# WARM UP PROJECT (TASK 2)

**Question:**

The word cloud for Top 5 news category:

- Category that has the most of news articles
- URL contains the category which the new belong

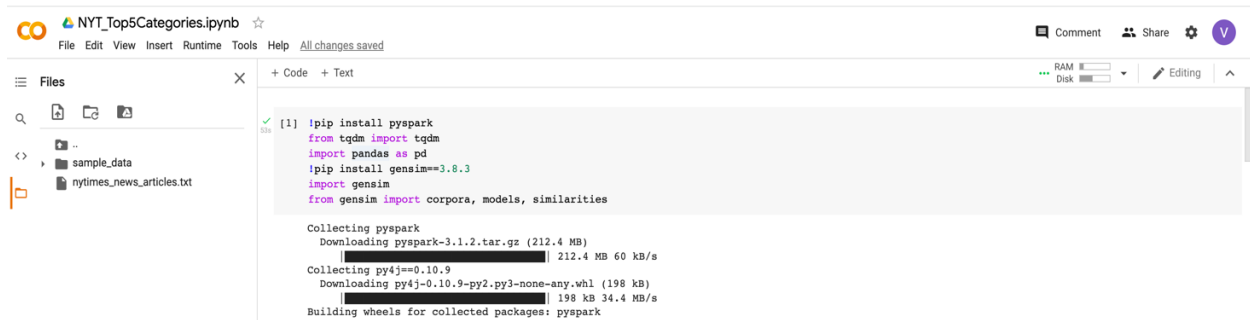For example:-

- -/us/politics/
- -/world/
- -/sports/
- -/arts/
- -and so on

**Step 1:**

1. Install and import all libraries

   <span style="color:red">I have used pyspark, tqdm, pandas, genism</span>



**Step 2:**

1. Import Sparkcontext and read the textfile (news article)

   <span style="color:red">'content/nytimes.new_articles.txt' → sc.textFile(newsarticle) using sparkcontext</span>

2. Ouptut of newsData collection



## Step 3: TOP 5 Catogories

1. Extracting Top 5 categories based on the URL string
2. For each value in newsdata → If 'URL: ' is present then split the url with '/' position 4 to extract words
3. Based on this 4th position value assigned the dictionary values in categories as word count and its corresponding article



```python
[4]  from tqdm import tqdm
     categories = {}
     for i, x in enumerate(tqdm(newsData)):
         # if i == 20000:
         #     break
         if 'URL:' in x:
             category = x.split('.com')[1].split('/')[4]
             if category == 'us':
                 category = f"{category}/{x.split('.com')[1].split('/')[5]}"
             if category not in categories:
                 categories[category] = {'count': 1, 'article': ''}
             else:
                 categories[category]['count'] +=1
         else:
             categories[category]['article'] += (', '+x)

     100%|████████| 192577/192577 [03:57<00:00, 810.87it/s]
```

4. With this key-value pair ('count', 'article'), with the help of reverse sort I got Top 5 categories which have highest value of frequency count.
5. I got 'sports' → 1268 , 'world' → 1211 , 'business' → 1041 , 'nyregion' → 663, 'arts' → 663
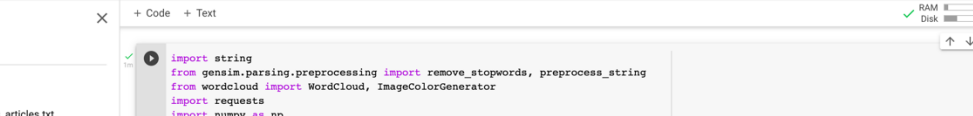


```python
[5]  cat_count = [(key, categories[key]['count']) for key in categories]
     cat_count.sort(reverse=True,key = lambda x: x[1])
     top5 = cat_count[:5]
     top5

     [('sports', 1268),
      ('world', 1211),
      ('business', 1041),
      ('nyregion', 663),
      ('arts', 663)]
```

## Step 4: WordCloud of Top 5 Categories Articles
1. Generated word cloud based on the top 5 category article using matplotlib and wordcloud
2. Preprocessed the text file → cleaned the text by removing stopwords and getting tokenized words
3. Removing special characters and punctuations
4. Plot Word Cloud



```python
import string
from gensim.parsing.preprocessing import remove_stopwords, preprocess_string
from wordcloud import WordCloud, ImageColorGenerator
import requests
import numpy as np
from PIL import import *
from PIL import Image as ImPIL
import matplotlib.pyplot as plt

# combining the image with the dataset
Mask = np.array(Image.open(requests.get('http://clipart-library.com/image_gallery2/Twitter-PNG-Image.png', stream=True).raw))

wordclouds = []
for x, _ in top5:
    mycat = categories[x]['article']
    # break
    sen = remove_stopwords(mycat)
    text_tokens = preprocess_string(sen)
    no_punctuation = [char.strip() for char in text_tokens if (char.strip() not in string.punctuation) and char.strip().isalnum()]
    filtered_sentence = ','.join(no_punctuation)
    # We use the ImageColorGenerator library from Wordcloud
    # Here we take the color of the image and impose it over our wordcloud
    # image_colors = ImageColorGenerator(Mask)
    filtered_sentence
    #Generating the wordcloud :
    wordcloud = WordCloud(background_color='black', height=1500, width=4000,mask=Mask).generate(filtered_sentence)
    wordclouds.append(wordcloud)
```

## WordCloud of SPORTS

```python
[8] # Sports
    #Plot the wordcloud :
    plt.figure(figsize = (12, 12))
    plt.imshow(wordclouds[0])
    #To remove the axis value :
    plt.axis("off")
    plt.show()
```

# WordCloud of World

```
# World
#Plot the wordcloud :
plt.figure(figsize = (12, 12))
plt.imshow(wordclouds[1])
#To remove the axis value :
plt.axis("off")
plt.show()
```



# WordCloud of Business

```
# business
#Plot the wordcloud :
plt.figure(figsize = (12, 12))
plt.imshow(wordclouds[2])
#To remove the axis value :
plt.axis("off")
plt.show()
```

# WordCloud of nyregion

```
# nyregion
#Plot the wordcloud :
plt.figure(figsize = (12, 12))
plt.imshow(wordclouds[3])
#To remove the axis value :
plt.axis("off")
plt.show()
```



# WordCloud of arts

NYT_Top5Categories.ipynb
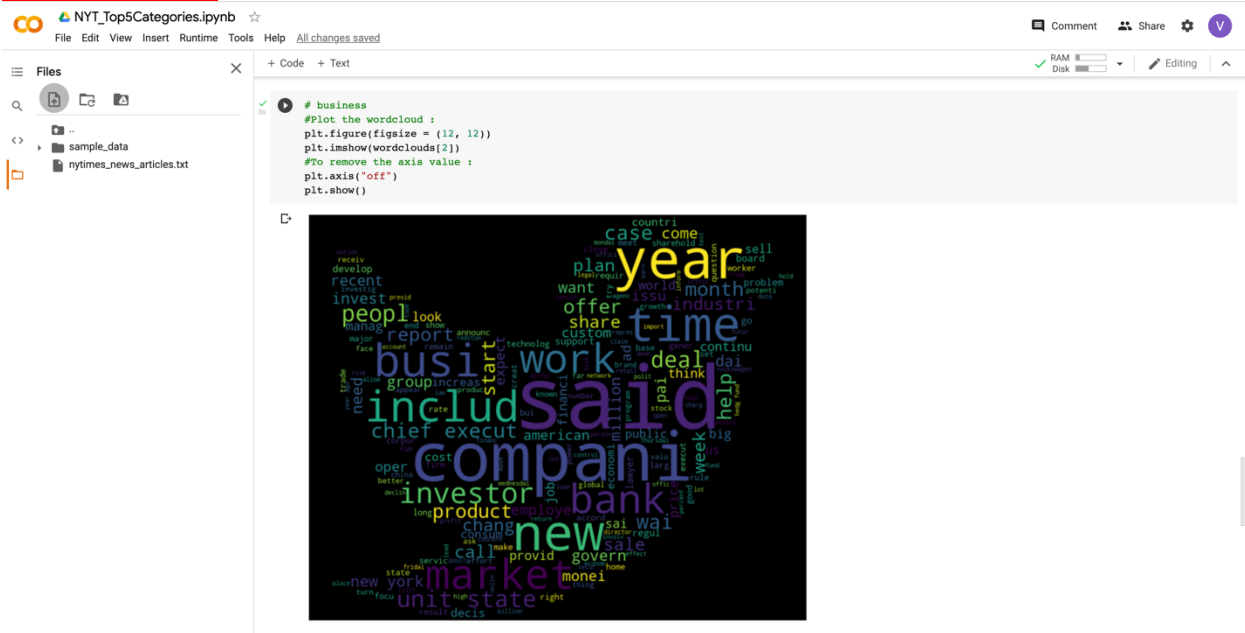File  Edit  View  Insert  Runtime  Tools  Help   All changes saved

Files

```
# arts
# business
#Plot the wordcloud :
plt.figure(figsize = (12, 12))
plt.imshow(wordclouds[4])
#To remove the axis value :
plt.axis("off")
plt.show()
```