## ▼ Chapter 6 - Exercise 5: WordClouds

## ▼ Part 1:

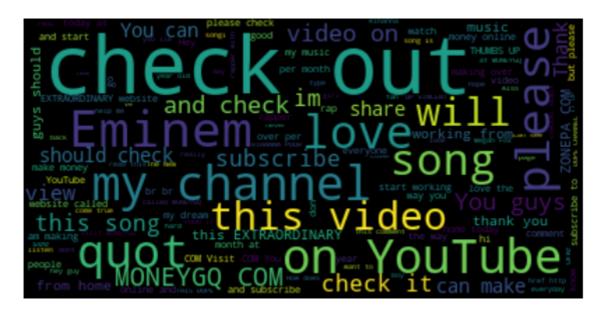
'am', 'an', 'and', 'any', 'are',

```
import pandas as pd
import matplotlib.pyplot as plt
from wordcloud import WordCloud, STOPWORDS
# Tạo dataframe df từ file dữ liệu Youtube04-Eminem.csv, in vài dòng dữ liệu đầu của df
df = pd.read_csv("data\Youtube04-Eminem.csv", encoding ="latin-1")
df.head()
df['CONTENT'].head(20)
                 +447935454150 lovely girl talk to me xxxi»¿
     1
           I always end up coming back to this song<br/>to />i»¿
     2
           my sister just received over 6,500 new <a rel=...
     3
                            Hello I'am from Palastine
     5
           Wow this video almost has a billion views! Did...
     6
           Go check out my rapping video called Four Whee...
     7
                                          Almost 1 billioni»¿
     8
                            Aslamu Lykum... From Pakistani»¿
     9
           Eminem is idol for very people in España and ...
     10
                              Help me get 50 subs please i»¿
     11
                                            i love song :) i»¿
           Alright ladies, if you like this song, then ch...
     12
     13
           The perfect example of abuse from husbands and...
           The boyfriend was Charlie from the TV show LOS...
           <a href="https://www.facebook.com/groups/10087...</pre>
     15
                    Take a look at this video on YouTube:
     16
     17
                   Check out our Channel for nice Beats!!i»¿
     18
             Rihanna and Eminem together are unstoppable.i»¿
     19
                      Check out this playlist on YouTube:
     Name: CONTENT, dtype: object
# Hãy bỏ các STOPWORD
stopwords = set(STOPWORDS)
stopwords
     {'a',
      'about',
      'above',
      'after',
      'again',
      'against',
      'all',
      'also',
```

```
"aren't",
      'as',
      'at',
      'be',
      'because',
      'been',
      'before',
      'being',
      'below',
      'between',
      'both',
      'but',
      'by',
      'can',
      "can't",
      'cannot',
      'com',
      'could',
      "couldn't",
      'did',
      "didn't",
      'do',
      'does',
      "doesn't",
      'doing',
      "don't",
      'down',
      'during',
      'each',
      'else',
      'ever',
      'few',
      'for',
      'from',
      'further',
      'get',
      'had',
      "hadn't",
      'has',
      "hasn't",
      'have',
      "haven't",
      'having',
      'he',
      "he'd",
      "he'll",
comment_words =' '.join(df['CONTENT'])
# Vẽ WordClouds cho dữ liệu này
wc = WordCloud(
    background_color='black',
    stopwords=stopwords
# generate the word cloud
wc.generate(comment_words)
# display the word clouds
```

)

```
plt.figure(figsize=(10, 12))
plt.imshow(wc, interpolation='bilinear')
plt.axis('off')
plt.show()
```



## ▼ Part 2:

```
# Sử dụng vietnamese stop word từ : https://github.com/stopwords/vietnamese-stopwords => t
stopwords = set()
f = open(r"data\vietnamese-stopwords.txt", "r", encoding='utf-8')
for line in f:
    word = f.readline()
    stopwords.add(word.replace('\n',''))
f.close()
stopwords
     { 'a ha',
      'ai ai',
      'ai đó',
      'amen',
      'anh ấy',
      'ba ba',
      'ba cùng',
      'ba ngày',
      'ba tăng',
      'bao lâu',
      'bao na',
      'biết',
      'biết bao nhiêu',
      'biết chừng nào',
      'biết mấy',
      'biết trước',
      'biết đâu',
      'biết đâu đấy',
      'buổi',
      'buổi mới',
      'buổi sớm',
```

```
'bà ấy',
      'bài bác',
      'bài cái',
      'bán',
      'bán dạ',
      'bây bẩy',
      'bây giờ',
      'bèn',
      'bên',
      'bên có',
      'bông',
      'bước khỏi',
      'bước đi',
      'bản',
      'bản riêng',
      'bản ý',
      'bất cứ',
      'bất kì',
      'bất kỳ',
      'bất ngờ',
      'bất quá',
      'bất thình lình',
      'bất đồ',
      'bấy chầy',
      'bấy giờ',
      'bấy lâu nay',
      'bấy nhiêu',
      'bập bõm',
      'bắt đầu từ',
      'bằng cứ',
      'bằng người',
      'bằng như',
      'bằng nấy',
      'bằng được',
      'bển',
      'bị',
      'bị vì',
      'bỏ bà',
# Đọc file ngon_tu_quang_cao.txt => đưa nội dung vào biến text
text = ''
with open(r'data\ngon_tu_quang_cao.txt', 'r', encoding='utf-8') as f:
    text = f.read()
text
     '\ufeffNGÔN Từ TRONG QUẢNG CÁO\nGiới thiệu\nViệt Nam đang từng bước hội nhập vào nền
\# Bổ sung thêm một số từ không quan trọng vào stopwords
list_of_words = ['và', 'một', 'của', 'có', 'đó', 'rất', 'nào', 'được',
                 'khi', 'thể', 'sự', 'tính', 'trong','cũng','cùng','cho','hay','chỉ']
for word in list_of_words:
    stopwords.add(word)
# Vē wordclouds
wc = WordCloud(
```

```
background_color='black',
    max_words=1000,
    stopwords=stopwords
)

wc.generate(text)

plt.figure(figsize=(10, 12))
plt.imshow(wc, interpolation='bilinear')
plt.axis('off')
plt.show()
```



```
import numpy as np
from PIL import Image

# Chon hinh lam wc_mask phu hop
wc_mask = np.array(Image.open('ad_s.png'))

plt.imshow(wc_mask, interpolation='bilinear')
plt.axis('off')
plt.show()
```

