

```
In [ ]: #Importando as bibliotecas necessárias para a análise
import pandas as pd
import matplotlib.pyplot as plt
plt.style.use("seaborn")
```

```
C:\Users\Filipe\AppData\Local\Temp\ipykernel_69336/1994541392.py:4: MatplotlibDeprecationWarning: The seaborn styles shipped by Matplotlib are deprecated since 3.6, as they no longer corre
spond to the styles shipped by seaborn. However, they will remain available as 'seaborn-v0.8+style'. Alternatively, directly use the seaborn API instead.
plt.style.use("seaborn")
```

```
In [ ]: #Importando o dataframe que iremos analisar
df = pd.read_csv('churn.csv', index_col='RowNumber')
```

```
In [ ]: df.head()
```

Out[]:	CustomerId	Surname	CreditScore	Geography	Gender	Age	Tenure	Balance	NumOfProducts	HasCrCard	IsActiveMember	EstimatedSalary	Exited
RowNumber													
1	15634602	Hargrave	619	France	Female	42	2	0.00	1	1	1	101348.88	1
2	15619304	Ono	502	France	Female	41	1	159660.80	1	0	1	119393.57	0
3	15619304	Ono	502	France	Female	42	8	159660.80	3	1	0	119393.57	1
4	15674012	Chu	645	Spain	Male	44	8	113755.78	2	1	0	149756.71	1
5	15656148	Ohnra	376	Germany	Female	29	4	115046.74	4	1	0	119346.88	1
17	15737452	Romeo	653	Germany	Male	58	1	132602.88	1	1	0	5097.67	1
...
9982	15672754	Burbridge	498	Germany	Male	42	3	152039.70	1	1	1	53445.17	1
9983	15769163	Giffin	655	Germany	Female	46	7	137145.12	1	1	0	115146.40	1
9992	15769959	Ajuchukwu	597	France	Female	53	4	88381.21	1	1	0	69384.71	1
9999	15584532	Liu	709	France	Female	36	7	0.00	1	0	1	42085.58	1
9999	15682355	Sablatira	772	Germany	Male	42	3	75075.31	2	1	0	92888.52	1

```
In [ ]: #Primeiro vamos verificar informações importantes do nosso dataframe
df.info()
<class 'pandas.core.frame.DataFrame'>
Int64Index: 10000 entries, 1 to 10000
Data columns (total 13 columns):
 #   Column                Non-Null Count  Dtype
---  --
 0   CustomerId            non-null        int64
 1   Surname               non-null        object
 2   CreditScore            non-null        int64
 3   NumOfProducts         non-null        int64
 4   Gender               non-null        object
 5   Age                  non-null        int64
 6   Tenure               non-null        int64
 7   Balance              non-null        float64
 8   NumOfProducts        non-null        int64
 9   HasCrCard            non-null        int64
10   IsActiveMember       non-null        int64
11   EstimatedSalary       non-null        float64
12   Exited               non-null        int64
dtypes: float64(2), int64(8), object(3)
memory usage: 1.1+ MB
```

Podemos verificar que não existem valores nulos em nosso dataframe, todos as 10 mil entradas estão devidamente preenchidas com algum valor.

Um resumo do que elas significam:

- CustomerId: identificação do cliente.
- Surname: sobrenome do cliente.
- CreditScore: pontuação de crédito. 0 alto risco de inadimplência e 1000 clientes com baixo risco de inadimplência.
- Geography: país que o serviço é oferecido.
- Gender: sexo do cliente.
- Age: idade do cliente.
- Tenure: um indicativo de estabilidade no emprego, em que 0 significa pouca estabilidade e 10 muita estabilidade.
- Balance: saldo da conta corrente.
- NumOfProducts: número de produtos bancários adquiridos.
- HasCrCard: se tem cartão de crédito ou não, (Sim = 1 e Não = 0).
- IsActiveMember: se é um cliente com conta ativa, (Ativo = 1).
- EstimatedSalary: salário estimado.
- Exited: cliente deixou de ser cliente do banco ou não (Churn = 1).

Como informado na introdução desse relatório, nossa análise passa entender o perfil dos clientes que deixaram o banco, para isso vamos filtrar nosso dataframe, retornando apenas clientes que saíram(Exited = 1). Para isso vamos utilizar o método query.

```
In [ ]: #Retornar todas pessoas que deixaram de ser clientes
df_query = df.query('Exited == 1')
```

Out[]:	CustomerId	Surname	CreditScore	Geography	Gender	Age	Tenure	Balance	NumOfProducts	HasCrCard	IsActiveMember	EstimatedSalary	Exited
RowNumber													
1	15634602	Hargrave	619	France	Female	42	2	0.00	1	1	1	101348.88	1
3	15619304	Ono	502	France	Female	42	8	159660.80	3	1	0	119393.57	1
6	15674012	Chu	645	Spain	Male	44	8	113755.78	2	1	0	149756.71	1
8	15656148	Ohnra	376	Germany	Female	29	4	115046.74	4	1	0	119346.88	1
17	15737452	Romeo	653	Germany	Male	58	1	132602.88	1	1	0	5097.67	1
...
9982	15672754	Burbridge	498	Germany	Male	42	3	152039.70	1	1	1	53445.17	1
9983	15769163	Giffin	655	Germany	Female	46	7	137145.12	1	1	0	115146.40	1
9992	15769959	Ajuchukwu	597	France	Female	53	4	88381.21	1	1	0	69384.71	1
9999	15584532	Liu	709	France	Female	36	7	0.00	1	0	1	42085.58	1
9999	15682355	Sablatira	772	Germany	Male	42	3	75075.31	2	1	0	92888.52	1

2037 rows x 13 columns

Podemos perceber que as pessoas que deixaram de ser clientes representam pouco mais de 20% do nosso dataframe. Como o foco da análise são os clientes que saíram(exited = 1), vou definir um novo dataframe com o nome df_churn apenas com esse filtro.

```
In [ ]: df_churn = df_query('Exited == 1')
df_churn
```

Out[]:	CustomerId	Surname	CreditScore	Geography	Gender	Age	Tenure	Balance	NumOfProducts	HasCrCard	IsActiveMember	EstimatedSalary	Exited
RowNumber													
1	15634602	Hargrave	619	France	Female	42	2	0.00	1	1	1	101348.88	1
3	15619304	Ono	502	France	Female	42	8	159660.80	3	1	0	119393.57	1
6	15674012	Chu	645	Spain	Male	44	8	113755.78	2	1	0	149756.71	1
8	15656148	Ohnra	376	Germany	Female	29	4	115046.74	4	1	0	119346.88	1
17	15737452	Romeo	653	Germany	Male	58	1	132602.88	1	1	0	5097.67	1
...
9982	15672754	Burbridge	498	Germany	Male	42	3	152039.70	1	1	1	53445.17	1
9983	15769163	Giffin	655	Germany	Female	46	7	137145.12	1	1	0	115146.40	1
9992	15769959	Ajuchukwu	597	France	Female	53	4	88381.21	1	1	0	69384.71	1
9999	15584532	Liu	709	France	Female	36	7	0.00	1	0	1	42085.58	1
9999	15682355	Sablatira	772	Germany	Male	42	3	75075.31	2	1	0	92888.52	1

2037 rows x 13 columns

A partir de agora vamos trabalhar com o dataframe acima

```
In [ ]: df_churn.info()
<class 'pandas.core.frame.DataFrame'>
Int64Index: 2037 entries, 1 to 9999
Data columns (total 13 columns):
 #   Column                Non-Null Count  Dtype
---  --
 0   CustomerId            non-null        int64
 1   Surname               non-null        object
 2   CreditScore            non-null        int64
 3   NumOfProducts         non-null        int64
 4   Gender               non-null        object
 5   Age                  non-null        int64
 6   Tenure               non-null        int64
 7   Balance              non-null        float64
 8   NumOfProducts        non-null        int64
 9   HasCrCard            non-null        int64
10   IsActiveMember       non-null        int64
11   EstimatedSalary       non-null        float64
12   Exited               non-null        int64
dtypes: float64(2), int64(8), object(3)
memory usage: 222.8+ KB
```

Com essas informações preliminares, vamos organizar os dados, separando as variáveis qualitativas e quantitativas. Iremos apagar a coluna CustomerId que é apenas uma identificação do cliente, o RowNumber já está servindo como nosso indexador, e também vou apagar o Surname, que propriamente dito é o sobrenome do cliente.

```
In [ ]: #Criando uma lista com as colunas que iremos deletar e deletando as mesmas em seguida
arr_delete = ['CustomerId', 'Surname']
df_churn.drop(arr_delete, axis=1, inplace = True)
```

#Criando listas com o nome de cada coluna

```
arr_qual = ['Geography', 'Gender', 'HasCrCard', 'IsActiveMember', 'Exited']
arr_quant1 = ['CreditScore', 'Age', 'Tenure', 'Balance', 'NumOfProducts', 'EstimatedSalary', 'Exited']
```

#Criando novos dataframes com base nas listas.

```
df_churn_qual1 = df_churn[arr_qual]
df_churn_quant1 = df_churn[arr_quant1]
```

```
C:\Users\Filipe\AppData\Local\Temp\ipykernel_69336/3519485126.py:15: SettingWithCopyWarning:
A value is being modified by inplace operations; this behavior has been deprecated.  In-place operations that modify pandas objects will raise the SettingWithCopyWarning in the future.  To avoid this warning in pandas, you can explicitly use df.copy() in place or use df = df.copy() on the right-hand-side.
```

See the caveats in the documentation: https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#returning-a-view-versus-a-copy

```
df_churn.drop(arr_delete, inplace = True)
```

Vamos analisar algumas métricas estatísticas com nossas variáveis quantitativas

```
In [ ]: df_churn_quant1.describe().round(2)
```

Out[]:	CreditScore	Age	Tenure	Balance	NumOfProducts	EstimatedSalary	Exited
count	2037.00	2037.00	2037.00	2037.00	2037.00	2037.00	2037.0
mean	645.35	44.84	4.93	91328.54	1.48	101465.68	1.0
std	100.82	9.76	2.84	98360.79	0.80	57922.42	0.0
min	292.00	35.00	0.00	0.00	1.00	11.56	1.0
25%	576.00	36.00	2.00	38340.02	1.00	51907.72	1.0
50%	646.00	45.00	5.00	109349.29	1.00	102460.84	1.0
75%	718.00	51.00	8.00	131433.33	2.00	154242.91	1.0
max	850.00	84.00	10.00	250898.09	4.00	199803.10	1.0

Alguns dados relevantes:

Média de idade bem próxima a mediana, baixa assimetria nessa variável(Média: 44.8 Mediana 45.0). Clientes que variam de 18 a 84 anos (A variável Age pode ser convertida para qualitativa, caso seja necessário criar um faixa de idade). Desvio padrão do salário(EstimatedSalary) relativamente alto, o que pode ser comum em uma carteira de clientes bancários.

```
In [ ]: #Clientes que saíram e que possuem idade acima da média arredondada(45)
df_churn_quant1.query('Age > 45')
```

Out[]:	CreditScore	Age	Tenure	Balance	NumOfProducts	EstimatedSalary	Exited
RowNumber							
17	653	58	1	132602.88	1	5097.67	1
36	475	45	0	134264.04	1	27822.99	1
42	465	51	8	12522.32	1	181297.65	1
44	834	49	2	131394.56	1	194365.76	1
59	511	66	4	0.00	1	1643.11	1
...
9925	632	50	5	107959.39	1	6985.34	1
9957	520	46	10	89216.61	1	117369.52	1
9976	610	50	1	113957.01	2	196526.55	1
9983	655	46	7	137145.12	1	115146.40	1
9992	597	53	4	88381.21	1	69384.71	1

1044 rows x 7 columns

```
In [ ]: #Clientes que saíram e que possuem idade abaixo da média(45)
df_churn_quant1.query('Age < 45')
```

Out[]:	CreditScore	Age	Tenure	Balance	NumOfProducts	EstimatedSalary	Exited
RowNumber							
1	619	42	2	0.00	1	101348.88	1
3	502	42	8	159660.80	3	113931.57	1
6	645	44	8	113755.78	2	149756.71	1
8	376	29	4	115046.74	4	119346.88	1
23	510	38	4	0.00	1	118913.53	1
...
9961	795	33	9	104852.72	1	120853.83	1
9982	702	44	9	0.00	1	96907.41	1
9982	498	42	3	152039.70	1	53445.17	1
9998	709	36	7	0.00	1	42085.58	1
9999	772	42	3	75075.31	2	92888.52	1

993 rows x 7 columns

```
In [ ]: #Vamos analisar os jovens
df_churn_quant1.query('Age > 17 & Age < 30')
```

Out[]:	CreditScore	Age	Tenure	Balance	NumOfProducts	EstimatedSalary	Exited
RowNumber							
8	376	29	4	115046.74	4	119346.88	1
47	829	27	9	132045.67	1	119708.21	1
87	750	22	3	121681.82	1	128643.35	1
115	721	28	9	154475.94	2	101300.94	1
165	683	29	0	133702.89	1	95802.54	1
...
9540	659	29	6	123102.12	1	56971.41	1
9955	779	29	3	46388.16	3	127939.26	1
9973	644	18	7	0.00	1	59645.24	1
9749	627	24	5	102773.20	2	56793.02	1
8826	487	28	10	126315.26	1	32349.29	1

124 rows x 7 columns

Podemos perceber que os jovens representam uma ínfima parcela dos antigos clientes do banco, 124 registros, aproximadamente 6%. Podemos visualizar um gráfico com a quantidade por idade

```
In [ ]: df_churn_quant1[["Age"]].value_counts().sort_values(ascending=True).plot.bar()
plt.title('Churn X Idade')
plt.xlabel('Idade')
plt.ylabel('Ex Clientes')
plt.figure(figsize=(9, 4))
```

```
Out[ ]: <Figure size 960x400 with 0 Axes>
```

Podemos notar que a maioria está na faixa dos 30 a 50 anos

Vamos retornar para o nosso dataframe principal e realizar algumas análises. Para isso vamos exibir uma amostra

```
In [ ]: df_churn.sample(18)
```

Out[]:	CreditScore	Geography	Gender	Age	Tenure	Balance	NumOfProducts	HasCrCard	IsActiveMember	EstimatedSalary	Exited
RowNumber											
23	510	Spain	Female	38	4	0.00	1	1	0	118913.53	1
8810	820	France	Female	39	1	104614.29	1	1	0	61538.43	1
9625	745	Germany	Female	44	0	11930.21	1	1	1	94266.08	1
6595	596	Germany	Male	48	2	13326.47	1	0	0	1146.02	1
82	777	France	Female	32	2	0.00	1	1	0	136458.19	1
1852	654	France	Male	53	8	144453.75	1	1	0	190986.96	1
2345	714	Germany	Female	49	4	93059.34	1	1	0	7571.51	1
9749	627	Germany	Male	24	5	102773.20	2	1	0	56793.02	1
4893	707	Germany	Female	51	10	89438.23	1	0	0	70778.63	1
6714	586	France	Female	46	0	0.00	3	0	1	131553.82	1

```
In [ ]: #quantidade de ex clientes por País
df_churn.groupby('value_counts')
```

Out[]:	Germany	France	Spain
count	824	810	413
Name: Geography, dtype: int64			

```
In [ ]: #visualizando esses dados em gráficos de barra
df_churn.groupby('value_counts').plot.bar()
plt
```