lineares. A partir disso, é possível notar a relação que o sistema apresenta com a álgebra linear, seus conceitos e aplicações.

Mas a Regressão Logística é uma regressão linear ou não?

A regressão logística não é uma regressão linear, mas uma regressão linear generalizada, pois o modelo (Regressão de Bernoulli) em si não é linear por parâmetros (coeficientes), mas pode ser transformado em uma regressão linear (através de uma função de ligação, no caso a *logit*).

- Um modelo linear generalizado é linear nos coeficientes dos dados observados e não nos dados em si.

Linear

$$\bullet \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_3 \dots$$

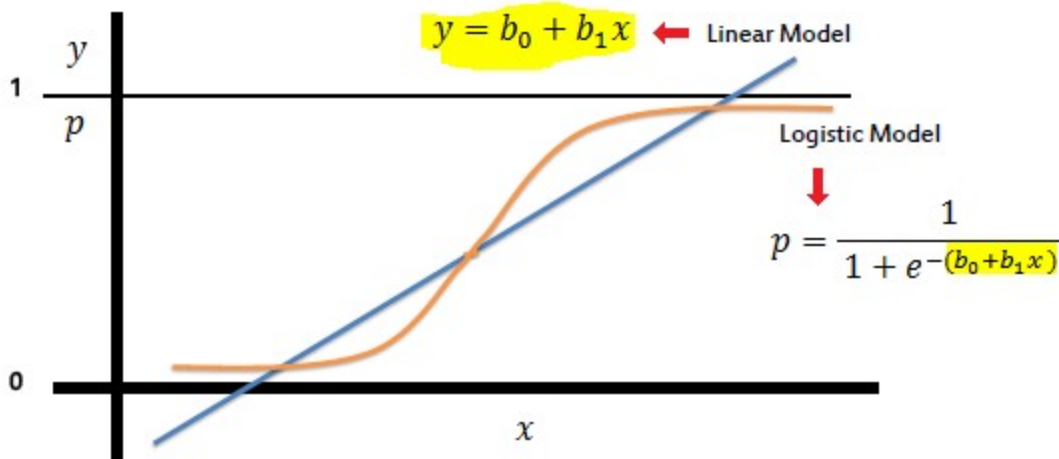
$$\bullet \beta_0 + \beta_1 x_1 + \beta_2 \sin(x_1) + \beta_2 e^{x_2} + \beta_3 \sqrt{x_1 x_2}$$

Não Linear

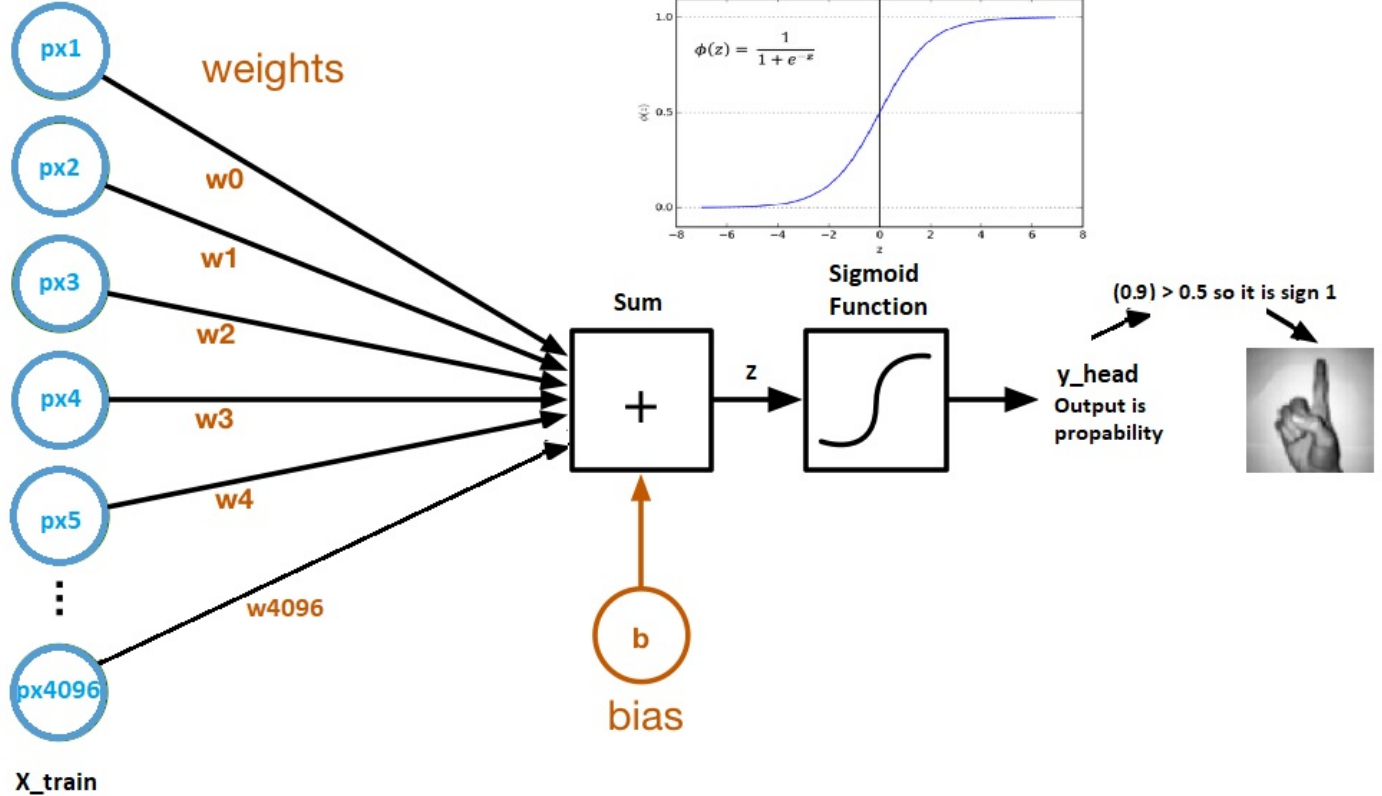
$$\bullet \beta_0 + \sin(\beta_1 x_1)$$

Tipo de regressão	utilização	exemplo de uso
Poisson	usada para modelar dados de contagem.	número de mortes em determinada região ou o número de consumidores que entram em um estabelecimento comercial.
Bernoulli	utilizada na modelagem de fenômenos que podem ser resumidos em uma variável binária, ou seja, se ocorreu ou não um evento.	modelos de concessão de crédito ou em pesquisas clínicas que tem como objetivo verificar os fatores de influência na ocorrência ou não de uma determinada doença.
Gama	usada para modelar dados positivos e assimétricos. A regressão Gama modela variáveis contínuas.	estudo dos fatores que influenciam no valor de um imóvel ou ainda os fatores que influenciam na demanda de produtos em diferentes centros de distribuição.

Fonte: <https://statplace.com.br/blog/o-que-sao-modelos-lineares-generalizados/>



Fonte: <https://www.researchgate.net/post/Is-Logistic-regression-a-linear-method>



Fonte: <https://www.kaggle.com/code/tanavbajaj/logistic-regression-math-behind-without-sklearn>

In [180...

```
from sklearn.linear_model import LogisticRegression
from sklearn.model_selection import train_test_split
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
```

CAR PRICE (<https://www.kaggle.com/datasets/shaistashaikh/carprice-assignment?resource=download>)

Importando Base

In [211...

```
#importando a base de dados
df = pd.read_csv("CarPrice_Assignment.csv")
df
```

Out[211...

	car_ID	symboling	CarName	fueltype	aspiration	doornumber	carbody	drivewheel	engine	location	wh
0	1	3	alfa-romero giulia	gas	std	two	convertible	rwd		front	
1	2	3	alfa-romero stelvio	gas	std	two	convertible	rwd		front	
2	3	1	alfa-romero Quadrifoglio	gas	std	two	hatchback	rwd		front	
3	4	2	audi 100 ls	gas	std	four	sedan	fwd		front	
4	5	2	audi 100ls	gas	std	four	sedan	4wd		front	
...	
200	201	-1	volvo 145e (sw)	gas	std	four	sedan	rwd		front	
201	202	-1	volvo 144ea	gas	turbo	four	sedan	rwd		front	

	car_ID	symboling	CarName	fueltype	aspiration	doornumber	carbody	drivewheel	engine	location	wheelbase
	202	203	-1	volvo 244dl	gas	std	four	sedan	rwd	front	98.6
	203	204	-1	volvo 246	diesel	turbo	four	sedan	rwd	front	94.5
	204	205	-1	volvo 264gl	gas	turbo	four	sedan	rwd	front	101.2

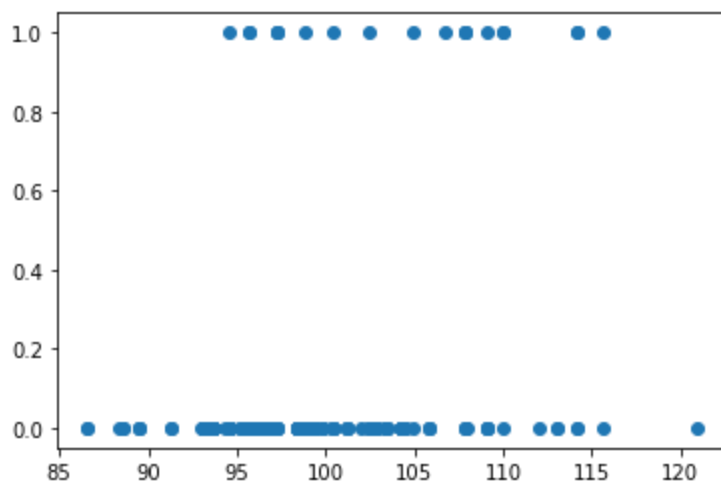
205 rows × 26 columns

```
In [182... #substituindo gas e diesel por 1 e 0
df2 = df.replace("diesel",1)
df3 = df2.replace("gas",0)
```

```
In [183... #definindo as variáveis independentes e a variável dependente sendo observada
y = df3["fueltype"]
x = df3.drop(df3.columns, axis=1)
```

Gráfico de Dispersão dos Dados

```
In [184... plt.figure()
plt.scatter(df3["wheelbase"], df3["fueltype"])
plt.show()
```



Definição de amostras para treinar o modelo e para testar o modelo

```
In [185... x_train, x_test, y_train, y_test = train_test_split(df3[["wheelbase"]], df3["fueltype"], train_size=0.8, random_state=42)
```

Dados de Teste

```
In [186... x_test
```

```
Out[186... wheelbase
62      98.8
163     94.5
10     101.2
```

wheelbase	
58	95.3
89	94.5
195	104.3
41	96.5
2	94.5
197	104.3
64	98.8
21	93.7
28	103.3
37	96.5
99	97.2
171	98.4
108	107.9
24	93.7
23	93.7
79	93.0
199	104.3
186	97.3
196	104.3
165	94.5
168	98.4
150	95.7
63	98.8
7	105.8
147	97.0
85	96.3
106	99.2
48	113.0
69	106.7
86	96.3
123	103.3
61	98.8
170	98.4
146	97.0
159	95.7
124	95.9

	wheelbase
45	94.5
126	89.5

Treinamento do modelo

In [187...

```
model = LogisticRegression()
model.fit(X_train, y_train)
print("Modelo treinado")
```

Modelo treinado

Previsão dos dados de teste

In [188...

```
y_predicted = model.predict(X_test)
print("Previsão realizada")
```

Previsão realizada

Acurácia do modelo

In [189...

```
acuracia = model.score(X_test, y_test)
print(f"A acurácia é de :{acuracia}")
```

A acurácia é de :0.9024390243902439

Comparação da previsão com a realidade dos dados de teste

0 - Gás

1 - Diesel

In [190...

```
df3_real = df3["fueltype"].filter(items=X_test.index)
previsao_realidade = pd.DataFrame((X_test))
previsao_realidade["previsao"] = y_predicted
previsao_realidade["realidade"] = df3_real
previsao_realidade
```

Out[190...

	wheelbase	previsao	realidade
62	98.8	0	0
163	94.5	0	0
10	101.2	0	0
58	95.3	0	0
89	94.5	0	0
195	104.3	0	0
41	96.5	0	0
2	94.5	0	0
197	104.3	0	0
64	98.8	0	0

	wheelbase	previsao	realidade
21	93.7	0	0
28	103.3	0	0
37	96.5	0	0
99	97.2	0	0
171	98.4	0	0
108	107.9	0	1
24	93.7	0	0
23	93.7	0	0
79	93.0	0	0
199	104.3	0	0
186	97.3	0	0
196	104.3	0	0
165	94.5	0	0
168	98.4	0	0
150	95.7	0	0
63	98.8	0	1
7	105.8	0	0
147	97.0	0	0
85	96.3	0	0
106	99.2	0	0
48	113.0	0	0
69	106.7	0	1
86	96.3	0	0
123	103.3	0	0
61	98.8	0	0
170	98.4	0	0
146	97.0	0	0
159	95.7	0	1
124	95.9	0	0
45	94.5	0	0
126	89.5	0	0

DIABETES (<https://www.kaggle.com/datasets/mathchi/diabetes-data-set>)

- Nº de Gravidezes X Diabetes

In [214...

```
diabetes = pd.read_csv("diabetes.csv")
display(diabetes)
```

```
diabetes.describe()
```

	Pregnancies	Glucose	BloodPressure	SkinThickness	Insulin	BMI	DiabetesPedigreeFunction	Age	Outcome
0	6	148	72	35	0	33.6	0.627	50	1
1	1	85	66	29	0	26.6	0.351	31	0
2	8	183	64	0	0	23.3	0.672	32	1
3	1	89	66	23	94	28.1	0.167	21	0
4	0	137	40	35	168	43.1	2.288	33	1
...
763	10	101	76	48	180	32.9	0.171	63	0
764	2	122	70	27	0	36.8	0.340	27	0
765	5	121	72	23	112	26.2	0.245	30	0
766	1	126	60	0	0	30.1	0.349	47	1
767	1	93	70	31	0	30.4	0.315	23	0

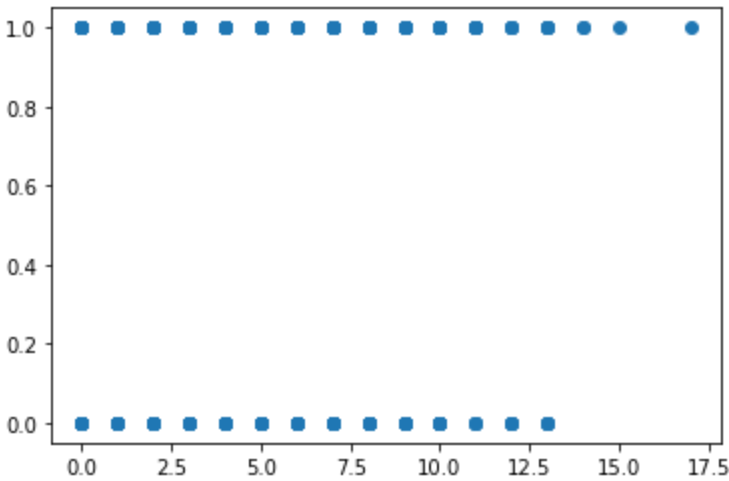
768 rows × 9 columns

Out[214...

	Pregnancies	Glucose	BloodPressure	SkinThickness	Insulin	BMI	DiabetesPedigreeFunction	
count	768.000000	768.000000	768.000000	768.000000	768.000000	768.000000	768.000000	768.000000
mean	3.845052	120.894531	69.105469	20.536458	79.799479	31.992578	0.471876	3.000000
std	3.369578	31.972618	19.355807	15.952218	115.244002	7.884160	0.331329	1.000000
min	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.078000	2.000000
25%	1.000000	99.000000	62.000000	0.000000	0.000000	27.300000	0.243750	2.000000
50%	3.000000	117.000000	72.000000	23.000000	30.500000	32.000000	0.372500	2.000000
75%	6.000000	140.250000	80.000000	32.000000	127.250000	36.600000	0.626250	4.000000
max	17.000000	199.000000	122.000000	99.000000	846.000000	67.100000	2.420000	8.000000

In [192...

```
plt.figure()
plt.scatter(diabetes["Pregnancies"], diabetes["Outcome"])
plt.show()
```



Definição de amostras para treinar o modelo e para testar o modelo

In [193...

```
x_train, x_test, y_train, y_test = train_test_split(diabetes[["Pregnancies"]],diabetes["Outcomes"],
```

Dados de Teste

In [194...

```
x_test
```

Out[194...

	Pregnancies
483	0
681	0
697	0
349	5
720	4
...	...
78	0
688	1
759	6
75	1
197	3

154 rows × 1 columns

Treinamento do modelo

In [195...

```
model = LogisticRegression()  
model.fit(x_train, y_train)  
print("Modelo treinado")
```

Modelo treinado

Previsão dos dados de teste

In [196...

```
y_predicted = model.predict(x_test)  
print("Previsão realizada")
```

Previsão realizada

Acurácia do modelo

In [197...

```
acuracia = model.score(x_test,y_test)  
print(f"A acurácia é de :{acuracia}")
```

A acurácia é de :0.6688311688311688

Comparação da previsão com a realidade dos dados de teste

1 - Diabética

0 - Não Diabética

In [198...

```
diabetes_real = diabetes["Outcome"].filter(items=X_test.index)
previsao_realidade = pd.DataFrame((X_test))
previsao_realidade["previsao"] = y_predicted
previsao_realidade["realidade"] = diabetes_real
previsao_realidade
```

Out[198...

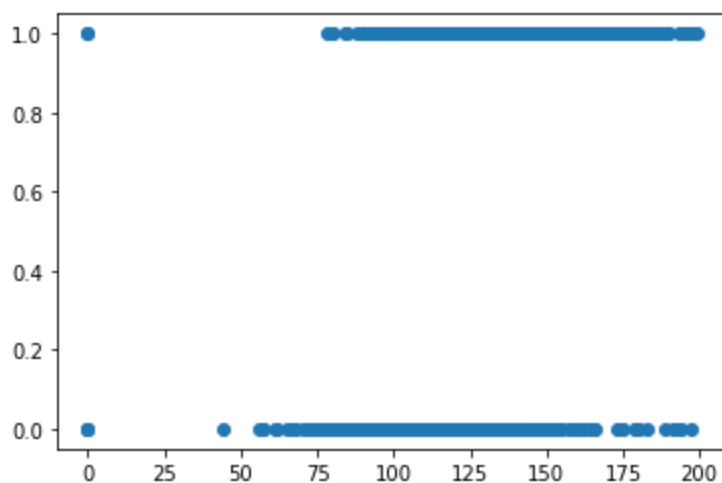
	Pregnancies	previsao	realidade
483	0	0	0
681	0	0	1
697	0	0	0
349	5	0	1
720	4	0	0
...
78	0	0	1
688	1	0	0
759	6	0	1
75	1	0	0
197	3	0	1

154 rows × 3 columns

- Glicose X Diabetes

In [199...

```
plt.figure()
plt.scatter(diabetes["Glucose"], diabetes["Outcome"])
plt.show()
```



Definição de amostras para treinar o modelo e para testar o modelo

In [200...

```
X_train, X_test, y_train, y_test = train_test_split(diabetes[["Glucose"]], diabetes["Outcor
```

Dados de Teste

In [201...

```
X_test
```

Out[201...

	Glucose
237	179
29	117
158	88
494	80
223	142
...	...
712	129
221	158
369	133
766	126
411	112

154 rows × 1 columns

Treinamento do modelo

In [202...

```
model = LogisticRegression()  
model.fit(X_train, y_train)  
print("Modelo treinado")
```

Modelo treinado

Previsão dos dados de teste

In [203...

```
y_predicted = model.predict(X_test)  
print("Previsão realizada")
```

Previsão realizada

Acurácia do modelo

In [204...

```
acuracia = model.score(X_test, y_test)  
print(f"A acurácia é de :{acuracia}")
```

A acurácia é de :0.7142857142857143

Comparação da previsão com a realidade dos dados de teste

1 - Diabética

0 - Não Diabética

In [205...

```
diabetes_real = diabetes["Outcome"].filter(items=X_test.index)  
previsao_realidade = pd.DataFrame((X_test))  
previsao_realidade["previsao"] = y_predicted
```

```
previsao_realidade["realidade"] = diabetes_real  
previsao_realidade
```

Out[205...

	Glucose	previsao	realidade
237	179	1	1
29	117	0	0
158	88	0	0
494	80	0	0
223	142	0	0
...
712	129	0	1
221	158	1	1
369	133	0	1
766	126	0	1
411	112	0	0

154 rows × 3 columns