

Notebook Jupyter 2_geraWordclouds

Geração das WordClouds das descrições

Os dados gerados após o tratamento, objeto de estudo, estão contidos em um campo de descrição, onde o importador ou seu representante faz a entrada livre de dados, sem qualquer exigência ou padronização, nossa análise exploratória foi convincente no sentido de que deveria ser feito um tratamento com a utilização de Programação de Linguagem Natural – PLN. Dentro dessa análise exploratória de dados, observou-se a ocorrência de palavras que, embora frequentes, serão irrelevantes para a determinação da diferenciação entre os itens que estão classificados dentro do agrupamento representado pela NCM em estudo. Comum no estudo de categorização de textos, o uso de stopwords é recomendado nesse caso, pois da mesma forma que artigos, pronomes e outras palavras tradicionalmente identificadas como não relevantes para a distinção, termos como kit, transmissão, corrente, coroa e pinhão também são igualmente irrelevantes para distinguir um item de outro. Na biblioteca Natural Language Processing Toolkit - nltk já existe uma lista de stopwords para a língua portuguesa. Como essa biblioteca já possui a funcionalidade de complementação dessa lista, adicionou-se os itens kit, transmissão, corrente, coroa e pinhão; termos que pela sua característica não distinguem os itens dentro do dataset estudado.

Importa os dados já tratados

```
In [1]: import pandas as pd
import time
import matplotlib.pyplot as plt
from wordcloud import WordCloud
```

```
In [2]: # Data e hora da execução do script
initot=time.time()
print(f'Código executado em {time.strftime("%d/%m/%Y às %H:%M", time.localtime(time.time()))}')

Código executado em 25/01/2022 às 11:25
```

```
In [3]: # Importa base de dados para um dataframe
df = pd.read_excel(r'./bases/dataframe.xlsx')
```

```
In [4]: # Verifica o tamanho do dataframe
df.shape
```

```
Out[4]: (18276, 3)
```

```
In [5]: # Mostra Linhas de exemplo do dataframe
df.sample(5)
```

Out[5]:

	PAIS DE ORIGEM	DESCRICAO DO PRODUTO	VALOR UN.PROD.DOLAR
3109	CHINA, REPUBLICA POP	KIT DE TRANSMISSAO , MARCA RIFFEL, TITANIUM (1...	6.708000
10980	CHINA, REPUBLICA POP	Y10821/i45" KIT DE TRANSMISSÃO, EM AÇO 1045, C...	2.950000
6336	CHINA, REPUBLICA POP	878193 - KIT DE TRANSMISSÃO, COMPOSTO DE CORRE...	23.945997
14741	CHINA, REPUBLICA POP	1104970 - DAROM / Kit de transmissao 428 - 110...	4.260000
1734	CHINA, REPUBLICA POP	80350 KIT DE TRANSMISSÃO, COMPOSTO DE CORRENTE...	5.925957

Importa as stopwords da língua portuguesa

```
In [6]: # Importar Lista de Stopwords
from nltk.corpus import stopwords
stopwords = set(stopwords.words('portuguese'))
```

```
In [7]: # Mostra tamanho da lista de stopwords
len(stopwords)
```

Out[7]: 204

```
In [8]: # Mostra toda a lista de stopwords
swtemp = list(stopwords)
swtemp.sort()
print(swtemp)
```

['a', 'ao', 'aos', 'aquela', 'aquelas', 'aquele', 'aqueles', 'aquilo', 'as', 'at é', 'com', 'como', 'da', 'das', 'de', 'dela', 'delas', 'dele', 'deles', 'depois', 'do', 'dos', 'e', 'ela', 'elas', 'ele', 'eles', 'em', 'entre', 'era', 'eram', 'essa', 'essas', 'esse', 'esses', 'esta', 'estamos', 'estas', 'estava', 'estavam', 'este', 'esteja', 'estejam', 'estejamos', 'estes', 'esteve', 'estive', 'estivemos', 'estiver', 'estivera', 'estiveram', 'estiverem', 'estivermos', 'estivesse', 'estivessem', 'estivéramos', 'estivéssemos', 'estou', 'está', 'estávamos', 'estão', 'eu', 'foi', 'fomos', 'for', 'fora', 'foram', 'forem', 'formos', 'fosse', 'fossem', 'fui', 'fôramos', 'fôssemos', 'haja', 'hajam', 'hajamos', 'hавemos', 'hei', 'houve', 'houvemos', 'houver', 'houvera', 'houveram', 'houverei', 'houverem', 'houveremos', 'houveria', 'houveriam', 'houvermos', 'houverá', 'houverão', 'houveríamos', 'houvesse', 'houvessem', 'houvéramos', 'houvéssemos', 'há', 'hão', 'isso', 'isto', 'já', 'lhe', 'lhes', 'mais', 'mas', 'me', 'mesmo', 'meu', 'meus', 'minha', 'minhas', 'muito', 'na', 'nas', 'nem', 'no', 'nos', 'nossa', 'nossas', 'nosso', 'nossos', 'num', 'numa', 'não', 'nós', 'o', 'os', 'ou', 'para', 'pela', 'pelas', 'pelo', 'pelos', 'por', 'qual', 'quando', 'que', 'quem', 'se', 'seja', 'sejam', 'sejamos', 'sem', 'serei', 'seremos', 'seria', 'seriam', 'será', 'serão', 'seríamos', 'seu', 'seus', 'somos', 'sou', 'sua', 'suas', 'são', 'só', 'também', 'te', 'tem', 'temos', 'tenha', 'tenham', 'tenhamos', 'tenho', 'terei', 'teremos', 'teria', 'teriam', 'terá', 'terão', 'teríamos', 'teu', 'teus', 'teve', 'tinha', 'tinham', 'tive', 'tivemos', 'tiver', 'tivera', 'tiveram', 'tiverem', 'tivermos', 'tivesse', 'tivessem', 'tivéramos', 'tivéssemos', 'tu', 'tua', 'tuas', 'tém', 'tínhamos', 'um', 'uma', 'você', 'vocês', 'vos', 'à', 'às', 'é', 'éramos']

Cria a string contendo todas as descrições

```
In [9]: # Mostra algumas descrições do dataset
df['DESCRICAO DO PRODUTO'].sample(5)
```

```
Out[9]: 12040 KIT TRANSMISSÃO COROA,CORRENTE E PINHAO PARA M...
16723 KIT DE TRANSMISSAO P/MOTOCICLETA MOD: XLR 125;...
16761 KIT DE TRANSMISSAO, MARCA RIFFEL, TITANIUM (10...
3674 item .12;Partes e peças para Motocicletas,Dest...
11420 980727 KIT DE TRANSMISSÃO, COMPOSTO DE CORRENT...
Name: DESCRICAO DO PRODUTO, dtype: object
```

```
In [10]: # Mescla todas as descrições como uma string usando espaço como separador
descricoes = " ".join(df['DESCRICAO DO PRODUTO']).lower()
```

Gera a WodCloud aplicando o filtro das stopwords

```
In [11]: # Define e gera a wordcloud para um máximo de 400 palavras de tamanho mínimo 2, sem termos duplos
```

```
wordcloud = WordCloud(stopwords=stopwords,  
                       background_color="white",  
                       width=1600, height=800,  
                       max_words=400,  
                       min_word_length=2,  
                       collocations=False,  
                       include_numbers=True).generate(descricoes)
```

```
In [12]: # Exibe a imagem da WordCloud gerada
fig, ax = plt.subplots(figsize=(20,8))
ax.imshow(wordcloud, interpolation='bilinear')
ax.set_axis_off()
plt.imshow(wordcloud)
```

```
Out[12]: <matplotlib.image.AxesImage at 0x266cd0df3c8>
```



```
In [13]: # Exporta para um arquivo
wordcloud.to_file(r"./imagens/wordcloud descricoes antes.png")
```

```
Out[13]: <wordcloud.wordcloud.WordCloud at 0x266cd169eb8>
```

Atualiza as stopwords contendo as palavras irrelevantes

Palavras a adicionar, tais como: kit, transmissao, transmissão, coroa, pinhão, etc. que não agregam nenhuma diferença aos itens da lista

```
In [14]: # Palavras a adicionar na lista de stopwords estão contidas em um arquivo csv externo
dfsw = pd.read_csv('./bases/stopwords.csv', encoding='ISO-8859-1')
stopwords_df=sorted(list(dfsw['stopword']))
swtemp = list(stopwords_df)
swtemp.sort()
print(swtemp)
```

```
['abaixo', 'acessorios', 'acessórios', 'aco', 'acondicionados', 'adaptavel', 'adaptável', 'allen', 'almas', 'alta', 'am', 'anel', 'ano', 'aplicacao', 'application', 'aplicavel', 'aplicação', 'aplicável', 'application', 'ate', 'atitanium', 'aç', 'aço', 'bicicleta', 'bicycle', 'bike', 'bravo', 'cada', 'caixa', 'caixas', 'cambio', 'carbono', 'certificado', 'cever', 'chain', 'chh', 'china', 'ciclomotor', 'ciclomotores', 'cilindrada', 'cilindradas', 'cod', 'code', 'codigo', 'comando', 'combustão', 'comercial', 'comercialmente', 'commodity', 'compativel', 'compatível', 'comp', 'completo', 'completos', 'composto', 'composto', 'compostopor', 'compostpo', 'comum', 'condicao', 'condicoes', 'condição', 'condições', 'confeccionado', 'conformidade', 'conhecido', 'conj', 'conjunto', 'conjuntos', 'constituído', 'constitutivo', 'constituído', 'contendo', 'coposto', 'coroa', 'coroaes', 'corr', 'corrent', 'corrente', 'correntee', 'correntes', 'corrnte', 'cx', 'câmbio', 'câmbio', 'código', 'decreto', 'denominada', 'dente', 'dentes', 'descricao', 'descrição', 'descriçao', 'descrição', 'destaque', 'destaque', 'destaques', 'detransmissão', 'diante', 'dimensao', 'dimensoes', 'dimensão', 'dimensões', 'diverso', 'diversos', 'dominado', 'durabilidade', 'elo', 'elos', 'embalagem', 'engine', 'engrenagem', 'engrenagens', 'epinhao', 'epinhão', 'especifico', 'específico', 'espessura', 'evol', 'exclusivo', 'fabr', 'fabri', 'fabricada', 'fabricado', 'final', 'foiproduzida', 'formado', 'funcao', 'funcao', 'função', 'função', 'função', 'função', 'gtin', 'hardt', 'ho', 'hp', 'imetro', 'in', 'incluindo', 'incluso', 'indicado', 'ingles', 'inmetro', 'inv', 'invoice', 'iron', 'item', 'jc', 'ki', 'kif', 'kit', 'kitr', 'kittr', 'kittr', 'kmc', 'ligacoes', 'ligações', 'ligações', 'ligações', 'marca', 'mark', 'match', 'material', 'materialdo', 'maxx', 'medida', 'medidas', 'medindo', 'metal', 'mini', 'mod', 'modelo', 'modelos', 'moto', 'motociclet', 'motocicleta', 'motocicletas', 'motoneta', 'motonetas', 'motoparts', 'motor', 'motorcicleta', 'motos', 'motos', 'movimento', 'nbsp', 'ncm', 'nome', 'nopinh', 'normais', 'nova', 'novo', 'nr', 'numero', 'número', 'onde', 'or', 'origem', 'oring', 'ox', 'papelao', 'papelão', 'part', 'parte', 'partes', 'parts', 'pc', 'pc-coroa', 'pc-corrente', 'pc-pinhao', 'pcs', 'pec', 'pecas', 'perfil', 'peç', 'peças', 'pinh', 'pinhao', 'pinhão', 'posição', 'premium', 'procedencia', 'procedência', 'prodepe', 'produto', 'próprio', 'pç', 'qdes', 'qtd', 'qtas', 'qty', 'quadriciclo', 'quadriciclos', 'quantidade', 're', 'ref', 'reforcada', 'reforçada', 'registro', 'rel', 'relacao', 'relação', 'reposicao', 'resp', 'respo', 'respons', 'responso', 'responsav', 'responsave', 'responsavel', 'responsáv', 'responsável', 'ret', 'retalho', 'retentor', 'riffel', 'ring', 'roda', 'sae', 'scud', 'semi', 'semi-', 'semi-kit', 'sendo', 'serve', 'set', 'shipping', 'sistema', 'sm', 'sprocket', 'standard', 'standart', 'standartt', 'std', 'stdmodelo', 'steel', 'tambem', 'também', 'tec', 'temp', 'temperado', 'tipo', 'tipos', 'titanio', 'titaniu', 'titanium', 'titaniun', 'tr', 'tracao', 'tracão', 'trans', 'transmis', 'transmisao', 'transmiss', 'transmissa', 'transmissao', 'transmission', 'transmissão', 'transmitir', 'traseira', 'tração', 'tração', 'und', 'unds', 'unid', 'unidade', 'unidade', 'unidades', 'unifort', 'uo', 'uso', 'utilizada', 'utilizadas', 'utilizado', 'utilizados', 'utilização', 'vem', 'venda', 'with', 'xy', 'year']
```

```
In [15]: # Atualizar stopwords
stopwords.update(stopwords_df)
```

In [16]: `# Mostra toda a Lista de stopwords`

```
swtemp = list(stopwords)
swtemp.sort()
print(swtemp)
```

['a', 'abaixo', 'acessorios', 'acessórios', 'aco', 'acondicionados', 'adaptavel', 'adaptável', 'allen', 'almas', 'alta', 'am', 'anel', 'ano', 'ao', 'aos', 'aplicacao', 'aplicação', 'aplicavel', 'aplicável', 'application', 'aquela', 'aquelas', 'aquele', 'aqueles', 'aquilo', 'as', 'ate', 'atitanium', 'até', 'aç', 'aço', 'bicicleta', 'bicycle', 'bike', 'bravo', 'cada', 'caixa', 'caixas', 'cambio', 'carbono', 'certificado', 'cever', 'chain', 'chh', 'china', 'ciclomotor', 'ciclomotores', 'cilindrada', 'cilindradas', 'cod', 'code', 'codigo', 'com', 'comand', 'combustão', 'comercial', 'comercialmente', 'commodity', 'como', 'compativel', 'compatível', 'compl', 'completo', 'completos', 'composto', 'compostopor', 'composto', 'comum', 'condicao', 'condicoes', 'condição', 'condições', 'confeccionado', 'conformidade', 'conhecido', 'conj', 'conjunto', 'conjuntos', 'constituído', 'constitutivo', 'constituído', 'contendo', 'coposto', 'coroa', 'coroaes', 'corr', 'corrent', 'corrente', 'correntee', 'correntes', 'corrente', 'cx', 'câmbio', 'câmbio', 'código', 'da', 'das', 'de', 'decreto', 'dela', 'delas', 'dele', 'deles', 'denominada', 'dente', 'dentes', 'depois', 'descricao', 'descrição', 'descrição', 'descrição', 'destaque', 'destaques', 'detransmissão', 'diante', 'dimensao', 'dimensoes', 'dimensão', 'dimensões', 'diverso', 'diversos', 'do', 'dominado', 'dos', 'durabilidade', 'e', 'ela', 'elas', 'ele', 'eles', 'elo', 'elos', 'em', 'embalagem', 'engine', 'engrenagem', 'engrenagens', 'entre', 'epinhao', 'epinhão', 'era', 'eram', 'especifico', 'específico', 'espessura', 'essa', 'essas', 'esse', 'esses', 'esta', 'estamos', 'estas', 'estava', 'estavam', 'este', 'esteja', 'estejam', 'estejamos', 'estes', 'estive', 'estive', 'estivemos', 'estiver', 'estivera', 'estiveram', 'estiverem', 'estivermos', 'estivesse', 'estivessem', 'estivérámos', 'estivéssemos', 'estou', 'está', 'estávamos', 'estão', 'eu', 'evol', 'exclusivo', 'fabr', 'fabri', 'fabricada', 'fabricado', 'final', 'foi', 'foiproduzida', 'fomos', 'for', 'fora', 'foram', 'forem', 'formado', 'formos', 'fosse', 'fossem', 'fui', 'funcao', 'função', 'função', 'função', 'fôramos', 'fôssemos', 'gtin', 'haja', 'hajam', 'hajamos', 'hardt', 'hавemos', 'hei', 'ho', 'houve', 'houvermos', 'houver', 'houvera', 'houveram', 'houverei', 'houverem', 'houveremos', 'houveria', 'houveriam', 'houvermos', 'houverá', 'houverão', 'houveríamos', 'houvesse', 'houvessem', 'houvéramos', 'houvéssemos', 'hp', 'há', 'hão', 'imetro', 'in', 'incluindo', 'incluso', 'indicado', 'ingles', 'inmetro', 'inv', 'invoice', 'iron', 'isso', 'isto', 'item', 'jc', 'já', 'ki', 'kif', 'kit', 'kitr', 'kittr', 'kmc', 'lhe', 'lhes', 'ligacoes', 'ligações', 'ligações', 'mais', 'marca', 'mark', 'mas', 'match', 'material', 'materialdo', 'maxx', 'me', 'medida', 'medidas', 'medindo', 'mesmo', 'metal', 'meu', 'meus', 'minha', 'minhas', 'mini', 'mod', 'modelo', 'modelos', 'moto', 'motociclet', 'motocicleta', 'motocicletas', 'motoneta', 'motonetas', 'motoparts', 'motor', 'motorcicleta', 'motos', 'movimento', 'muito', 'na', 'nas', 'nbsp', 'ncm', 'nem', 'no', 'nome', 'nopinh', 'normais', 'nos', 'nossa', 'nossas', 'nosso', 'nossos', 'nova', 'novos', 'nr', 'num', 'numa', 'numero', 'não', 'nós', 'número', 'o', 'onde', 'or', 'origem', 'oring', 'os', 'ou', 'ox', 'papela', 'papela', 'papela', 'para', 'part', 'parte', 'partes', 'parts', 'pc', 'pc-coroa', 'pc-corrente', 'pc-pinhao', 'pcs', 'pec', 'peças', 'pela', 'pelas', 'pelo', 'pelos', 'perfil', 'peç', 'peças', 'pinh', 'pinhao', 'pinhão', 'por', 'posição', 'premium', 'procedencia', 'procedência', 'prodepe', 'produto', 'próprio', 'pç', 'qdes', 'qtd', 'qtds', 'qty', 'quadriciclo', 'quadriciclos', 'qual', 'quando', 'quantidade', 'quantidade', 'que', 'quem', 're', 'ref', 'reforcada', 'reforcada', 'registro', 'rel', 'relacao', 'relação', 'reposicao', 'resp', 'respo', 'respos', 'responso', 'responsav', 'responsave', 'responsavel', 'responsável', 'ret', 'retalho', 'retentor', 'riffel', 'ring', 'roda', 'sa', 'scud', 'se', 'seja', 'sejam', 'sejamos', 'sem', 'semi', 'semi-', 'semi-kit', 'sendo', 'serei', 'seremos', 'seria', 'seriam', 'serve', 'será', 'serão', 'seríamos', 'set', 'seu', 'seus', 'shipping', 'sistema', 'sm', 'somos', 'sou', 'sprocket', 'standard', 'standart', 'standartt', 'std', 'stdmodelo', 'steel', 'sua', 'suas', 'são', 'só', 'tambem', 'também', 'te', 'tec', 'tem', 'temos', 'temp', 'temperado', 'tenha', 'tenham', 'tenhamos', 'tenho', 'terei', 'teremos', 'teria', 'teriam', 'terá', 'terão', 'teríamos', 'teu', 'teus', 'teve', 'tinha', 'tinham', 'tipo', 'tipos', 'titania', 'titaniu', 'titanium', 'titaniun', 'tive', 'tivemos', 'tiver', 'tivera', 'tiveram', 'tiverem', 'tivermos', 'tivesse', 'tivessem', 'tivérámos', 'tivéssemos', 'tr', 'tracao', 'tracão', 'trans', 'transmis', 'transmisao', 'transmiss', 'transmissa', 'transmissao', 'transmission', 'transmissão', 'transmitir', 'traseira', 'tração', 'tração', 'tu', 'tua', 'tuas', 'tém', 'tínhamos', 'um', 'uma', 'unds', 'unid', 'unidade', 'unidades', 'unifort', 'uo', 'uso', 'utilizada', 'utilizadas', 'utilizado', 'utilizados', 'utilização', 'vem', 'venda', 'você', 'você', 'vos', 'with', 'xy', 'year', 'à', 'às', 'é', 'éramos']

Gera a WodCloud aplicando o filtro das stopwords atualizado

[illegible]

```
In [18]: # Exibe a imagem da nova WordCloud gerada
fig, ax = plt.subplots(figsize=(20,8))
ax.imshow(wordcloud, interpolation='bilinear')
ax.set_axis_off()
plt.imshow(wordcloud)
```

```
Out[18]: <matplotlib.image.AxesImage at 0x266ccdf9080>
```



```
In [19]: # Exporta para um arquivo
wordcloud.to_file("./imagens/wordcloud_descricoes_depois.png")
```

```
Out[19]: <wordcloud.wordcloud.WordCloud at 0x266ccdd27b8>
```

```
In [20]: tempotot=time.time()-initot
if tempotot>60:
    print(f'Tempo total de execução: {tempotot/60:.2f} minutos.')
else:
    print(f'Tempo total de execução: {tempotot:.2f} segundos.')
```

Tempo total de execução: 13.05 segundos.