Classificação de sentimento no E-Commerce



com Regressão Logística













Contexto

Analise de Sentimento

Análise de sentimento é um termo que se refere ao uso de **processamento de linguagem natural**, análise de texto e linguística computacional para determinar a atitude de um falante ou escritor em relação a um tópico específico.

Basicamente, ajuda a determinar se um texto expressa sentimentos positivos, negativos ou neutros. A análise de sentimento é uma excelente maneira de descobrir como as pessoas, principalmente os consumidores, se sentem sobre um determinado tópico, produto ou ideia.



Origem dos Dados

Esses dados foram disponibilizados no site da Kaggle. Nessa base de dados há milhares de comentários feitos por clientes do Brasil inteiro em sites E-commerce.

Vamos minerar esses dados e criar um modelo de classificação de emoção.

https://www.kaggle.com/olistbr/brazilian-ecommerce



@Odemir Depieri Jr

Vamos utilizar uma base de dados da Kaggle

https://www.kaggle.com/olistbr/brazilian-ecommerce

Vamos importar as bibliotecas necessárias

```
[1] # Lib para Modelagem de Dados
       import pandas as pd
       # Lib para matrizes
       import numpy as np
       # Lib para plotagem gráfica
       import matplotlib.pyplot as plt
       # Lib para expressões regulares
       import re
       # Lib para trabalhar com palavras
       import string
       # Lib para processamento de linguageum
       import nltk
       # Função para extrair o radical da Palavra
       from nltk.stem import RSLPStemmer
       # Função para remover palavras sem relevânica
       from nltk.corpus import stopwords
       # Libs para uso de ML
       # Função do modelo logistico
       from sklearn.linear_model import LogisticRegression
       # Funções Estimadores
       from sklearn.base import BaseEstimator, TransformerMixin
       # Função de transformadores
       from sklearn.pipeline import Pipeline
       # Função para matriz
       from sklearn.feature_extraction.text import CountVectorizer
       # Funções para avaliação do modelo
       from sklearn.model_selection import cross_val_score
       # Ignorar avisos
       warnings.simplefilter("ignore")
```

Vamos instalar alguns pacotes e ler os dados

e64fb393e7b32834bb789ff8bb30750e 658677c97b385a9be170737859d3511b

[2] # Download de alguns pacores da biblioteca NLTK

f7c4243c7fe1938f181bec41a392bdeb

```
# Pacote Radical da palavra
     nltk.download('rslp')
     # Pacote palavras sem relevância
     nltk.download('stopwords')
     # Grupo de Palavras sem relevância
     Palavras_Sem_Relancia = stopwords.words('portuguese')
     [nltk_data] Downloading package rslp to /root/nltk_data...
[nltk_data] Package rslp is already up-to-date!
[nltk_data] Downloading package stopwords to /root/nltk_data...
[nltk_data] Package stopwords is already up-to-date!
[3] # Lendo a Base de Dados
     Base = pd.read_csv('olist_order_reviews_dataset.csv')
      # Verificando as primeiras linhas
     Base.head()
                                     review id
                                                                                order_id review_score review_comment_title review_comment_message review_creation
      0 7bc2406110b926393aa56f80a40eba40 73fc7af87114b39712e6da79b0a377eb
                                                                                                         4
                                                                                                                                NaN
                                                                                                                                                             NaN
                                                                                                                                                                        2018-01-18 0
          80e641a11e56f04c1ad469d5645fdfde a548910a1c6147796b98fdf73dbeba33
                                                                                                                                                                        2018-03-10 0
      2 228ce5500dc1d8e020d8d1322874b6f0 f9e4b658b201a9f2ecdecbb34bed034b
                                                                                                         5
                                                                                                                                                             NaN
                                                                                                                                                                        2018-02-17 0
                                                                                                                               NaN
                                                                                                                                            Recebi bem antes do
```

Nessa base de dados vamos utilizar apenas o Comentário e Avaliação do cliente.

8e6bfb81e283fa7e4f11123a3fb894f1

NaN

2017-04-21 0

2018-03-01 0

prazo estipulado. Parabéns lojas lannister

adorei comprar pela I...

@Odemir Depieri Jr

Vamos filtrar apenas as colunas essenciais

	Comentário	Pontuação
0	NaN	4
1	NaN	5
2	NaN	5
3	Recebi bem antes do prazo estipulado.	5
4	Parabéns lojas lannister adorei comprar pela I	5

Retirando os valores nulos

```
[5] # Retirando os valores Nulos
Base_Filtrada = Base_Filtrada.loc[
    ( Base_Filtrada['Comentário'].notnull() ) &
    ( Base_Filtrada['Pontuação'].notnull() ) ]

# Verificando os ajustes de registros
print('Antes:', len(Base) )
print('Depois:', len(Base_Filtrada) )
Antes: 100000
```

Antes: 100000 Depois: 41753

Vamos criar a classe, a variável que vamos tentar prever

```
[6] # Vamos utilizar a função 'CUT' do pandas para gerar nossa variavel de classificação
      # Use cut quando precisar segmentar e classificar os valores dos dados em compartimentos
      # Função CUT ira ter 2 Parametros
      # Pontos [0 a 5]
      # Classe [0 ou 1]
      # Bassicamente o pandas irá classificar entre 0 e 1 nossas frases com base na pontuação do cliente
# Assim vamos ter as classes para classificar nosso moldeo
      # Lista com escala dos pontos
      Pontos = [0, 2, 5]
      # Classificação [ 0 - Comentarios Ruins ]
      # Classificação [ 1 - Comentarios Bons ]
      Label = [0, 1]
      # Aplicando a função para gerar nossa Classe
      Base_Filtrada['Classe'] = pd.cut( Base_Filtrada['Pontuação'],
                                             bins=Pontos, # Modelo de score
labels=Label ) # Classificação
      # Verificando
      print( 'Registros Superiores', '\n', Base_Filtrada.head(10), '\n' )
print( 'Registros Inferiores', '\n', Base_Filtrada.tail(10), '\n' )
```

Para gerar a classe, usei o **seguinte racional**. Classificação **maior** que 3 será considerado 'Positivo = 1', **menor** que isso será 'negativo = 0'

Vamos criar 3 Funções para tratar nossos dados

```
[7] # Classe com funções para limpar os Textos
     class Funcao_Limpeza(BaseEstimator, TransformerMixin):
          # Função para ativar na chamada
         def fit(self, x, y=None):
    return self
          # Função para transformar os dados
          def transform(self, x, y=None):
              # Lista que ira receber os dados Transformado
              Texto_Transformado = []
              # Loop no Comentário
              for Texto in x:
                   # retirando os paragrafos, numeros e afins
                  Texto = re.sub('\n', ' ', Texto)
Texto = re.sub('\r', ' ', Texto)
                  Texto = re.sub(r'\d+(?:\.\d*(?:[eE]\d+))?', ' numero ', Texto)
Texto = re.sub(r'\d*(?:[eE]\d+))?', ' numero ', Texto)
Texto = re.sub(r'\d*(', Texto)
Texto = re.sub(r'\d*(', Texto)
Texto = re.sub(r'\s+', ' ', Texto)
                  # Caso não exista uma URL vai passar
                  if len(URL_Texto) == 0:
                       pass
                   # No caso de ter URL no Texto
                       # Loop no Texto
                       for url in URL_Texto:
                            # Loop na frase
                            for link in url:
                                # Retirando o Link
                                Texto = Texto.replace(link, '')
                       # Retirando possveis ':' e
Texto = Texto.replace(':',
                       Texto = Texto.replace('/', '')
                   # Salvando a info na Lista
                   Texto_Transformado.append(Texto)
               # Retorno da Função
              return Texto_Transformado
```

Essa função irá retirar os 'resíduos' das frases.

Caracteres especiais, espaços, link e afins, não agregam em nada nos modelos ao contrário, podem piorar.

Assim vamos fazer essa limpeza.

Vamos Extrair o Radical da Palavra

```
[8] # Classe com funções para Extrair o Radical da Palabras
     class Funcao_Radical(BaseEstimator, TransformerMixin):
         # Função para ativar na chamada
         def fit(self, x, y=None):
             return self
         # Extraindo o Radical da Palavra
         def Aplicando_Radical(self, Loop):
             # Função do Radical
             Radical = RSLPStemmer()
              # Retornando o radical da palavra
              return list(map(lambda x: Radical.stem(x), [Palavra for Palavra in Loop.split()]))
         def transform(self, x, y=None):
             # Retornando aplicando a função do Radical
             Texto_Transformado = list(map(lambda c: self.Aplicando_Radical(c), x))
             # Unificando as palavrvas
             Texto_Transformado = list(map(lambda x: ' '.join(x), Texto_Transformado))
              # Retornando o comentario tratado
             {\tt return\ Texto\_Transformado}
```

Caso não saiba o que é o Radical da palavra.

https://querobolsa.com.br/enem/portugues/radical

Vamos retirar pontuações e palavras sem relevâncias

```
( [9] # Retirar palavras sem relevancia e pontuações
      class Funcao_Palavras(BaseEstimator, TransformerMixin):
        def fit(self, x, y=None):
        # Revomendo Stops Words e Pontuações
        def Remover_Palavras(self, Texto):
            # Verificando se existe pontuações na Palavraas
            Sem Pontuacao = [Palavra for Palavra in Texto if Palavra not in string.punctuation]
            # Unificando o texto
            Sem_Pontuacao = ''.join(Sem_Pontuacao)
            # Retornando a frase sem as stops words
            return [Sem_Relancia for Sem_Relancia in Sem_Pontuacao.split() if Sem_Relancia.lower() not in Palavras_Sem_Relancia]
        # Aplicando a transformação na frase
        def transform(self, x, y=None):
            # Removendo palavras e pontuações com a função 'Remover_Palavras'
            Texto_Transformado = list(map(lambda c: self.Remover_Palavras(c), x))
            Texto_Transformado = list(map(lambda x: ' '.join(x), Texto_Transformado))
            # Retorno da Função
            return Texto Transformado
```

Explicação do que seria as Stops Words

https://virtuati.com.br/cliente/knowledgebase/25/Lista-de-StopWords.html

Vamos aplicar essas 3 funções nos dados

O que aconteceu ate aqui?

Veja que as palavras ficaram com uma **certa abreviação**, certo? Esse processo é referente ao radical da palavra.

Outro detalhe que as frases **ficaram mais curtas**, isso se deve ao processo da remoção das stops words.

Quando treinamos um modelo de NLP (Linguagem de processamento natural) é essencial fazer esse tipo de treinamento para melhorarmos a performance do modelo.

Toda vez que for trabalhar com NLP é primordial fazer esse tipo de tratamento.

Vamos aplicar a vetorização nas Frases

É denominado **vetorização** o processo geral de converter a coleção de textos em **vetores numéricos**. A estratégia de *tokenizar* (separar textos ou palavras em "blocos"), *contar* e *normalizar* (como já fizemos anteriormente) é chamado de representação **Bag of Words (BOW)**, ou na tradução literal **bolsa de palavras**, no qual são descritos pelas ocorrências das palavras enquanto são ignoradas suas posições relativas em todos os textos. Segue exemplo abaixo de como fica a estrutura, com a frequencia absoluta de cada palavra:

	- 1	love	dogs	hate	and	knitting	is	my	hobby	passion
Doc 1	1	1	1							
Doc 2	1		1	1	1	1				
Doc 3					1	1	1	2	1	1

fig.1 - Imagem adaptada de text-analytics-beginners-nltk bootcamp

Sugiro a leitura caso não entenda sobre Vetorização

https://carlos-bonfim.medium.com/machine-learning-e-processamento-de-linguagem-natural-pln-com-modelos-lineares-d5aaaaf0efa5

Vamos treinar o modelo e medir sua Acuracia

0.8677952823375319

Modelo foi muito bem !!!

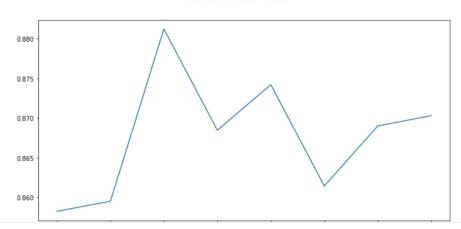


@Odemir Depieri Jr

Vamos plotar essas acurácias

```
# Definindo uma figura
fig = plt.figure( figsize=(12,6) )
# Titulo
fig.suptitle('Acuracia do Modelo - Testes')
# Plotando os Valores
plt.plot(Acuracia)
# Mostrando os gráfico
plt.show()
```

Acuracia do Modelo - Testes



Hora dos Testes

```
[17] # Função testar as frases
def Testando_Modelo(Seu_Texto):
    # Variavel_para receber a classe
    Classificação = ''

# Aplicando a transformação
    Processamento_Texto = Processamento_Geral.fit_transform( Seu_Texto )

# Aplicando a Vetorização
    Texto_Transformado = Vetorizacao.transform( Processamento_Texto ).toarray()

# Fazendo a Previsão
    Previsão = Funcao_Logistica.predict( Texto_Transformado )

# Condição para verificar Positivo ou False
if Previsao[0] == 1:
    Classificação = 'Positivo'
else:
    Classificação = 'Negativo'

return print('Avalição do Comentário: ', Classificação )
```

```
[18] Testando_Modelo(['Excelente produto!'])
```

Avalição do Comentário: Positivo

[19] Testando_Modelo(['produto de péssima qualidade!!!'])

Avalição do Comentário: Negativo



Mandamos muito bem !!!

Final

Esse guia foi elaborado para classificação de comentários.

Link do Colab

https://colab.research.google.com/drive/1CmC4SbPpjdmlcD2WAunoadXq6_rsjH1A?usp=sharing















Odemir Depieri Jr

Data Intelligence Analyst Sr Tech Lead Specialization AI