

Análise de correlações

In [1]:	<pre>import pandas as pd import matplotlib.pyplot as plt import numpy as np from mixtend.preprocessing import TransactionEncoder from mixtend.frequent_patterns import apriori from mixtend.frequent_patterns import association_rules import warnings warnings.filterwarnings('ignore') pd.options.display.float_format = "{:,.2f}".format pd.options.display.max_rows = 10</pre>
In [2]:	<pre>df = pd.read_excel("Online Retail.xlsx")</pre>

Conhecendo o dataset

Colunas

In [3]:	<pre>df.columns</pre>
Out[3]:	Index(['InvoiceNo', 'StockCode', 'Description', 'Quantity', 'InvoiceDate', 'UnitPrice', 'CustomerID', 'Country'], dtype='object')

Tipos de dados das colunas

In [4]:	<pre>df.dtypes</pre>
Out[4]:	InvoiceNo object StockCode object Description object Quantity int64 InvoiceDate datetime64[ns] UnitPrice float64 CustomerID float64 Country object dtype: object

Dataset

<

Limpeza dos dados

Removendo linhas vazias

In [6]:	<pre>df.dropna(axis='index', how='all', inplace=True)</pre>
---------	---

Removendo linhas onde a coluna 'InvoiceNo' esteja vazia

In [7]:	<pre>df.dropna(axis='index', subset=['InvoiceNo'], inplace=True)</pre>
---------	--

Exibindo o resultado

Removendo espaços errados (começo ou final) das strings

```
df['InvoiceNo'] = df['InvoiceNo'].str.strip()
df['StockCode'] = df['StockCode'].str.strip()
df['Description'] = df['Description'].str.strip()
df['Country'] = df['Country'].str.strip()
```

Filtrando resultados com dados de 'Quantity' e 'UnitPrice' coerentes com a realidade

```
df = df[(df['Quantity'] >= 0) & (df['UnitPrice'] >= 0)]
```

df

	InvoiceNo	StockCode	Description	Quantity	InvoiceDate	UnitPrice	CustomerID	Country
0	536365	85123A	WHITE HANGING HEART T-LIGHT HOLDER	6	2010-12-01 08:26:00	2.55	17850.00	United Kingdom

Convertendo as colunas que contenham texto para string

In [9]:	<pre>df['InvoiceNo'] = df['InvoiceNo'].astype('str') df['StockCode'] = df['StockCode'].astype('str') df['Description'] = df['Description'].astype('str') df['Country'] = df['Country'].astype('str')</pre>
---------	--

Removendo espaços errados (começo ou final) das strings

In [10]:	<pre>df['InvoiceNo'] = df['InvoiceNo'].str.strip() df['StockCode'] = df['StockCode'].str.strip() df['Description'] = df['Description'].str.strip() df['Country'] = df['Country'].str.strip()</pre>
----------	--

Filtrando resultados com dados de 'Quantity' e 'UnitPrice' coerentes com a realidade

Checando a quantidade de 'InvoiceNo' presente em cada país

```

df.groupby(['country'])['InvoiceNo'].count()

```

country	
Australia	1185
Austria	398
Bahrain	16
Belgium	2831
Brazil	32
...	
Switzerland	1967
USA	179
United Arab Emirates	68
United Kingdom	486284
Unspecified	446

Name: InvoiceNo, Length: 38, dtype: int64

Escolhendo a 'Germany' por possuir uma quantidade significativa de dados

Escolhendo país para análise

Checando a quantidade de 'InvoiceNo' presente em cada país

In [13]:	<pre>df.groupby(['Country'])['InvoiceNo'].count()</pre>
Out[13]:	Country Australia 1185 Austria 398 Bahrain 18 Belgium 2831 Brazil 32 Switzerland 1967 USA 279 United Arab Emirates 68 United Kingdom 486284 Unspecified 446 Name: InvoiceNo, Length: 38, dtype: int64

Escolhendo a 'Germany' por possuir uma quantidade significativa de dados

In [14]:	<pre>pais = 'Germany'</pre>
In [15]:	<pre>df_germany = df[df['Country'] == pais]</pre>

Preparação dos dados para construir a correlação

Agrupando por 'InvoiceNo' (id da compra/transação) e 'Description' para separar o dataset em transações detalhadas (carrinhos de compras)

In [16]:	<pre>df_grouped_germany = df_germany.groupby(['InvoiceNo', 'Description'])['Quantity'].sum()</pre>
----------	--

Resultado do agrupamento

InvoiceNo																			
536527	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
536840	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
536881	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
536967	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
536983	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
...
581266	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
581494	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
581570	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
581574	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
581578	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0

Usando função 'pandas.unstack' para selecionar uma coluna e transformar todas as suas linhas em coluna.

In [18]:	<pre>carrinho_compras_germany = df_grouped_germany.unstack(level=1, fill_value=0)</pre>
----------	---

Resultado da transformação, observe que agora todas as linhas são equivalentes as compras (carrinhos) e que cada coluna é um tipo de produto e o valor a quantidade comprada por aquele cliente.

```
carrinho_compras_germany.reset_index().set_index('InvoiceNo')
```

```
# carrinho_compras_germany.index
```

Convertendo as quantidades dos itens comprados para booleano, pois para o algoritmo basta estar ou não no carrinho.

```
carrinho_compras_germany = carrinho_compras_germany.applymap(lambda item: 1 if item > 0 else 0)
```

```
carrinho_compras_germany
```

	10	10	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12
--	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----

Checando se a coluna 'InvoiceNo' é o índice do dataframe

In [20]:	<pre>carrinho_compras_germany.index</pre>
Out[20]:	Index(['536527', '536840', '536861', '536967', '536983', '537197', '537198', '537201', '537212', '537259', '...', '580648', '581890', '581179', '581183', '581184', '581266', '581494', '581570', '581574', '581578'], dtype='object', name='InvoiceNo', length=457)

Caso não fosse, bastaria utilizar 'set_index'

In [21]:	<pre># carrinho_compras_germany.reset_index().set_index('InvoiceNo')</pre>
In [22]:	<pre># carrinho_compras_germany.index</pre>

Convertendo as quantidades dos itens comprados para booleano, pois para o algoritmo basta estar ou não no carrinho.

In [23]:	<pre>carrinho_compras_germany = carrinho_compras_germany.applymap(lambda item: 1 if item > 0 else 0)</pre>
In [24]:	<pre>carrinho_compras_germany</pre>

Description	10 COLOUR SPACEBOY PEN	10 COLOURED PARTY BALLOONS	12 IVORY ROSE PEG PLACES SETTINGS	12 MESSAGE CARDS WITH ENVELOPES	12 PENCIL SMALL TUBE WOODLAND	12 PENCILS SMALL TUBE RED RETROSPOT	12 PENCILS SMALL TUBE SKULL	12 PENCILS TALL TUBE POSY	12 PENCILS TALL TUBE RED RETROSPOT	12 PENCILS TALL TUBE RED TUBE SKULLS	...	YULETIDE IMAGES GIFT WRAP SET	ZINC HEART T-LIGHT HOLDER	ZINC STAR T-LIGHT HOLDER	ZINC BOX SIGN HOME	ZINC FOLKART SLEIGH BELLS	ZINC HEART LATTICE T-LIGHT HOLDER
InvoiceNo																	
536527	0	0	0	0	0	0	0	0	0	0	...	0	0	0	0	0	0
536840	0	0	0	0	0	0	0	0	0	0	...	0	0	0	0	0	0
536861	0	0	0	0	0	0	0	0	0	0	...	0	0	0	0	0	0
536967	0	0	0	0	0	0	0	0	0	0	...	0	0	0	0	0	0
536983	0	0	0	0	0	0	0	0	0	0	...	0	0	0	0	0	0
...
581266	0	0	0	0	0	0	0	0	0	0	...	0	0	0	0	0	0
581494	0	0	0	0	0	0	0	0	0	0	...	0	0	0	0	0	0
581570	0	0	0	0	0	0	0	0	0	0	...	0	0	0	0	0	0
581574	0	0	0	0	0	0	0	0	0	0	...	0	0	0	0	0	0
581578	0	0	0	0	0	0	0	0	0	0	...	0	0	0	0	0	0

457 rows x 1694 columns

Removendo a coluna 'POSTAGE' que representa envio e não um produto

In [25]:	<pre>carrinho_compras_germany.drop(labels=['POSTAGE'], inplace=True, axis='columns')</pre>
In [26]:	<pre>carrinho_compras_germany</pre>

Description	10 COLOUR SPACEBOY PEN	10 COLOURED PARTY BALLOONS	12 IVORY ROSE PEG PLACES SETTINGS	12 MESSAGE CARDS WITH ENVELOPES	12 PENCIL SMALL TUBE WOODLAND	12 PENCILS SMALL TUBE RED RETROSPOT	12 PENCILS SMALL TUBE SKULL	12 PENCILS TALL TUBE POSY	12 PENCILS TALL TUBE RED RETROSPOT	12 PENCILS TALL TUBE RED TUBE SKULLS	...	YULETIDE IMAGES GIFT WRAP SET	ZINC HEART T-LIGHT HOLDER	ZINC STAR T-LIGHT HOLDER	ZINC BOX SIGN HOME	ZINC FOLKART SLEIGH BELLS	ZINC HEART LATTICE T-LIGHT HOLDER
InvoiceNo																	
536527	0	0	0	0	0	0	0	0	0	0	...	0	0	0	0	0	0
536840	0	0	0	0	0	0	0	0	0	0	...	0	0	0	0	0	0
536861	0	0	0	0	0	0	0	0	0	0	...	0	0	0	0	0	0
536967	0	0	0	0	0	0	0	0	0	0	...	0	0	0	0	0	0
536983	0	0	0	0	0	0	0	0	0	0	...	0	0	0	0	0	0
...
581266	0	0	0	0	0	0	0	0	0	0	...	0	0	0	0	0	0
581494	0	0	0	0	0	0	0	0	0	0	...	0	0	0	0	0	0
581570	0	0	0	0	0	0	0	0	0	0	...	0	0	0	0	0	0
581574	0	0	0	0	0	0	0	0	0	0	...	0	0	0	0	0	0
581578	0	0	0	0	0	0	0	0	0	0	...	0	0	0	0	0	0

457 rows x 1694 columns

Algoritmo APRIORI

In [27]:	<pre>itemsets_mais_frequentes_germany = apriori(carrinho_compras_germany, min_support=0.05, use_colnames=True)</pre>
----------	--

Suportes para cada produto

In [28]:	<pre>itemsets_mais_frequentes_germany</pre>
----------	---