

```

1 # funcoesTCC.py
2 #
3 # !/usr/bin/env python
4 # coding: utf-8
5
6 # # Classificação dos modelos de motocicleta a partir da descrição do produto
7
8 import pandas as pd, numpy as np
9 import re
10 import pickle
11
12 # ### Importa as stopwords da língua portuguesa
13 # Importar lista de Stopwords
14 from nltk.corpus import stopwords
15 stopwords = set(stopwords.words('portuguese'))
16 # Palavras a adicionar na lista de stopwords estão contidas em um arquivo csv
17 # externo
18 dfsw = pd.read_csv('./bases/stopwords.csv', encoding='ISO-8859-1')
19 stopwords_df=sorted(list(dfsw['stopword']))
20 # Atualizar stopwords
21 stopwords.update(stopwords_df)
22
23 # ### Carrega lista de aplicações e cria palavrasChave
24 # carrega a lista de marcas de motos do arquivo
25 dfaplicacoes=pd.read_csv('./bases/Aplicacoes.csv')
26 # remove caracteres especiais ou soltos e termos duplicados, salvando na lista
27 palavrasChave=sorted(set(re.sub(r"\b \w \b",
28                                ' ',
29                                re.sub(r"[<>()|\\+\\$%&#@\\'\\\" ]+",
30                                        " ",
31                                        " ").join(dfaplicacoes['APLICACOES'].tolist()))).split()))
32
33 # ##### Função de limpeza de dados irrelevantes para a classificação e remoção de
34 # stopwords
35 def limpaDescricao(descricao): #
36     descricao=descricao.lower() #transformar em minúsculas
37     # remove top (1045) e variantes
38     descricao=re.sub(r'\b[ (-]*top \(\s*(1045\s*)\)[- ]*\b',' ',descricao)
39     # remove códigos numéricos entre parênteses com -*/
40     descricao=re.sub(r"\(\s*\d*\[\\/*\s*\d*\s*\)", ' ', descricao)
41     # remove a ocorrência de "código e etc." e o termo seguinte começado com
42     # número
43     # att: (alguns tem hífen ou asterisco) (colocar antes de remover pontuação)
44     descricao=re.sub(r"\b(invoice|código|codigo|cod|cód|(certificado|cert)(no|nr)|ref)[0-9a-z/\s*\.\:]*\s*\d+[^\s,;)]", ' ', descricao)
45     # remove identificação de referência de engrenagens dos kits (antes da
46     # pontuação)
47     descricao=re.sub(r"([^\s,;)])(ho|uo|h|l|t|ktd|sm|m|d| x|elos )\d{1,}[\s,;)]", ' ', descricao)
48     descricao=re.sub(r"\d*(ho|uo|h|l|t|ktd|sm|m|d|elos )\d{1,}[\s,;)]", ' ', descricao)
49     descricao=re.sub(r"\d*(ho|uo|h|l|z|t|ktd|m|d| dentes| elos)", ' ', descricao)
50     # substitui os termos "s/re" e "s/ret" por "sem retentor"
51     descricao=re.sub(r"\b(s/re|s/ret)\b", 'sem retentor', descricao)
52     # substitui os termos "c/re" ou "c/ret" por "com retentor"
53     descricao=re.sub(r"\b(c/re|c/ret)\b", 'com retentor', descricao)
54     # substitui o termo "aplicação" e "modelo" emendado com outro

```

```

50     descricao=re.sub(r"aplicacao", "aplicacao ", descricao)
51     descricao=re.sub(r"modelo", "modelo ", descricao)
52     # remove códigos no início da descrição
53     descricao=re.sub(r"^\b\d{2,}[^]*\b", ' ', descricao)
54     descricao=re.sub(r"^\k{^}+", ' ', descricao)
55     descricao=re.sub(r"- | -|[\\\+,.;!/?/]+", ' ', descricao) #remover pontuação
(att: "- " ou " -")
56     #correção de erros de digitação comuns
57     termos={'titan': ['titian', 'tita', 'tintan', 'tit'],
58             'honda': ['hond', 'hnda', 'hon'],
59             'twister': ['twist', 'twiste'],
60             'dafra kansas': ['dafra kan'],
61             'tenere': ['tener', 'tenerre'],
62             'broz': ['bros', 'bross'],
63             'titan 150': ['titan150'],
64             'broz 150':
['bross125.', 'bros125.', 'broz125', 'bross150.', 'bros150.', 'broz150'],
65             'pop 100': ['pop100'],
66             'phoenix': ['phoeni', 'phenix'],
67             'c100': ['c 100']}
68     for termo in termos:
69         for termoerrado in termos[termo]:
70             descricao=re.sub(r"\b"+termoerrado+"\b", termo, descricao)
71             descricao=re.sub(r"/<>()|+\\$%&#@'\""+", ' ', descricao) #remover
carcteres especiais
72             # remove a ocorrência de medidas tipo 00x000x00 ou 000x0000
73             descricao=re.sub(r"\b\d{1,}(x|\\*)\d{1,}(x|\\*)\d{1,}|\\d{1,}(x|\\*)\\d{1,}\\b", '
', descricao)
74             # remove identificação de quantidades, unidades, peças e conjuntos
75             descricao=re.sub(r"\b(\\d* *(conj|und|uni|pc|pç|pec|peç)( \\w|\\w)+?)\\b", ' ',
descricao)
76             # remove identificação de mais de 4 dígitos com ou sem letras no início e no
final
77             descricao=re.sub(r"\\w+\\d{4,}\\w+", ' ', descricao)
78             # remove números de 4 dígitos ou mais começados de 2 a 9
79             descricao=re.sub(r"\\b[02-9]\\d{3,}\\b", ' ', descricao)
80             # remove identificação de termos começados por zero
81             descricao=re.sub(r"\\b0\\d*\\w+(?=\\b)", ' ', descricao)
82             # remove a ocorrência de "marca " e o termo na lista até o próximo espaço
83             for marca in ['kmc *gold', 'am *gold', 'king', 'bravo *racing', 'riffel *top']:
84                 descricao=re.sub(r"\\bmarca[ :\\./*]+str(marca)+r"^[^]*", ' ', descricao)
# colocar antes das stopwords
85             descricao=re.sub(r"marca[ :\\./*]*\\w+", ' ', descricao)
86             descricao=re.sub(r"^(^|-|-)", ' ', descricao)
87             # remove stopwords mantendo a ordem original da descrição
88             descricao=list(dict.fromkeys(descricao.split())) # cria lista com termos
únicos
89             descricao=" ".join([x for x in descricao if x not in set(stopwords)]) #
exclui stopwords
90             # limpa os número que não estão na lista de aplicações (colocar depois das
stopwords)
91             desc=descricao.upper().split() # quebra a descrição
92             dif=list(set(descricao.upper().split()).difference(palavrasChave)) # pega os
termos diferentes de palavrasChave
93             [desc.remove(x) for x in desc if (x in dif and x.isnumeric())] # exclui de
desc os termos numéricos diferentes
94             descricao=" ".join(desc).lower() # volta para texto
95             #remover hífen, letras ou números soltos (deixar duplicado mesmo)
96             descricao=re.sub(r"^(^|-|-|\\b\\w\\b)", ' ', descricao)

```

```

97     descricao=re.sub(r"^(^|-|-|\b\w\b)", ' ', descricao)
98     #substitui remove o i das cilindradas: ex.: 125i por 125
99     termos=re.findall(r"\d{1,}i\b",descricao)
100     if termos:
101         for termo in termos:descricao=descricao.replace(termo,termo[:-1])
102     # remove espaços em excesso (colocar no final)
103     descricao=re.sub(r" {2,}", ' ', descricao)
104     descricao=descricao.strip()
105     # retorna a descricao como saída da função
106     return descricao # retorna a descrição
107
108 def achaPalavraChave(descricao):
109     palavras=[]
110     descricao=descricao.upper()
111     desc=descricao.split()
112     for palavra in palavrasChave:
113         if palavra in desc:
114             palavras.append(palavra)
115         else:
116             if palavra.isnumeric():
117                 pat=r"[0-9]*"+str(palavra)+r"[0-9]*"
118             elif palavra.isalpha():
119                 pat=r"[A-Z]*"+str(palavra)+r"[A-Z]*"
120             else:
121                 pat=r"\b"+palavra+r"\b"
122             a = re.findall(pat,descricao)
123             if len(a)>0:
124                 # adiciona resultado nas palavras se o resultado estiver em
palavrasChave
125                 palavras+=[a[i] for i in range(len(a)) if a[i] in palavrasChave]
126                 palavras=list(set(palavras)) # remove duplicados
127                 palavras=" ".join(palavras) # converte para string
128                 return palavras.lower()
129
130 # termos que iniciam item da descrição correspondem a marca
131 # As que começam com espaço devem permanecer assim, pois há outros modelos com o
mesmo final
132 Marcas = {'HONDA': ['CG', 'CD', 'CBX', 'CB', 'CBR', 'CRF', 'BIZ', 'BROS', 'BROZ', 'XL', '
FAN', 'XR', 'XRE'
133             'DREAM', 'TITAN', 'TODAY', 'TWIN', 'POP', 'NX', 'NXR', 'TWISTER',
'HORNET',
134             'AMERICA', 'BOLDOR', 'DUTY', 'FIREBLADE', 'FURY', 'WING', 'LEAD', 'MAGNA', 'NL',
'NC', 'NSR', 'NC', 'NXR', 'PACIFIC', 'COAST', 'SHADOW', '
135 STRADA', 'STUNNER', 'HAWK',
136             'SUPERBLACKBIRD', 'TORNADO', 'TURUNA', 'XRV', 'AFRICA', 'VALKYRIE', 'VARADERO',
'VFR', 'VLR', 'VTR', 'VTX', 'TRANSALP'],
137             'YAMAHA': ['AEROX', 'ALBA', 'AXIS', 'BWS', 'DRAG ', 'DT', 'FZ', 'FJ', '
RD', 'TENERE',
138             'MT', 'XF', 'XJ', 'XS', 'XT', 'XZ', 'YF', 'YZ', 'LANDER', 'GLADIATOR', 'GRIZZLY',
'YBR', 'YZ', 'VIRAGO', 'FACTOR', 'EC', 'CRYPTON', 'FAZER', 'JOG', '
139 LANDER',
140             'FROG', 'LIBERO', 'MAJESTY', 'MEST', 'MIDNIGHT', 'MORPH', 'NEO', 'PASSOL'],
141             'DAFRA':
142             ['APACHE', 'CITYCOM', 'KANSAS', 'LASER', 'NEXT', 'RIVA', 'ROADWIN', 'ZIG', 'SPEED'],
143             'SUZUKI': ['KATANA', 'YES', 'INTRUDER'],

```

```

144         'ZONGSHEN': ['ZS'],
145         'KASINSKI': ['COMET', 'MIRAGE'],
146         'POLARIS': ['SPORTSMAN', 'RZR', 'RANGER'],
147         'KAWASAKI':
148     ['NINJA', 'VERSYS', 'VOYAGER', 'GTR', 'KDX', 'KL', 'KX', 'KZ', 'ZR', 'ZZ', 'ER6N', 'ER6F'],
149     'DAYANG': ['DY1', 'DY2', 'DY5'],
150     'SUNDOWN': ['WEB', 'FIFITY', 'PALIO', 'PGO', 'STX', 'VBLADE', 'EVO', 'HUNTER
151     MAX'],
152     'SHINERAY':
153     ['BIKE', 'BRAVO', 'DISCOVER', 'EAGLE', 'INDIANAPOLIS', 'JET', 'NEW', 'WAVE',
154     'STRONG', 'SUPER SMART', 'VENICE', 'XY']}
155
156 # Função para pegar a chave pelo valor, dado que valor é único.
157 def pegaChave(v, dict):
158     for chave, valores in dict.items():
159         if type(valores) != type([1, 2]):
160             valores = [valores]
161         for valor in valores:
162             if v == valor:
163                 return chave
164     return "Não existe chave para esse valor."
165
166 def acrescentaMarca(descricao):
167     for marca in Marcas:
168         if re.search(marca, descricao.upper()):
169             descricao += " " + marca
170         for termo in Marcas[marca]:
171             t1 = termo.split()
172             if len(t1) > 1:
173                 pat = r"(?:" + t1[0] + r"|" + t1[1] + r").*(?:" + t1[0] + r"|" + t1[1] + r")"
174             elif len(termo) < 3:
175                 pat = termo + r"([0-9]{1,}|\b)"
176             else:
177                 pat = termo
178             resultados = re.findall(pat, descricao.upper())
179             if resultados:
180                 descricao += " " + marca
181                 descricao += " " + " ".join(resultados)
182                 descricao += " " + termo
183     descricao = " ".join(sorted(set(descricao.lower().split())))
184     return descricao
185
186 # ### Função final que transforma a DESCRICAO DO PRODUTO em Modelo para
187 classificar
188 def criaModelo(descricao):
189     descricao = limpaDescricao(descricao)
190     descricao = achaPalavraChave(descricao)
191     descricao = acrescentaMarca(descricao)
192     return descricao
193
194 # ### Função que determina a existência de retentor no kit e retorna True|False
195 def retentorAux(descricao):
196     # define o padrão de busca
197     padrao = r'c/ *ret|com *ret' #r"(com|c/) *ret"
198     descricao = descricao.lower()
199     busca = re.findall(padrao, descricao)
200     if busca:
201         descricao = busca[0]

```

```
198         return True
199     else:
200         descricao=''
201         return False
202
203 # ## Carrega o modelo Linear SVC
204 # ### Carrega arquivos do pickle
205 with open(r'./pickle/clfsvc.pkl', 'rb') as file:
206     clfsvc = pickle.load(file)
207     file.close()
208 with open(r'./pickle/cvt.pkl', 'rb') as file:
209     cvt = pickle.load(file)
210     file.close()
211
212 # ### Define a função de classificação
213 def classificaAplicacao(descricao):
214     modelo = criaModelo(descricao)
215     novo_cvt = cvt.transform(pd.Series(modelo))
216     aplicacao = clfsvc.predict(novo_cvt)[0]
217     return aplicacao
```