

## Elaborado com o Notebook Python Google Colab

```

1 # Load the Drive helper and mount
2 from google.colab import drive
3
4 # This will prompt for authorization.
5 drive.mount('/content/drive')

```

↳ Drive already mounted at /content/drive; to attempt to forcibly remount, call drive.m

```

1 !ls "/content/drive/My Drive/Databases"

```

↳

```

'Arvore Piratininga 2017.xlsx'    NoThemeTweets.csv
blogset-br.csv.gz                'Sentiment Analysis Dataset.csv'
imdb-reviews-pt-br.csv

```

```

1 import pandas as pd
2
3 dados = pd.read_csv("/content/drive/My Drive/Databases/NoThemeTweets.csv")

```

```

1 dados

```

↳

	id	tweet_text	tweet_date	sentiment	query_used
0	1031761728445530112	@Tixaa23 14 para eu ir :)	Tue Aug 21 04:35:39 +0000 2018	Positivo	:)
1	1031761040462278656	@drexalvarez O meu like eu já dei na época :)	Tue Aug 21 04:32:55 +0000 2018	Positivo	:)
2	1031760962372689920	Eu só queria conseguir comer alguma coisa pra ...	Tue Aug 21 04:32:37 +0000 2018	Positivo	:)
3	1031760948250456066	:D que lindo dia !	Tue Aug 21 04:32:33 +0000 2018	Positivo	:)
4	1031760895985246208	@Primo_Resmungao Pq da pr jeito!!é uma "oferta...	Tue Aug 21 04:32:21 +0000 2018	Positivo	:)
...	...	...	...	...	...
		Acordear 8 horas é tão	Fri Oct 12		

```

1 classificacao = dados["sentiment"].replace(["Negativo", "Positivo"], [0, 1])

```

```

1 dados["classificacao"] = classificacao

```

```

1 dados.groupby('classificacao').count()

```



	id	tweet_text	tweet_date	sentiment	query_used
<b>classificacao</b>					

0	522707	522707	522707	522707	522707
1	263107	263107	263107	263107	263107

1 dados



	id	tweet_text	tweet_date	sentiment	query_used
0	1031761728445530112	@Tixaa23 14 para eu ir :)	Tue Aug 21 04:35:39 +0000 2018	Positivo	:)
1	1031761040462278656	@drexalvarez O meu like eu já dei na época :)	Tue Aug 21 04:32:55 +0000 2018	Positivo	:)
2	1031760962372689920	Eu só queria conseguir comer alguma coisa pra ...	Tue Aug 21 04:32:37 +0000 2018	Positivo	:)
3	1031760948250456066	:D que lindo dia !	Tue Aug 21 04:32:33 +0000 2018	Positivo	:)
4	1031760895985246208	@Primo_Resmungao Pq da pr jeito!!é uma "oferta...	Tue Aug 21 04:32:21 +0000 2018	Positivo	:)
...	...	...	...	...	...
		Acordar 8 horas á tão	Fri Oct 12		

```

1 import nltk
2 from nltk import tokenize
3
4 nltk.download('stopwords')
5
6 token_espaco = tokenize.WhitespaceTokenizer()
7
8 token_pontuacao = tokenize.WordPunctTokenizer()

```



```

[nltk_data] Downloading package stopwords to /root/nltk_data...
[nltk_data] Package stopwords is already up-to-date!

```

```

1 palavras_irrelevantes = nltk.corpus.stopwords.words("portuguese")
2
3 from string import punctuation
4
5 pontuacao = list()
6 for ponto in punctuation:
7     pontuacao.append(ponto)
8
9 excluir = [":)", ":", "(", ":", ":", ":", "https"]

```

```

10
11 pontuacao_stopwords = pontuacao + palavras_irrelevantes + excluir
12
13 frase_processada = list()
14 for tweet in dados.tweet_text:
15     nova_frase = list()
16     tweet = tweet.lower()
17     palavras_texto = token_espaco.tokenize(tweet)
18     for palavra in palavras_texto:
19         if '@' not in palavra:
20             if palavra not in pontuacao_stopwords:
21                 nova_frase.append(palavra)
22     frase_processada.append(' '.join(nova_frase))
23
24 dados["tratamento_1"] = frase_processada

1 dados.head()

```

	id	tweet_text	tweet_date	sentiment	query_used	classi
0	1031761728445530112	@Tixaa23 14 para eu ir :)	Tue Aug 21 04:35:39 +0000 2018	Positivo	:)	
1	1031761040462278656	@drexalvarez O meu like eu já dei na época :)	Tue Aug 21 04:32:55 +0000 2018	Positivo	:)	
2	1031760962372689920	Eu só queria conseguir comer alguma coisa pra ...	Tue Aug 21 04:32:37 +0000 2018	Positivo	:)	
Tue Aug 21						

```
1 !pip install unicode
```

```
↳ Requirement already satisfied: unicode in /usr/local/lib/python3.6/dist-packages (1
```

```

1 import unicode
2
3 sem_acentos = [unicode.unidecode(tweet) for tweet in dados["tratamento_1"]]

```

```
1 dados["tratamento_2"] = sem_acentos
```

```
1 dados.head()
```

```
↳
```

	id	tweet_text	tweet_date	sentiment	query_used	classi
0	1031761728445530112	@Tixaa23 14 para eu ir :)	Tue Aug 21 04:35:39 +0000 2018	Positivo	:	
1	1031761040462278656	@drexalvarez O meu like eu já dei na época :)	Tue Aug 21 04:32:55 +0000 2018	Positivo	:	

```
1 nltk.download('rslp')
```

```
[nltk_data] Downloading package rslp to /root/nltk_data...
[nltk_data] Package rslp is already up-to-date!
True
```

Tue Aug 21

```
1 stopwords_sem_acento = [unicode.unidecode(texto) for texto in pontuacao_stopwords]
```

```
1 stemer = nltk.RSLPStemmer()
2
3 frase_processada = list()
4 for tweet in dados["tratamento_2"]:
5     nova_frase = list()
6     palavras_texto = token_pontuacao.tokenize(tweet)
7     for palavra in palavras_texto:
8         if palavra not in stopwords_sem_acento:
9             nova_frase.append(stemer.stem(palavra))
10    frase_processada.append(' '.join(nova_frase))
11
12 dados["tratamento_3"] = frase_processada
```

```
1 dados
```

```
[nltk_data]
```

	id	tweet_text	tweet_date	sentiment	query_used
0	1031761728445530112	@Tixaa23 14 para eu ir :)	Tue Aug 21 04:35:39 +0000 2018	Positivo	:)
1	1031761040462278656	@drexalvarez O meu like eu já dei na época :)	Tue Aug 21 04:32:55 +0000 2018	Positivo	:)
2	1031760962372689920	Eu só queria conseguir comer alguma coisa pra ...	Tue Aug 21 04:32:37 +0000 2018	Positivo	:)
3	1031760948250456066	:D que lindo dia !	Tue Aug 21 04:32:33 +0000 2018	Positivo	:)
4	1031760895985246208	@Primo_Resmungao Pq da pr jeito!!é uma	Tue Aug 21 04:32:21	Positivo	:)

```

1 import matplotlib.pyplot as plt
2 from wordcloud import WordCloud
3 import seaborn as sns

```

```

785800 1050705144007007000 Acordar 8 horas é tao

```

```

1 def nuvem_palavras_neg(texto, coluna_texto):
2     texto_negativo = texto.query("sentiment == 'Negativo'")
3     todas_palavras = ' '.join([texto for texto in texto_negativo[coluna_texto]])
4
5     nuvem_palavras = WordCloud(width=800, height=500, max_font_size=110, collocations=F
6
7     plt.figure(figsize=(10,7))
8     plt.imshow(nuvem_palavras, interpolation='bilinear')
9     plt.axis("off")
10    plt.show()

```

```

@andrebraga2806 Fri Oct 12

```

```

1 def nuvem_palavras_pos(texto, coluna_texto):
2     texto_positivo = texto.query("sentiment == 'Positivo'")
3     todas_palavras = ' '.join([texto for texto in texto_positivo[coluna_texto]])
4
5     nuvem_palavras = WordCloud(width=800, height=500, max_font_size=110, collocations=F
6
7     plt.figure(figsize=(10,7))
8     plt.imshow(nuvem_palavras, interpolation='bilinear')
9     plt.axis("off")
10    plt.show()

```

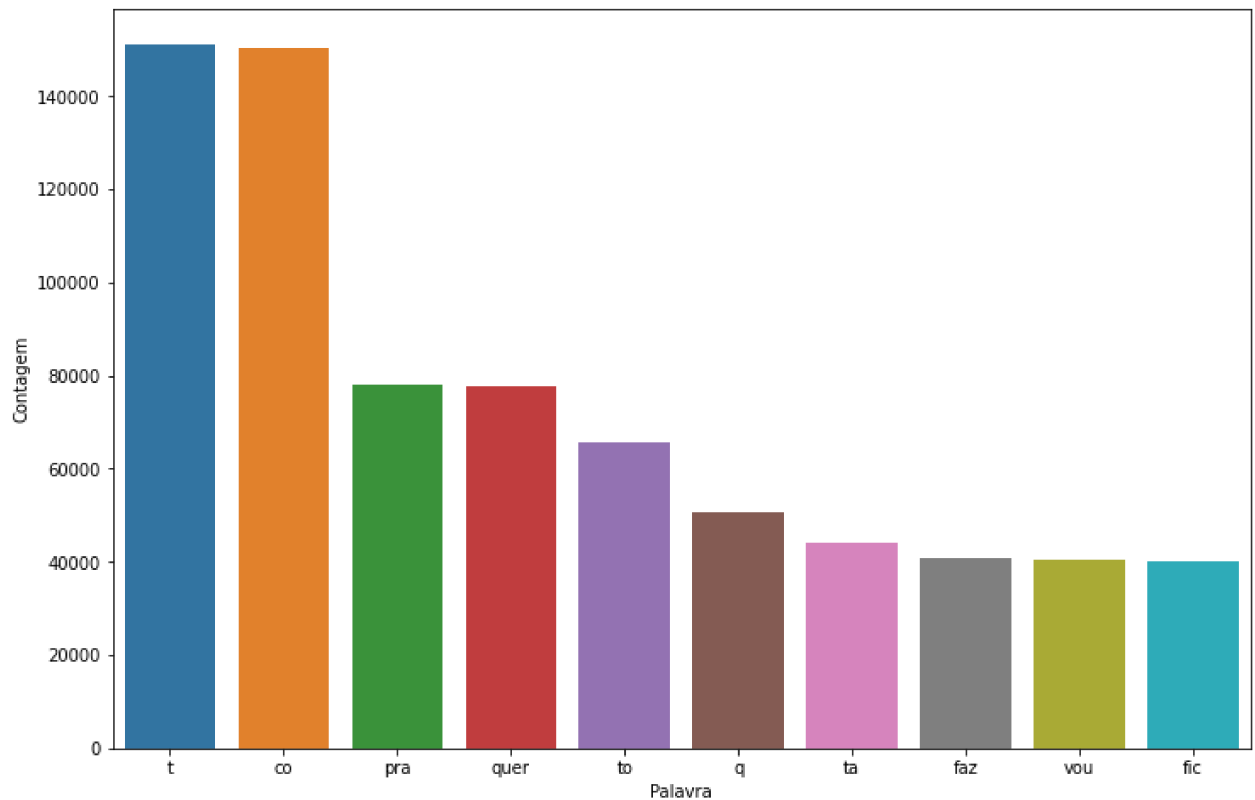
```

1 def pareto(texto, coluna_texto, quantidade):
2     todas_palavras = ' '.join([texto for texto in texto[coluna_texto]])
3
4     token_frase = token_espaco.tokenize(todas_palavras)
5     frequencia = nltk.FreqDist(token_frase)
6     df_frequencia = pd.DataFrame({"Palavra":list(frequencia.keys()), "Frequência":list(f
7     df_frequencia = df_frequencia.nlargest(quantidade, "Frequência")

```

```
1  nuvem_palavras_neg(dados, "tratamento_3")
```

[illegible]



```

1 def trata_frase(frase):
2     frase_tratada = list()
3     nova_frase = list()
4     frase = frase.lower()
5     palavras_texto = token_espaco.tokenize(frase)
6     for palavra in palavras_texto:
7         if '@' not in palavra:
8             if palavra not in pontuacao_stopwords:
9                 nova_frase.append(stemmer.stem(palavra))
10    frase_tratada.append(' '.join(nova_frase))
11    frase_tratada = [unicode.unidecode(frase) for frase in frase_tratada]
12
13    return frase_tratada

```

```

1 from sklearn.feature_extraction.text import TfidfVectorizer
2
3 from sklearn.linear_model import LogisticRegression
4
5 from sklearn.model_selection import train_test_split
6
7 tfidf1 = TfidfVectorizer(lowercase = False, max_features=50)
8
9 tfidf_bruto = tfidf1.fit_transform(dados["tweet_text"])
10
11 treino, teste, classe_treino, classe_teste = train_test_split(tfidf_bruto, dados["cla
12
13 regressao_logistica1 = LogisticRegression(solver = "lbfgs")
14
15 regressao_logistica1.fit(treino, classe_treino)

```

```

16
17 acuracia_tfidf_bruto = regressao_logistica1.score(teste, classe_teste)
18
19 print(acuracia_tfidf_bruto)

```

0.6881763669866737

```

1 tfidf2 = TfidfVectorizer(lowercase = False, max_features=50)
2
3 tfidf_tratado = tfidf2.fit_transform(dados["tratamento_3"])
4
5 treino, teste, classe_treino, classe_teste = train_test_split(tfidf_tratado, dados["c
6
7 regressao_logistica2 = LogisticRegression(solver = "lbfgs")
8
9 regressao_logistica2.fit(treino, classe_treino)
10
11 acuracia_tfidf_tratado = regressao_logistica2.score(teste, classe_teste)
12
13 print(acuracia_tfidf_tratado)

```

0.6885021429953068

```

1 tfidf3 = TfidfVectorizer(lowercase=False,ngram_range=(1,2))
2
3 tfidf_ngrams = tfidf3.fit_transform(dados["tratamento_3"])
4 treino, teste, classe_treino, classe_teste = train_test_split(tfidf_ngrams, dados["cl
5 regressao_logistica3 = LogisticRegression(solver = "lbfgs", max_iter=1000)
6 regressao_logistica3.fit(treino, classe_treino)
7 acuracia_tfidf_ngrams = regressao_logistica3.score(teste, classe_teste)
8
9 print(acuracia_tfidf_ngrams)

```

0.794929092815621

```

1 print(tfidf1.get_feature_names())

```

['Eu', 'amo', 'as', 'bem', 'co', 'com', 'da', 'de', 'dia', 'do', 'ela', 'ele', 'em',

```

1 trata_frase("Faltá Frase")

```

['falt fras']

```

1 def testa_frase(frase):
2     print("Regressão 1")
3     print(regressao_logistica1.predict(tfidf1.transform(trata_frase(frase))))
4     print(regressao_logistica1.predict_proba(tfidf1.transform(trata_frase(frase))))
5     print("Regressão 2")
6     print(regressao_logistica2.predict(tfidf2.transform(trata_frase(frase))))
7     print(regressao_logistica2.predict_proba(tfidf2.transform(trata_frase(frase))))
8     print("Regressão 3")
9     print(regressao_logistica3.predict(tfidf3.transform(trata_frase(frase))))

```



```
print(regressao_logistica3.predict_proba(tfidf3.transform(trata_frase(frase))))

10 print(regressao_logistica3.predict_proba(tfidf3.transform(trata_frase(frase))))

1 testa_frase("Um dia ruim")
2 testa_frase("Um dia bom")

↳ Regressão 1
[1]
[[0.36962036 0.63037964]]
Regressão 2
[1]
[[0.47902906 0.52097094]]
Regressão 3
[0]
[[0.87927611 0.12072389]]
Regressão 1
[1]
[[0.36962036 0.63037964]]
Regressão 2
[1]
[[0.22172233 0.77827767]]
Regressão 3
[1]
[[0.08002752 0.91997248]]
```

## Resultados Encontrados

O treinamento sem o tratamento dos dados gerou uma taxa de acerto de 0.6881, já com o tratamento dos dados a taxa vai para 0.6885, o que não representa um ganho significativo. Utilizando-se um bigram a taxa de acerto foi de 0.7949, que representa um aumento de 15,5% com relação aos dados sem tratamento. Foi possível observar também que a confiança do resultado é significativamente maior utilizando bigrams.