

## **Twitter Sentiment Analysis**

Jefté Lopes G, Lucas Lopes, Edivaldo Araújo Jr.

### ***Resumo***

O objetivo desse estudo de caso é criar um modelo que analisa um ou mais Tweets para prever o sentimento(Positivo ou Negativo) presente em cada Tweet é usado o Processamento de linguagem natural (NLP) juntamente com aprendizagem de maquina para construir esse modelo. Buscamos apresentar tudo com uma didática simplificada para maximizar a medida do possível o entendimento do assunto, que a princípio pode parecer complexo.

### ***Abstract***

The purpose of this case study is to create a model that analyzes one or more Tweets to predict the feeling (Positive or Negative) present in each Tweet. Natural language processing (NLP) is used in conjunction with machine learning to build this model. We seek to present everything with a simplified didactic to maximize the extent of the possible understanding of the subject, which at first may seem complex.

### **1. Problema**

Classificar os sentimentos do Twitter como Negativo ou positivo. O objetivo desse estudo de caso é criar um modelo que analisa um ou mais Tweets para prever o sentimento(Positivo ou Negativo) presente em cada Tweet.

### 1.1. 1ª Solução e 2ª Solução

Escolhemos realizar uma segunda solução pelo fato de que os resultados obtidos e aprendidos durante a primeira solução não foram satisfatórios, portanto, buscamos uma nova solução que em tese traria um melhor resultado a explicação das duas soluções estão detalhadas no documento.

1ª Solução	2ª Solução
Base de dados	Base de dados
Idioma Inglês	Idioma Português-Brasil
Utilizando tweets	Utilizando tweets
Baixada no kaggle	Criada
NLP(Natural Language Processing)	NLP(Natural Language Processing)
Bag of Word	Bag of Word
Naive Baye	Naive Baye
Visualização dos Resultado	Visualização dos Resultado

#### 1. 1ª Solução

- Entender a Declaração do Problema e o caso de negócios.
- Importar bibliotecas e conjuntos de dados.
- Executar a análise exploratória dos dados.
- Plotar a nuvem de palavras.
- Executar a limpeza de dados - remover pontuação.
- Executar a limpeza de dados - remover palavras de parada(stop words).
- Executar vetorização de contagem (Tokenization).
- Criar um pipeline para remover palavras irrelevantes, pontuação e realizar tokenização.
- Compreender a teoria e a intuição por trás dos classificadores Naive Bayes.
- Treinar um Classificador Naive Bayes.
- Avaliar o desempenho do modelo treinado.

#### 2. 2ª Solução

- Entender a Declaração do Problema e o caso de negócios.
- Importar bibliotecas e conjuntos dos dados.
- Executar a análise exploratória de dados.
- Plotar a nuvem de palavras.
- Executar vetorização de contagem (Tokenization).
- Treinar um Classificador Naive Bayes.
- Avaliar o desempenho do modelo treinado.
- Salvar o modelo treinado.

## 2. Entender a Declaração do Problema e o caso de negócios

Esse documento foi criado, para informar o passo a passo feito com base no nosso notebook que é usado o Processamento de linguagem natural(NLP) juntamente com aprendizagem de maquina para construir um modelo que analisa milhares de Tweets para prever o sentimentos das pessoas.

A Inteligencia artificial e a analise de sentimentos baseadas em aprendizado de maquina é crucial para empresas, visto que, os insight revelado pela analise visa indicar o grau de qualidade dos serviços e/ou produtos da empresa de acordo com os clientes.(Trabalhos Futuros)

Esse projeto é diretamente aplicável a praticamente qualquer empresa que disponha de meios online(Twitter, Instagran, Facebook, WebSite) para interagir com seus clientes.(Trabalhos Futuros)

Os algoritmos podem ser usados para detectar e possivelmente sinalizar automaticamente tweets de ódio e racismo.(Trabalhos Futuros)

### 2.1. Importar bibliotecas e conjuntos de dados

Com base no nosso notebook explicarei passo a passo sobre a 1ª Solução, nela foram utilizados Pacotes essenciais para analise numéricas, manipulação de data frames e visualização de dados.

#### 1. Bibliotecas de apoio:

---

```
1 #Pacotes essenciais para analise numericas
2 #manipulação de data frames e visualização de dados.
3 import pandas as pd
4 import numpy as np
5 import seaborn as sns
6 import matplotlib.pyplot as plt
```

---

pandas para carregar e manipular os dados.

Numpy é usado para trabalhar com matrizes. Também possui funções para trabalhar no domínio da álgebra linear, transformada de Fourier e matrizes.

pyplot from matplotlib é usado para visualizar os resultados.

Seaborn é uma biblioteca de visualização de dados Python baseada em matplotlib. Ele fornece uma interface de alto nível para desenhar gráficos estatísticos atraentes e informativos.

## 2. Pacote para a criação de nuvens de palavras:

```
1 In [ ]: # Pacote para a criação de nuvens de palavras
2         !pip install WordCloud
3         from wordcloud import WordCloud
```

A base de dados desse exemplo está no meu repositório do github e pode ser baixada no link abaixo.

[https://raw.githubusercontent.com/JefteLG/Twitter\\_Sentiment\\_Analysis/main/Data\\_Set/data\\_base\\_estudos/twitter.csv](https://raw.githubusercontent.com/JefteLG/Twitter_Sentiment_Analysis/main/Data_Set/data_base_estudos/twitter.csv)

## 3. Utilizando o pandas para ler o arquivo CSV

```
1 #estruturando em um data frame na variavel `tweets_df` por meio do metodo read_csv.
2 tweets_df = pd.read_csv('/tmp/twitter.csv')
3
4 #visualização1.
5 tweets_df.head()
6
7 #visualização2.
8 tweets_df
9
10 #Informações sobre a quantidade de Tweets, memoria, tipo de dados e dados faltantes
11 tweets_df.info()
```

### 3. Executar a análise exploratória de dados

Esta é uma função no nível dos eixos e desenhará o mapa de calor para os eixos ativos no momento. Nesse caso a função verifica se existe dados faltantes.

```
1 sns.heatmap(tweets_df.isnull(), yticklabels=False, cbar=False, cmap="Blues")
```

```
<matplotlib.axes._subplots.AxesSubplot at 0x7f3fdbd2e208>
```



O método `sns.countplot()` é usado para mostrar as contagens de observações em cada categoria categórica usando barras (Histograma).

---

```
1 sns.countplot(data=tweets_df, x='label', palette='Set2')
2 plt.show()
```

---

---

```
1 #Nova coluna com o tamanho dos tweets
2 tweets_df['length'] = tweets_df['tweet'].apply(len)
```

---

#### 1. Teste:

---

```
1 from nltk.tokenize import TweetTokenizer
2 tweets_df[tweets_df['id']==31958]['tweet'].iloc[0]
3 tweetTK = TweetTokenizer()
4 tweetTK.tokenize(tweets_df[tweets_df['id']==31958]['tweet'].iloc[0])
```

---

#### 2. Fim Teste:

---

```
1 #Descobrir o tamanho maximo, minimo e medio dos tweets.
2 tweets_df['length'].plot(bins=100, kind='hist', figsize=(12,8), color='g')
3 tweets_df.describe()
4
```

---

```

5  # Selecionar a menor frase
6  tweets_df[tweets_df['length']==11]['tweet'].iloc[0]
7
8  # Separa o DataFrame em dois Dataframes, um com sentimentos
9  #positivos e o outro com sentimentos negativos.
10 positive_df = tweets_df[tweets_df['label']==0]
11 negative_df = tweets_df[tweets_df['label']==1]
12
13 positive_df
14
15 negative_df

```

---

### 3.1. Plotar a nuvem de palavras

Wordcloud – em português, nuvem de palavras ou nuvem de tags – é um tipo de visualização de dados muito poderoso e ferramenta de Data Science usado quando estamos trabalhando com textos, documentos, pesquisas, entre outras.

Resumidamente, é como se você estivesse contando a frequência com que cada palavra aparece em um texto. Com essa frequência, você define tamanhos proporcionais às palavras e quanto maior for a ocorrência da palavra no texto maior será na nuvem ou desenho.

Segue abaixo a implementação:

#### 1. Criar uma lista de tweets:

```

1  sentences = tweets_df['tweet'].tolist()

```

---

#### 2. O método join() pega todos os itens em um iterável e os une em uma string é usado o espaço como separador

```

1  sentences_as_one_string = " ".join(sentences)

```

---

#### 3. Nuvem de palavras de todos os tweets NEGATIVOS:

```

1  negative_sentences = negative_df['tweet'].tolist()
2  negative_as_one_string = " ".join(negative_sentences)
3
4  plt.figure(figsize=(10,10))
5  plt.imshow(WordCloud().generate(negative_as_one_string))

```

---

#### 4. Nuvem de palavras de todos os tweets POSITIVOS:

```

1  positive_sentences = positive_df['tweet'].tolist()
2  positive_as_one_string = " ".join(positive_sentences)
3
4  plt.figure(figsize=(10,10))
5  plt.imshow(WordCloud().generate(positive_as_one_string))

```

---

## 4. Processamento de Linguagem Natural

Processamento de linguagem natural (PLN) é uma vertente da inteligência artificial que ajuda computadores a entender, interpretar e manipular a linguagem humana.

A linguagem humana é surpreendentemente complexa e diversa. Nós nos expressamos de infinitas maneiras, tanto verbalmente quanto por escrito. Não apenas existem centenas de idiomas e dialetos, como há também um conjunto único de regras gramaticais e de sintaxe, expressões e gírias dentro de cada um deles. Quando escrevemos, costumamos cometer erros ou abreviar palavras, ou omitimos pontuações quando falamos, carregamos sotaques regionais, tendemos a murmurar e emprestamos termos de outros idiomas.

Embora o aprendizado supervisionado, o aprendizado não-supervisionado e, especificamente, o deep learning sejam hoje amplamente utilizados para modelar a linguagem humana, há também a necessidade de compreensão sintática e semântica, além de domínio, que não estão necessariamente presentes nessas abordagens de machine learning. O PLN é importante porque ajuda a resolver a ambiguidade na linguagem e adiciona uma estrutura numérica útil aos dados para muitas aplicações downstream, como reconhecimento de fala ou análise de texto.

### 4.1. Executar a limpeza de dados - remover pontuação

Vamos criar uma função para tratar nossos dados, removendo símbolos e caracteres especiais, essa remoção tem como objetivo ter uma melhor precisão nas previsões.

Nesta etapa vamos começar a “limpar” nosso texto, removendo tudo que não tem valor para nosso objetivo.

Para fazermos essa limpeza existem vários métodos, nós utilizamos cinco que será mostrado a implementação logo abaixo:

---

```
1 import string
2 text = 'Good morning beautiful people :)...
3 'I am having fun learning Machine learning and artificial intelligence'
```

---

#### 1. Método Nº1:

```
1 test_punc_remove = [char for char in text if char not in
2     string.punctuation]
3 test_punc_remove
4 test_punc_remove = ''.join(test_punc_remove)
5 test_punc_remove
```

#### 2. Método Nº1.1:

```
1 text_remove_punct = []
2 for char in text:
3     if char not in string.punctuation:
4         text_remove_punct.append(char)
5
6 new_text = ''.join(text_remove_punct)
7
8 new_text
```

### 3. Método N°2:

```
1 punct = string.punctuation + string.digits
2 table_tst = str.maketrans('','',punct)
3 newtext = text.translate(table_tst)
4 newtext
```

### 4. Método N°3:

```
1 punct = string.punctuation + string.digits
2 table_ = str.maketrans(punct, ' '*len(punct))
3 newtext = ' '.join(text.translate(table_).split())
4 newtext
```

### 5. Método N°4:

```
1 punct = string.punctuation + string.digits
2     for s in punct:
3         text = text.replace(s, '')
4     text
```

### 6. Método N°5:

```
1 import re
2 newtext = re.sub(r'[A-Za-z]+', ' ', text)
3 newtext
4 'Good morning beautiful people I am having fun learning
5 'Machine learning and artificial intelligence'
```

## 4.2. Executar a limpeza de dados - remover palavras de parada(stop words)

As palavras irrelevantes são frequentemente removidas do texto antes do treinamento de modelos de aprendizado profundo e de aprendizado de máquina, já que as palavras irrelevantes ocorrem em abundância, fornecendo, portanto, pouca ou nenhuma informação exclusiva que pode ser usada para classificação ou agrupamento.

**Usando a biblioteca NLTK do Python** é uma das bibliotecas Python mais antigas e mais comumente usadas para processamento de linguagem natural. O NLTK oferece suporte à remoção de palavras de interrupção e você pode encontrar a lista de palavras de interrupção no corpusmódulo. Para remover palavras irrelevantes de uma frase, você pode dividir seu texto em palavras e, em seguida, remover a palavra se ela sair da lista de palavras irrelevantes fornecida pelo NLTK.

Vamos ver na prática:

```
1 #Pacotes para facilitar o processamento de linguagem natural
2 import nltk
3 import string
4 nltk.download('stopwords')
5
6 from nltk.corpus import stopwords
7 stopwords.words('english')
8
```



```

9 text = 'Good morning beautiful people :)... I am having fun
10 'learning Machine learning and artificial intelligence in 2020'
11 #caracteres para remover
12 puncts = string.digits + string.punctuation
13
14 #remoção de caracteres
15 text_remove_punct = [char for char in text if char not in puncts]
16
17 #pega todos os itens em um iterável e os une em uma string.
18 text_remove_punct = ''.join(text_remove_punct)
19
20 text_remove_punct
21
22 #remoção de stopwords
23 text_remove_stopword = [word for word in text_remove_punct.split()
24 if word.lower() not in stopwords.words('english')]
25
26 text_remove_stopword

```

**Pipeline remover pontuações e stopwords**, Pipelines são interessantes para reduzir código e automatizar fluxos é também executa uma série de ações que normalmente são definidas por funções e/ou métodos, a execução ocorre de maneira linear.

```

1 def pipeline(text):
2     #frase sem nenhum tratamento
3     print(text)
4     #remover as pontuações
5     text_remove_punct = [char for char in text if char not in puncts]
6     #unir todos os elemntos da lista
7     text_remove_punct = ''.join(text_remove_punct)
8     #frase com pontuações removidas
9     print(text_remove_punct)
10
11     # remover as stopwords
12     text_remove_stopword = [word for word in text_remove_punct.split()
13 if word.lower() not in stopwords.words('english')]
14     text_remove_stopword = ' '.join(text_remove_stopword)
15     #frase com stopwords removidas
16     print(text_remove_stopword)
17
18     pipeline(text)
19

```

### 4.3. Executar vetorização de contagem(Tokenization)

A tokenização, também conhecida como segmentação de palavras, quebra a sequência de caracteres em um texto localizando o limite de cada palavra, ou seja, os pontos onde uma palavra termina e outra começa [Palmer, 2010]. Para fins de linguística computacional, as palavras assim identificadas são frequentemente chamadas de *tokens*.

Quando a linguagem escrita é armazenada em um arquivo de computador, ela é normalmente representada por meio de uma sequência ou string (do inglês, "cadeia") de caracteres. Isto é, em um arquivo de texto padrão, as palavras são strings, as sentenças são strings e o próprio texto não passa, no fundo, de uma longa string. Os caracteres de uma string não precisam ser necessariamente alfanuméricos; estas também podem incluir caracteres especiais que representem os espaços, as tabulações, os sinais de nova linha, etc.

Muito do processamento computacional é realizado acima do plano dos caracteres. Quando compilamos uma linguagem de programação, por exemplo, o compilador espera que o input (os dados de entrada) seja uma sequência de tokens com os quais ele seja capaz de lidar; por exemplo, as classes dos identificadores, constantes textuais e alfanuméricas. De forma análoga, um parser espera que seu input seja uma sequência de tokens de palavras e não uma sequência de caracteres individuais.

Em sua forma mais simples, a tokenização de um texto envolve a busca das posições da string que contem os chamados "caracteres em branco"(espaços, tabulações ou sinais de nova linha) ou sinais de pontuação específicos, separando a string em tokens nestas posições.

**Metodologia: Bag of Word**, Bag of words é um método usado para extrair recursos de documentos de texto. Esses recursos são utilizados para treinar algoritmos de aprendizagem de máquina. Ele cria um vocabulário de todas as palavras do conjunto de treinamento.

Um dos maiores problemas com o texto é que ele é confuso e não estruturado, e os algoritmos de aprendizado de máquina preferem entradas de comprimento fixo bem definidas e estruturadas, utilizando a técnica Bag-of-Words, podemos converter textos, em um comprimento fixo vetor. Além disso, os modelos de aprendizado máquina funcionam com dados numéricos em vez de dados textuais, então para ser mais específico, usando a técnica do saco de palavras (BoW), é possível converter texto em vetores para facilitar o trabalho com números equivalentes.

Veremos agora na prática:

```
1  #Classe responsável por Converter uma coleção de documentos de texto
2  #em uma matriz de contagens de tokens
3  from sklearn.feature_extraction.text import CountVectorizer
4  sample_data=[
5      'Porcaria de produto',
6      'Obrigado pelo retorno estou muito satisfeito com seu trabalho?',
7      'Voces ainda vão me responder?',
8      'O pior serviço de atendimento de todos!!!'
9  ]
10
11 sample_data2=[
```

```

12         'Hello World',
13         'Hello Hello hello World',
14         'Hello World world world'
15     ]
16
17     # Aprender o vocabulário do texto e retornar um array, [n_samples, n_features]
18     vectorizer = CountVectorizer()
19     X = vectorizer.fit_transform(sample_data)
20     print(vectorizer.get_feature_names())
21
22
23     # versão tokenizada das strings no dataset
24     # (id palavra, id frase)  frequência da palavra na frase
25     c = X
26     print(c)
27
28     print(vectorizer.vocabulary_)
29     {'porcaria': 11, 'de': 4, 'produto': 12, 'barato': 2, 'ruim': 15, 'obrigado': 8,
30     'pelo': 9, 'retorno': 14, 'estou': 5, 'muito': 7, 'satisfeito': 16, 'com': 3,
31     'seu': 18, 'trabalho': 20, 'você': 21, 'ainda': 0, 'vão': 22, 'me': 6,
32     'responder': 13, 'pior': 10, 'serviço': 17, 'atendimento': 1, 'todos': 19}
33
34     #Retorne uma representação densa desta matriz.
35     print(X.toarray())

```

#### 4.4. Criar um pipeline para remover palavras irrelevantes, pontuação e realizar tokenização

1. Limpar as mensagens:

---

```

1 def message_cleaning(msg):

```

---

2. Remoção das pontuações:

---

```

1 text_remove_punct = [char for char in msg if char not in char_remove]
2 text_remove_punct = ''.join(text_remove_punct)

```

---

3. Remoção das palavras de parada:

---

```

1 text_remove_stopword = [word for word in text_remove_punct.split()
2 if word.lower() not in stopwords.words('english')]
3 # text_remove_stopword = ' '.join(text_remove_stopword)
4 return text_remove_stopword

```

---

#### 4.5. Treinar um Classificador Naive Bayes

Em termos simples, um classificador Naive Bayes assume que a presença de uma característica particular em uma classe não está relacionada com a presença de qualquer

outro recurso. O modelo Naive Bayes é fácil de construir e particularmente útil para grandes volumes de dados.

Classificadores Naive Bayes utilizados principalmente em classificação de textos (devido a um melhor resultado em problemas de classes múltiplas e regra de independência) têm maior taxa de sucesso em comparação com outros algoritmos. E bastante utilizado em Análise de Sentimento que é o nosso caso, para identificar sentimentos positivos e negativos.

Segue a implementação de treino:

---

```
1 X.shape
2 (31962, 47386)
3 y.shape
4 (31962,)
```

---

#### 1. Separação dos dados de treino e teste:

```
1 from sklearn.model_selection import train_test_split
2 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size = 0.2)
```

#### 2. Utilização do algoritmo de Naive Bayes para o treinamento:

```
1 from sklearn.naive_bayes import MultinomialNB
2 NB_classifier = MultinomialNB()
3 NB_classifier.fit(X_train, y_train)
```

### 5. Avaliar o desempenho do modelo treinado

---

```
1 from sklearn.metrics import classification_report, confusion_matrix
2 import seaborn as sns
3 y_predict_test = NB_classifier.predict(X_test)
4 cm = confusion_matrix(y_test, y_predict_test)
5 cm
6 # sns.heatmap(cm, annot=True)
7 print(classification_report(y_test, y_predict_test))
```

---

## **6. Tópicos relacionado a 2ª Solução**

A partir de agora serão explicados os tópicos relacionados apenas a 2ª solução, os tópicos abordados serão Criação da base de dados (feita por nós), Análise exploratória dos dados, Avaliar o desempenho do modelo treinado e Salvar o modelo treinado.

### **6.1. Criação da base de dados utilizando a API do Twitter com python pelo tweepy**

Para ser possível acessar a api é necessário ter uma conta de desenvolvedor e a conta em questão foi criada por fins acadêmicos no intuito de desenvolver um modelo de aprendizagem para reconhecimentos de sentimentos positivos e negativos a partir de tweets.

Sobre o acesso às informações coletadas segundo as diretrizes de privacidade do twitter não é possível disponibilizar dados referentes aos usuários sendo que o não cumprimento do mesmo é passivo a revogação da licença de desenvolvedor é perdendo assim o acesso a api.

### **6.2. Sobre o desenvolvimento do algoritmo para a criação do dataset:**

Todo o algoritmo foi elaborado em python, e primeiramente foi elaborado com uma integração a um banco de dados mysql, porém pela demora na inserção e validação dos dados decidimos utilizar-se do pandas para a criação do dataset em um documento CSV.

Sobre o método de tratamento dos dados após ser coletados do twitter utilizamos de expressões regulares para fazer uma limpeza em cada tweet retirando links, nomes de usuários, pontuações, espaços em excesso e por fim transformando cada caractere para caixa baixa(lower case).

### **6.3. Definição dos termos a serem buscados no twitter:**

A princípio foi definido somente emojis para a classificação de tweets sendo :) para positivo e :( para negativos e com a melhor compreensão da api as buscas foram feitas em cima das seguintes combinações:

#### **1. Para tweets positivos:**

- :) feliz
- :) alegre
- \ø/ empolgado
- :) amor
- \ø/ confiante
- :) apaixonado
- \ø/ otimista
- resiliência

#### **2. Para tweets negativos**

- :( triste
- :-( triste

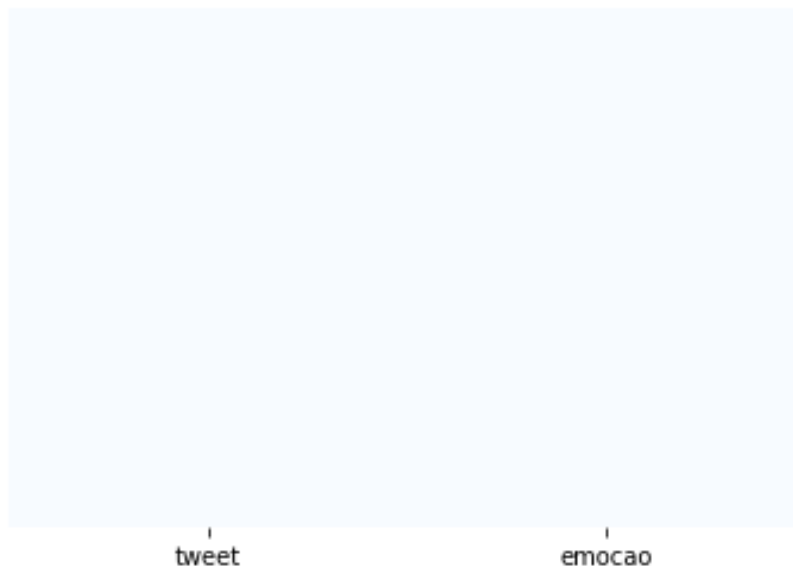
- :( mau
- :( morrer

Grande parte dos tweets eram em sua maioria negativos, sendo necessário incluir mais tweets positivos para igualar a base de dados e acabando tendo mais instruções de pesquisa positivas do que negativas.

## 7. Análise exploratória dos dados

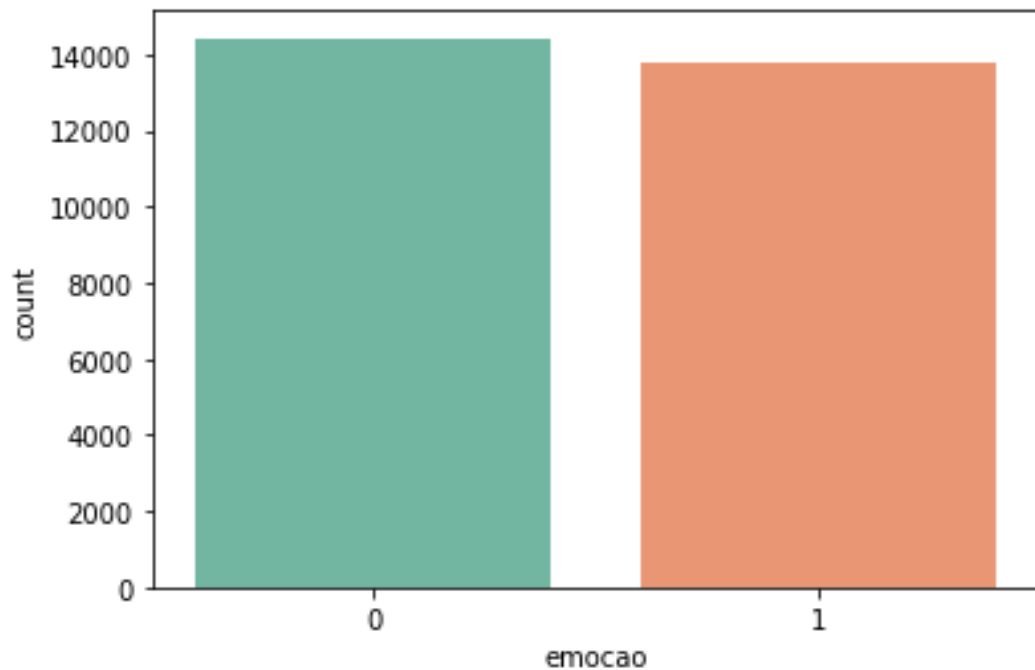
1. Esta é uma função no nível dos eixos e desenhará o mapa de calor para os eixos ativos no momento. Nesse caso a função verifica se existe dados faltantes.

```
1 sns.heatmap(tweets_df_dataframe.isnull(), yticklabels=False, cbar=False, cmap="
<matplotlib.axes._subplots.AxesSubplot at 0x7f459cc58438>
```



2. O método `sns.countplot()` é usado para mostrar as contagens de observações em cada categoria categórica usando barras (Histograma).

```
1 sns.countplot(data=tweets_df_dataframe, x='emocao', palette='Set2')
2 plt.show()
```



3. Separar o DataFrame em dois Dataframes, um com sentimentos positivos e o outro com sentimentos negativos.

```
1 positive_df = tweets_df_dataframe[tweets_df_dataframe['emocao']==1]
2 negative_df = tweets_df_dataframe[tweets_df_dataframe['emocao']==0]
```

## 8. Avaliar o desempenho do modelo treinado

```
1 from sklearn.metrics import classification_report, confusion_matrix
2 import seaborn as sns
```

1. Visualização da matriz de confusão:

```
1 y_predict_test = NB_classifier.predict(X_test)
2 cm = confusion_matrix(y_test, y_predict_test)
3 cm
4 # sns.heatmap(cm, annot=True)
5 array([[2779, 73],
6        [ 201, 2585]])
```

2. Métricas:

```
1 print(classification_report(y_test, y_predict_test))
2 precision    recall  f1-score   support
3
4         0       0.93      0.97      0.95      2852
5         1       0.97      0.93      0.95      2786
6
7    accuracy                   0.95      5638
8    macro avg       0.95      0.95      0.95      5638
9    weighted avg       0.95      0.95      0.95      5638
```

## 8.1. Testes e Classificações

### 1. Quantidade de classificações certas:

```
1 tweet_correto = (y_test == NB_classifier.predict(X_test)).sum()
2 print(f'{tweet_correto} classificações certas')
3
4 5364 classificações certas
```

### 2. Quantidade de classificações incorretas:

```
1 tweet_incorreto = (y_test.size - tweet_correto)
2 print(f'{tweet_incorreto} classificações incorretas')
3 274 classificações incorretas
```

### 3. Porcentagem de accuracy:

```
1 error = tweet_incorreto/(tweet_correto + tweet_incorreto)
2 print(f'accuracy do modelo {1-error:.2%}')
3 accuracy do modelo 95.14%
```

### 4. Accuracy:

```
1 NB_classifier.score(X_test, y_test)
2 0.9547711954593827
```

## Exemplos

```
1 exemplos = [
2     'O meu dia hoje foi pessimo',
3     'Minha cabeça está doendo muito',
4     'Eu acho que me dei muito mal na minha entrevista de emprego hoje',
5     'acordei feliz mas meu dia acabou muito triste',
6     'eu quero morrer',
7     'vou matar o',
8     'acordei mó triste mas meu dia acabou muito feliz',
9     'estou bastante otimista hoje',
10    'Vou ter um irmãozinhoooooo eba',
11    'Vou ser mamãe, estou muito contente',
12    'Me desejem sorte',
13    'o é meu melhor amigo',
14    'fiquei na bad',
15    'estou hypado',
16    'estou no hype'
17 ]
18
19 exemplos2 = [
20     'Gosto de um bom churrasco'
21 ]
```



### 1. Tokenizar os exemplos:

```
1 exemplos = vectorizer.transform(exemplos)
2 exemplos2 = vectorizer.transform(exemplos2)
3 print(exemplos.toarray())
4
5 t = exemplos.toarray()
6 print(t[0][10:40])
```

### 2. classificar o sentimento dos exemplos:

```
1 NB_classifier.predict(exemplos)
2 array([0, 0, 0, 0, 0, 0, 0, 0, 1, 1, 1, 1, 1, 0, 0, 1, 1])
3
4 NB_classifier.predict(exemplos)
5 array([0, 0, 0, 0, 0, 0, 0, 0, 1, 1, 0, 1, 1, 0, 0, 1, 1])
```

### 3. TESTE 1:

```
[39] y_predict_test = NB_classifier.predict(X_test)
      cm = confusion_matrix(y_test, y_predict_test)
      cm
      # sns.heatmap(cm, annot=True)

array([[2790,  76],
       [ 170, 2606]])

[33] print(classification_report(y_test, y_predict_test))
```

	precision	recall	f1-score	support
0	0.94	0.98	0.96	2881
1	0.98	0.94	0.96	2761
accuracy			0.96	5642
macro avg	0.96	0.96	0.96	5642
weighted avg	0.96	0.96	0.96	5642

```
[44] NB_classifier.predict(exemplos)

array([1, 0, 0, 0, 0, 0, 0, 0, 1, 1, 1, 1, 1, 0, 0, 1, 1])

[24] NB_classifier.predict(exemplos)

array([0, 0, 0, 0, 0, 0, 0, 0, 1, 1, 0, 1, 1, 0, 0, 1, 1])
```

#### 4. TESTE 2:

```
[80] y_predict_test = NB_classifier.predict(X_test)
      cm = confusion_matrix(y_test, y_predict_test)
      cm
      # sns.heatmap(cm, annot=True)
```

```
array([[2828,   70],
       [ 174, 2570]])
```

```
[88] print(classification_report(y_test, y_predict_test))
```

	precision	recall	f1-score	support
0	0.94	0.98	0.96	2898
1	0.97	0.94	0.95	2744
accuracy			0.96	5642
macro avg	0.96	0.96	0.96	5642
weighted avg	0.96	0.96	0.96	5642

```
[94] NB_classifier.predict(exemplos)
```

```
array([0, 0, 0, 0, 0, 0, 0, 1, 1, 1, 1, 1, 0, 0, 1, 1])
```

```
[24] NB_classifier.predict(exemplos)
```

```
array([0, 0, 0, 0, 0, 0, 0, 1, 1, 0, 1, 1, 0, 0, 1, 1])
```

## 9. Salvar o modelo treinado

Salvamos o modelos de machine learning com Pickle, ele permite a serialização de objetos, ou seja, os transforma em sequências de bytes. Armazenando toda informação necessária, consegue reconstruir objetos aplicando uma lógica interna.

Para salvar um modelo, usa-se o método `dump()` – basta informar a variável referente ao objeto e o nome do arquivo a ser criado.

```
1 import pickle
2 #salvar o modelo Naive Bayes (NB_classifier) no arquivo model_1.pkl
3 with open('/tmp/model_2.pkl', 'wb') as file:
4     pickle.dump(NB_classifier, file)
```

1. Carregando modelos de machine learning com Pickle, já que o modelo está treinado e validado, uma vez carregado ele já esta pronto para ser usado. No caso da célula abaixo, o modelo de machine learning foi atribuído à variável `model`.

```
1 with open('/tmp/model_1.pkl', 'rb') as f:
2     model = pickle.load(f)
```

2. Classificar o sentimento:

```
1 model.predict(exemplos2)
```

## **10. Conclusão**

Neste trabalho foi apresentada uma ferramenta para análise de sentimentos em mensagens enviadas através do Twitter para conhecer o que os usuários da rede social estão sentindo sobre determinados assuntos em relação aos sentimentos positivos e negativos presentes nas mensagens. A principal função da ferramenta é comparar os sentimentos das mensagens entre duas consultas sobre o sentimento presente nas mensagens (positivo e negativo).

Ainda existe muito espaço para a ferramenta evoluir, mas existem dois pontos principais que precisam ser melhorados. O primeiro está relacionado à qualidade do classificador. Um estudo maior precisa ser realizado para que se consiga resultados melhores, pois a análise de sentimentos ainda é um problema difícil para os classificadores que nós utilizamos neste estudo. O segundo ponto que precisa ser melhorado é a infraestrutura para suportar muitas consultas, pois o Twitter impõe limitações de uso de sua API.

## **11. Possíveis trabalhos futuros**

A Inteligência artificial e a análise de sentimentos baseadas em aprendizado de máquina é crucial para empresas, visto que, o insight revelado pela análise visa indicar o grau de qualidade dos serviços e/ou produtos da empresa de acordo com os clientes.

Esse projeto é diretamente aplicável a praticamente qualquer empresa que disponha de meios online (Twitter, Instagram, Facebook, Website) para interagir com seus clientes.

Os algoritmos podem ser usados para detectar e possivelmente sinalizar automaticamente tweets de ódio e racismo.

## 12. Bibliografia

- **Pandas Disponível em:**

<https://pandas.pydata.org/docs/>

- **Numpy Disponível em:**

<https://numpy.org/doc/stable/>

- **Seaborn Disponível em:**

<https://seaborn.pydata.org/introduction.html>

- **Matplotlib Disponível em:**

<https://matplotlib.org/3.3.3/contents.html>

- **Scikit-Learn Disponível em:**

<https://scikit-learn.org/stable/tutorial/index.html>

- **Natural Language Toolkit (NLTK) Disponível em:**

<https://www.nltk.org/>

- **pickle Disponível em:**

<https://docs.python.org/3/library/pickle.html>

- **python Disponível em:**

<https://www.python.org/doc/>

- **Bag of Words (BoW) Disponível em:**

<https://www.mygreatlearning.com/blog/bag-of-words/>

- **Naive Bayes Disponível em:**

[https://scikit-learn.org/stable/modules/naive\\_bayes.html](https://scikit-learn.org/stable/modules/naive_bayes.html)