

```

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, os
11
12 # ### Importa as stopwords da língua portuguesa
13 # Importar lista de Stopwords
14 import nltk
15 from nltk.corpus import stopwords
16 nltk.download('stopwords')
17 stopwords = set(stopwords.words('portuguese'))
18 # Palavras a adicionar na lista de stopwords estão contidas em um arquivo csv
19 # externo
20 dfsw = pd.read_csv('./bases/stopwords.csv', encoding='ISO-8859-1')
21 stopwords_df=sorted(list(dfsw['stopword']))
22 # Atualizar stopwords
23 stopwords.update(stopwords_df)
24
25 # ### Carrega lista de aplicações e cria palavrasChave
26 # carrega a lista de marcas de motos do arquivo
27 dfaplicacoes=pd.read_csv('./bases/Aplicacoes.csv')
28 # remove caracteres especiais ou soltos e termos duplicados, salvando na lista
29 palavrasChave=sorted(set(re.sub(r"\b \w \b",
30                                ' ',
31                                re.sub(r"[<>()|\\+\\$%&#@'\""]+",
32                                " ",
33                                " ").join(dfaplicacoes['APLICACOES'].tolist()))).split()))
34
35 # ##### Função de limpeza de dados irrelevantes para a classificação e remoção de
36 # stopwords
37 def limpaDescricao(descricao): #
38     descricao=descricao.lower() #transformar em minúsculas
39     # remove top (1045) e variantes
40     descricao=re.sub(r'\b[ (-]*top \(\s*(1045\s*)\)[- ]*\b',' ',descricao)
41     # remove códigos numéricos entre parênteses com -*/
42     descricao=re.sub(r'\(\s*\d*\[\\/*\s*\d*\s*\)', ' ', descricao)
43     # remove a ocorrência de "código e etc." e o termo seguinte começado com
44     # número
45     # att: (alguns tem hífen ou asterisco) (colocar antes de remover pontuação)
46     descricao=re.sub(r"\b(invoice|código|codigo|cod|cód|(certificado|cert)(
47     no|nr|)|ref)[0-9a-z/\s*\.\:]* *\d[^\s]+", ' ', descricao)
48     # remove identificação de referência de engrenagens dos kits (antes da
49     # pontuação)
50     descricao=re.sub(r"([\^a-z]|)(ho|uo|h|l|t|ktd|sm|m|d| x|elos )\d{1,}[ \s-
51     \/,;].|[ \s-\\/(]*\d{1,}(ho|uo|h|l|z|t|ktd|m|d|x| dentes| elos)[ \s-\\/,;]", ' ',
52     descricao) # 00h
53     descricao=re.sub(r"\d*(ho|uo|h|l|t|ktd|sm|m|d|elos )\d{1,}[ \s-\\/,;].|[ \s-
54     \\/(]*\d{1,}(ho|uo|h|l|z|t|ktd|m|d| dentes| elos)", ' ', descricao) # 00h
55     # substitui os termos "s/re" e "s/ret" por "sem retentor"
56     descricao=re.sub(r"\b(s\\re|s\\ret)\b", 'sem retentor', descricao)
57     # substitui os termos "c/re" ou "c/ret" por "com retentor"

```

```

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

```

```

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

```

```

145 'FROG', 'LIBERO', 'MAJESTY', 'MEST', 'MIDNIGHT', 'MORPH', 'NEO', 'PASSOL'],
146     'DAFRA':
147 ['APACHE', 'CITYCOM', 'KANSAS', 'LASER', 'NEXT', 'RIVA', 'ROADWIN', 'ZIG', 'SPEED'],
148     'SUZUKI': ['KATANA', 'YES', 'INTRUDER'],
149     'ZONGSHEN': ['ZS'],
150     'KASINSKI': ['COMET', 'MIRAGE'],
151     'POLARIS': ['SPORTSMAN', 'RZR', 'RANGER'],
152     'KAWASAKI':
153 ['NINJA', 'VERSYS', 'VOYAGER', 'GTR', 'KDX', 'KL', 'KX', 'KZ', 'ZR', 'ZZ', 'ER6N', 'ER6F'],
154     'DAYANG': ['DY1', 'DY2', 'DY5'],
155     'SUNDOWN': ['WEB', 'FIFITY', 'PALIO', 'PGO', 'STX', 'VBLADE', 'EVO', 'HUNTER
156 MAX'],
157     'SHINERAY':
158 ['BIKE', 'BRAVO', 'DISCOVER', 'EAGLE', 'INDIANAPOLIS', 'JET', 'NEW', 'WAVE',
159     'STRONG', 'SUPER SMART', 'VENICE', 'XY']}
160
161 # Função para pegar a chave pelo valor, dado que valor é único.
162 def pegaChave(v, dict):
163     for chave, valores in dict.items():
164         if type(valores) != type([1, 2]):
165             valores = [valores]
166         for valor in valores:
167             if v == valor:
168                 return chave
169     return "Não existe chave para esse valor."
170
171 def acrescentaMarca(descricao):
172     for marca in Marcas:
173         if re.search(marca, descricao.upper()):
174             descricao += " " + marca
175         for termo in Marcas[marca]:
176             t1 = termo.split()
177             if len(t1) > 1:
178                 pat = r"(?:" + t1[0] + r"|" + t1[1] + r").*(?:" + t1[0] + r"|" + t1[1] + r")"
179             elif len(termo) < 3:
180                 pat = termo + r"([0-9]{1,}|\b)"
181             else:
182                 pat = termo
183             resultados = re.findall(pat, descricao.upper())
184             if resultados:
185                 descricao += " " + marca
186                 descricao += " " + " ".join(resultados)
187                 descricao += " " + termo
188     descricao = " ".join(sorted(set(descricao.lower().split())))
189     return descricao
190
191 # ### Função final que transforma a DESCRICAO DO PRODUTO em Modelo para
192 classificar
193 def criaModelo(descricao):
194     descricao = limpaDescricao(descricao)
195     descricao = achaPalavraChave(descricao)
196     descricao = acrescentaMarca(descricao)
197     return descricao
198
199 # ### Função que determina a existência de retentor no kit e retorna True|False
200 def retentorAux(descricao):
201     # 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]
202         return True
203     else:
204         descricao=''
205         return False
206
207 # ### Define a função de classificação
208 def classificaAplicacao(descricao):
209     modelo = criaModelo(descricao)
210     novo_cvt = cvt.transform(pd.Series(modelo))
211     aplicacao = clfsvc.predict(novo_cvt)[0]
212     return aplicacao
213
214 # ## Carrega o modelo Linear SVC
215 # ### Carrega arquivos do pickle
216 # Caso não existam ou a versão for diferente, ignora
217 try:
218     with open(r'./pickle/clfsvc.pkl', 'rb') as file:
219         clfsvc = pickle.load(file)
220         file.close()
221 except:
222     print("Não foi possível carregar clfsvc.pkl.")
223
224 try:
225     with open(r'./pickle/cvt.pkl', 'rb') as file:
226         cvt = pickle.load(file)
227         file.close()
228 except:
229     print("Não foi possível carregar cvt.pkl.")
```