Atividade Prática 3 - Base de dados Trip Advisor Hotel Reviews

Summary

- 1. Limpeza da base de dados: Remover stopwords, caracteres especiais e aplicar stemming.
- Transformação em atributos numéricos: Utilizar CountVectorizer para converter o texto em um vetor de termos/term frequency, e então aplicar TfidfTransformer para normalizar os dados.
- 3. Avaliação com classificadores: Utilizar árvore de decisão e random forest para avaliar os dados transformados.

```
In [ ]: import pandas as pd
        import numpy as np
         import matplotlib.pyplot as plt
         import seaborn as sns
         import nltk
         from nltk.corpus import stopwords
         from nltk.stem import SnowballStemmer
         from nltk.tokenize import word tokenize
         import re
         from sklearn.model selection import train test split
In [ ]: # Carregar os dados
        df = pd.read csv('tripadvisor hotel reviews.csv')
Out[]:
                                              Review Rating
             nice hotel expensive parking got good deal sta...
         1 ok nothing special charge diamond member hilto...
         2 nice rooms not 4* experience hotel monaco seat...
            unique, great stay, wonderful time hotel monac...
         3
         4 great stay great stay, went seahawk game aweso...
In [ ]: # Função atualizada de limpeza de dados
        def Limpeza Dados(Texto):
             # Remoção de links
             etapa_01 = re.sub('www\S+|http\S+', '', Texto)
             # Transformar para minúsculo
             etapa 02 = etapa 01.lower()
             # Remoção de caracteres especiais
             etapa 03 = re.sub(r'[!~@#$%^&*()+=|{}[\]:;<.>?/\'\\",-]', '', etapa
             # Remoção de números
             etapa 04 = re.sub('[0-9]', '', etapa 03)
             # Remove espaços extras
```

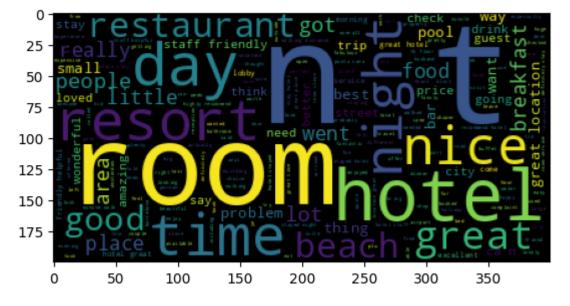
1 of 9 2/28/24, 15:46

etapa $05 = re.sub(r'\s+', ' ', etapa 04)$

```
return etapa 05
         # Aplicar a limpeza nos dados
         df['Review cleaned'] = df['Review'].apply(Limpeza Dados)
In [ ]: df.head(5)
Out[]:
                                    Review Rating
                                                                      Review_cleaned
                nice hotel expensive parking got
                                                     nice hotel expensive parking got good
         0
                              good deal sta...
                                                                            deal sta...
              ok nothing special charge diamond
                                                        ok nothing special charge diamond
         1
                                                 2
                              member hilto...
                                                                        member hilto...
              nice rooms not 4* experience hotel
                                                          nice rooms not experience hotel
         2
                                                 3
                              monaco seat...
                                                                      monaco seattle...
                                                    unique great stay wonderful time hotel
              unique, great stay, wonderful time
         3
                               hotel monac...
                                                                            monaco...
             great stay great stay, went seahawk
                                                        great stay great stay went seahawk
                                                 5
                               game aweso...
                                                                       game awesom...
In [ ]: exemplo = df['Review'][10]
         print(exemplo)
         Limpeza Dados(exemplo)
        poor value stayed monaco seattle july, nice hotel priced 100- 150 night no
        t, hotel takes beating quotient, experience simply average, nothing except
        ional paying 300+ n't ca n't terribly disappointed, wife stayed nicest sui
        tes 200/night felt like overpaying, best advice shop, quality-wise league
       average marriott nice best western,
Out[ ]: 'poor value stayed monaco seattle july nice hotel priced night not hotel
         takes beating quotient experience simply average nothing exceptional pay
         ing n t ca n t terribly disappointed wife stayed nicest suites night fel
         t like overpaying best advice shop quality wise league average marriott
         nice best western '
In [ ]: df.drop('Review', axis=1, inplace=True)
In [ ]: df.shape
Out[]: (20491, 2)
In [ ]: df.rename(columns ={"Rating":"class","Review cleaned":"text"}, inplace =
In [ ]: df.columns
Out[]: Index(['class', 'text'], dtype='object')
In [ ]: | df['class'].isnull().sum()
Out[]: 0
In [ ]: df['class'].value counts()
```

```
Out[]: class
         5
              9054
         4
              6039
         3
              2184
         2
              1793
              1421
        Name: count, dtype: int64
In [ ]: # Reduzir target para variável binaria
        def change_rating(x):
            if x > 3:
                return 'Positive Rating'
            return 'Negative Rating'
In [ ]: # ALtera na Base de Dados (Dataframe)
        df['class'] = df['class'].apply(change_rating)
In [ ]: | df['class'].value counts()
Out[]: class
         Positive Rating
                            15093
        Negative Rating
                             5398
        Name: count, dtype: int64
In [ ]: from wordcloud import WordCloud
In [ ]: # Concatenar todas as palavras em uma única string
        all words = ''.join(df['text'])
In [ ]: # Objeto WorCloud()
        Nuvem_Palavras = WordCloud().generate(all_words1)
        Fig, eixo = plt.subplots()
        eixo.imshow(Nuvem Palavras)
Out[]: <matplotlib.image.AxesImage at 0x7f7e253a8e10>
         50
         75
        100
                          valet ann
        125 -
        150 -
        175 -
            0
                   50
                           100
                                   150
                                           200
                                                   250
                                                           300
                                                                   350
```

```
In [ ]: df['class'].unique()
Out[ ]: array(['Positive Rating', 'Negative Rating'], dtype=object)
In []: # Avaliação por review positivo e negativo
        review_positive = ''
        review negative = ''
        # Verifying each class (positive and negative)
        for c in df['class'].unique():
            # Filter the dataframe for the current class
            filter df = df[df['class'] == c]
            # Iterate over the filtered dataframe
            for text in filter df['text']:
                # If the class is 'Positive Rating', concatenate to review positi
                if c == 'Positive Rating':
                    review positive += ' ' + ' '.join(text.split())
                # If the class is 'Negative Rating', concatenate to review negati
                elif c == 'Negative Rating':
                    review negative += ' ' + ' '.join(text.split())
        # Trim the leading spaces
        review positive = review positive.strip()
        review_negative = review_negative.strip()
In [ ]: print(
            'Quantidade de Palavras: \n',
            f'Best review: { len(review positive) }',
            f'Worst_review: { len(review_negative) }'
       Quantidade de Palavras:
        Best review: 10177358 Worst review: 4223226
In []: # Classe Positiva
        # Instanciando a Nuvem de Palavras
        Nuvem Palavras = WordCloud().generate(review positive)
        # Tamanho
        Figura, Eixo =plt.subplots()
        # Plot
        Eixo.imshow(Nuvem Palavras)
Out[]: <matplotlib.image.AxesImage at 0x7f7dff3e5190>
```

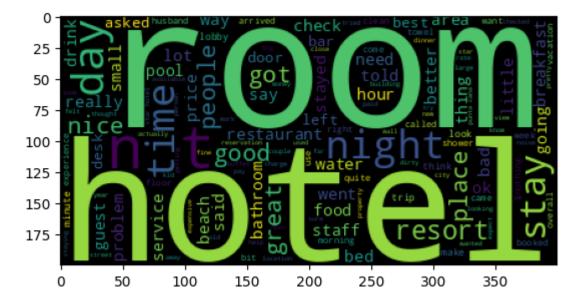


```
In []: # Classe Negativo

# Instanciando a Nuvem de Palavras
Nuvem_Palavras = WordCloud().generate(review_negative)

# Tamanho
Figura, Eixo = plt.subplots()
# Plot
Eixo.imshow(Nuvem_Palavras)
```

Out[]: <matplotlib.image.AxesImage at 0x7f7dfefa5350>



```
In []: df.shape
Out[]: (20491, 2)
In []: df.drop_duplicates(['text'], inplace=True)
In []: #no duplicates
df.shape
Out[]: (20491, 2)
```

StopWords

```
In [ ]: | nltk.download('stopwords')
       [nltk data] Downloading package stopwords to /home/a10/nltk data...
       [nltk_data] Package stopwords is already up-to-date!
Out[]: True
In [ ]: df['text'][1000]
Out[ ]: 'shame hotel wasnt good restaurant arrived clift late afternoon struggle
        luggage bags reception staff unhelpful uninterested eventually managed s
        orted shown room th floor room suite tried make separate living room put
        ting curtain inbetween bedroom living room bathroom tiny dirty stayed mu
        m unfortunatley night didnt feel suffering bad foot decided phone recept
        ion ask doctor come hotel told ther wasnt local receptionist closest tol
        d phone eventually decided hospital just safe came hospital evening door
        men talking girls outside let following night ate hotel restaurant aisa
        cuba fantastic think hotel intrest restaurant bar end day sleeping ignor
        ed wouldnt stay '
In [ ]: lista stopwords = nltk.corpus.stopwords.words('english')
In [ ]: def Remover StopWords(Texto):
          Lista Palavras = Texto.split()
          # Texto sem as stopwords
          nova frase = ''
          for word in Lista Palavras:
            if word not in lista stopwords:
              nova frase = nova frase + ' ' + word
          return nova frase
In [ ]: | exemplo = df['text'][10]
        print(len(exemplo))
        saida = Remover_StopWords(exemplo)
        print(len(saida))
       306
       298
In [ ]: | df['text'] = df['text'].apply(Remover StopWords)
```

Extração do Radical

```
radical = Stem.stem(word)
  nova_frase = nova_frase + ' ' + radical
  return nova_frase
```

```
In [ ]: exemplo = df['text'][10]
    print(exemplo)
    Extrair_Radical(exemplo)
```

poor value stayed monaco seattle july nice hotel priced night hotel takes beating quotient experience simply average nothing exceptional paying n can terribly disappointed wife stayed nicest suites night felt like overpaying best advice shop quality wise league average marriott nice best western

Out[]: 'po valu stayed monac seattl july nic hotel priced night hotel tak beat ing quotient experienc simply averag nothing except paying n ca n terrib ly disappointed wif stayed nicest suit night felt lik overpaying best ad vic shop quality wis leag averag marriott nic best western'

```
In [ ]: df['text'] = df['text'].apply(Extrair_Radical)
```

Tokenização

Modelo

```
In [ ]: # Separar entre previsores e classe
    X = df['text']
    y = df['class']

In [ ]: # Dividir os dados em conjuntos de treino e teste
    X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.3,

In [ ]: X_train.shape

Out[ ]: (14343,)

In [ ]: X_test.shape

Out[ ]: (6148,)
```

```
In [ ]: from sklearn.feature extraction.text import CountVectorizer, TfidfTransfo
        # Inicializar CountVectorizer
        vect = CountVectorizer()
        vect.fit(X_train)
        X train vect = vect.transform(X train)
        X test vect = vect.transform(X test)
        # Inicializar TfidfTransformer
        Tfidf = TfidfTransformer()
        X train vect = Tfidf.fit transform(X train vect)
        X test vect = Tfidf.fit transform(X test vect)
In [ ]: X train vect.shape
Out[]: (14343, 35046)
In [ ]: X test vect.shape
Out[]: (6148, 35046)
In [ ]: vect.get feature names out()
Out[]: array(['__', '___', '___', ..., 'üè', 'üè__', 'üèc'], dtype=object)
In [ ]: from sklearn.tree import DecisionTreeClassifier
        from sklearn.ensemble import RandomForestClassifier
        from sklearn.metrics import classification report, accuracy score
        # Inicializar e treinar o classificador de árvore de decisão
        dt clf = DecisionTreeClassifier()
        dt clf.fit(X train vect, y train)
        dt predictions = dt clf.predict(X test vect)
        # Avaliar árvore de decisão
        print("Árvore de Decisão - Relatório de Classificação:\n", classification
        # Inicializar e treinar o classificador Random Forest
        rf clf = RandomForestClassifier()
        rf clf.fit(X train vect, y train)
        rf predictions = rf clf.predict(X test vect)
        # Avaliar Random Forest
        print("Random Forest - Relatório de Classificação:\n", classification rep
```

Arvore de Decisão - Relatório de Classificação:				
р	recision	recall	f1-score	support
Negative Rating	0.56	0.58	0.57	1600
Positive Rating	0.85	0.84	0.84	4548
accuracy			0.77	6148
macro avg	0.70	0.71	0.71	6148
weighted avg	0.77	0.77	0.77	6148
Random Forest - Rel	atório de	Classific	ação:	
	recision		f1-score	support
Negative Rating	0.96	0.41	0.57	1600
Positive Rating	0.83	0.99	0.90	4548
accuracy			0.84	6148
macro avg	0.89	0.70	0.74	6148
weighted avg	0.86	0.84	0.82	6148

In []: