

Customer Churn?



Customer churn é uma métrica que informa quanto uma empresa perdeu de clientes.

Entender os fatores que levam um ou mais clientes é de extrema importância para a saúde e planejamento das empresas

Colunas presentes no dataset

Feature name	Description			
customerID	Client ID			
gender	Whether the customer is a male or a female			
SeniorCitizen	Whether the customer is a senior citizen or not (1, 0)			
Partner	Whether the customer has a partner or not (Yes, No)			
Dependents	Whether the customer has dependents or not (Yes, No)			
tenure	Number of months the customer has stayed with the company			
PhoneService	Whether the customer has a phone service or not (Yes, No)			
MultipleLines	Whether the customer has multiple lines or not (Yes, No, No phone service)			
InternetService	Customer's internet service provider (DSL, Fiber optic, No)			
OnlineSecurity	Whether the customer has online security or not (Yes, No, No internet service)			
OnlineBackup	Whether the customer has online backup or not (Yes, No, No internet service)			
DeviceProtection	Whether the customer has device protection or not (Yes, No, No internet service)			
TechSupport	Whether the customer has tech support or not (Yes, No, No internet service)			
StreamingTV	Whether the customer has streaming TV or not (Yes, No, No internet service)			
StreamingMovies	Whether the customer has streaming movies or not (Yes, No, No internet service)			
Contract	The contract term of the customer (Month-to-month, One year, Two year)			
PaperlessBilling	Whether the customer has paperless billing or not (Yes, No)			
PaymentMethod	The customer's payment method (Electronic check, Mailed check, Bank transfer (automatic), Credit card (automatic			
MonthlyCharges	The amount charged to the customer monthly			
TotalCharges	The total amount charged to the customer			
Churn	Whether the customer churned or not (Yes or No)			

Ao lado estão exibidas as colunas disponíveis no dataset.

Como pode-se observar na descrição, muitas das variáveis disponíveis são do tipo categórica.

Tipo das colunas presentes no dataset

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 7043 entries, 0 to 7042
Data columns (total 21 columns):
    Column
                      Non-Null Count
                                       Dtype
                       7043 non-null
                                       object
     customerID
    gender
                       7043 non-null
                                       object
    SeniorCitizen
                       7043 non-null
                                       int64
                       7043 non-null
                                       object
    Partner
    Dependents
                       7043 non-null
                                       object
                      7043 non-null
    tenure
                                       int64
    PhoneService
                      7043 non-null
                                       object
    MultipleLines
                       7043 non-null
                                       object
    InternetService
                       7043 non-null
                                       object
    OnlineSecurity
                       7043 non-null
                                       object
    OnlineBackup
                       7043 non-null
                                       object
    DeviceProtection 7043 non-null
                                       object
    TechSupport
                       7043 non-null
                                       object
    StreamingTV
                       7043 non-null
                                       object
    StreamingMovies
                       7043 non-null
                                       object
    Contract
                       7043 non-null
                                       object
    PaperlessBilling
                       7043 non-null
                                       object
    PaymentMethod
                                       object
                       7043 non-null
    MonthlyCharges
                       7043 non-null
                                       float64
    TotalCharges
                                       object
                       7043 non-null
    Churn
                       7043 non-null
                                       object
dtypes: float64(1), int64(2), object(18)
memory usage: 1.1+ MB
```

Como já era de se imaginar, a maioria dos dados presentes nesse dataset são do tipo "**object**", esse tipo de dado requer um cuidado especial para conseguir extrair informações quantitativas do mesmo.

Precisamos converter esse dado para float, já que é um número.

Primeira análise estatística da base de dados

	SeniorCitizen	tenure	MonthlyCharges	TotalCharges
count	7043.000000	7043.000000	7043.000000	7032.000000
mean	0.162147	32.371149	64.761692	2283.300441
std	0.368612	24.559481	30.090047	2266.771362
min	0.000000	0.000000	18.250000	18.800000
25%	0.000000	9.000000	35.500000	401.450000
50%	0.000000	29.000000	70.350000	1397.475000
75%	0.000000	55.000000	89.850000	3794.737500
max	1.000000	72.000000	118.750000	8684.800000

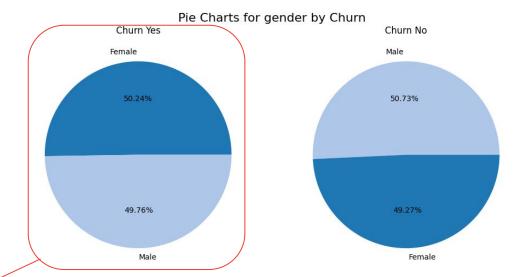


Sem o devido tratamento nos dados, temos uma perda de 75% das informações contidas no dataset. Uma vez que das 20 colunas iniciais só obtivemos informação numérica de 4.

Análise quantitativa da base de dados

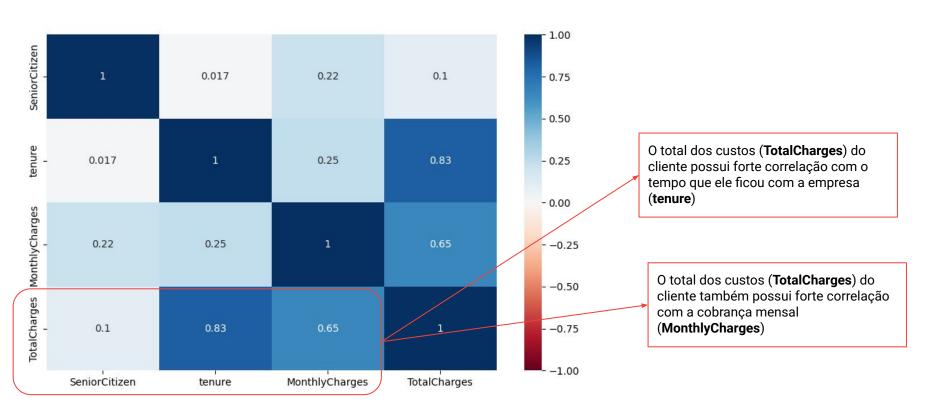
	count	unique	top	freq
customerID	7043	7043	7590-VHVEG	
gender	7043	2	Male	3555
Partner	7043	2	No	3641
Dependents	7043	2	No	4933
PhoneService	7043	2	Yes	6361
MultipleLines	7043	3	No	3390
InternetService	7043	3	Fiber optic	3096
Online Security	7043	3	No	3498
OnlineBackup	7043	3	No	3088
DeviceProtection	7043	3	No	3095
Tech Support	7043	3	No	3473
StreamingTV	7043	3	No	2810
StreamingMovies	7043	3	No	2785
Contract	7043	3	Month-to-month	3875
PaperlessBilling	7043	2	Yes	4171
PaymentMethod	7043	4	Electronic check	2365
Churn	7043	2	No	5174

Como a base de dados apresenta uma predominância inicial de dados categóricos, se faz necessário uma análise quantitativa da mesma.



Há praticamente um equilíbrio entre a quantidade de homens e mulheres que deixaram a empresa

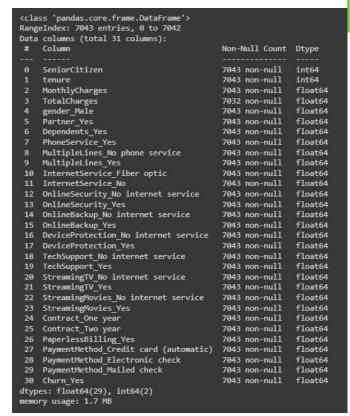
Análise da correlação entre os dados numéricos pré-existentes



Transformação dos dados categóricos

Podemos transformar os dados categóricos em dados numéricos de duas maneiras: **Dummy Encoding** e

OneHotEncoder

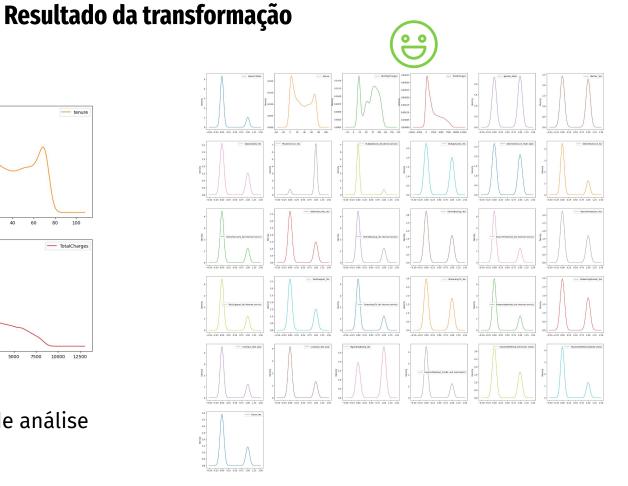


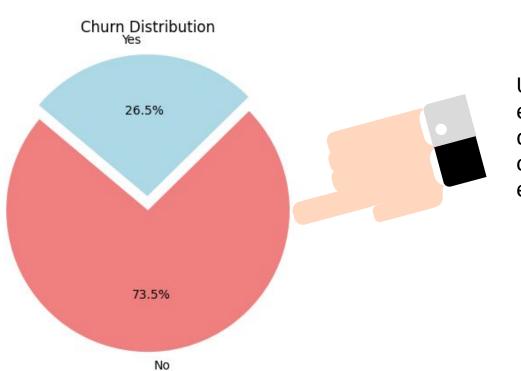


- SeniorCitizen 0.020 를 0.010 0.005 0.000 -0.50 -0.25 0.00 0.25 0.50 0.75 1.00 1.25 1.50 -40 0.00035 --- MonthlyCharges 0.0150 0.00030 0.0125 0.00025 0.0100 0.00020 0.0075 Ö 0.00015 0.0050 0.00010 0.0025 0.00005

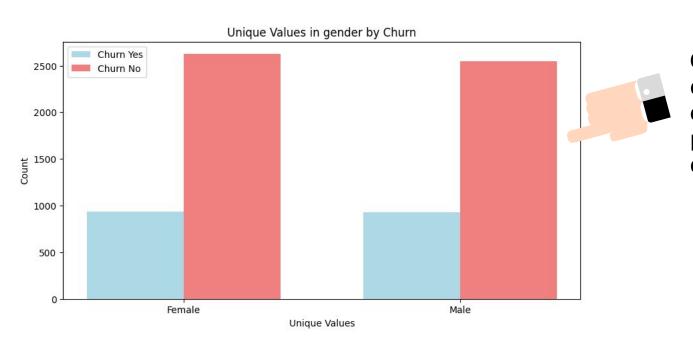
0.0000

Ampliamos nosso poder de análise gráfica na base de dados.

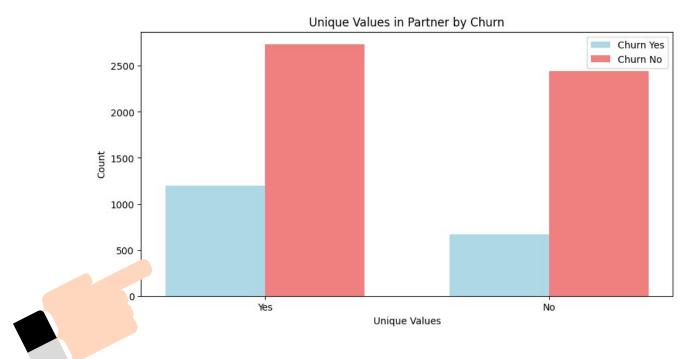




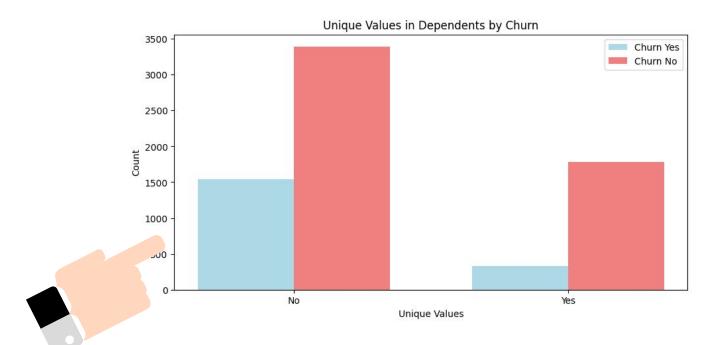
Uma base de dados extremamente desbalanceada, o que pode elevar a dificuldade da criação de modelos para fazer essa classificação.



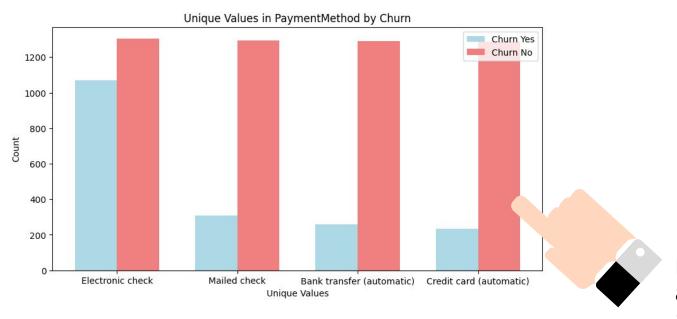
Com a relação de distribuição de gênero, essa base possui praticamente uma divisão de 50%.



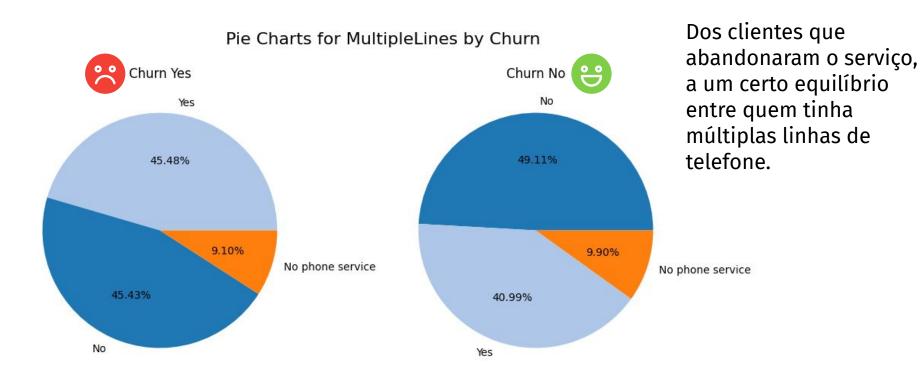
Na base de dados há uma predominância de solteiros.



Dos clientes que abandonaram o serviço, a maioria não tinha dependentes.

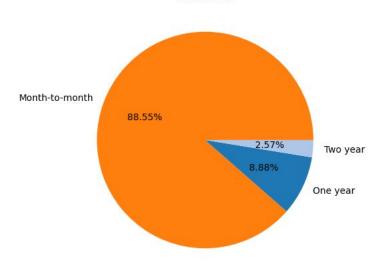


Dos clientes que abandonaram o serviço, a maioria utilizava cheque eletrônico com método de pagamento



Pie Charts for Contract by Churn





Dos clientes que abandonaram o serviço, 88.55% tinham contratos do tipo mensal.

Uma análise mais detalhada pode ser encontrada no notebook

