# STA323 Project 1 report

SID: 12110821
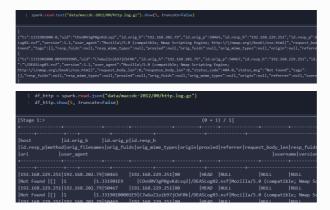
Name: ZHANG Chi

## Solution for Q1

The data given is about an an existing network infrastructure, which is stored in multiple directories.

### (1)

Before we read the whole data, we need to know the structure of the data. I use `spark.read.text` to read one `http.log.gz` of the first directory, and we can find that each record is in `json` format. So we can use `spark.read.json` to read the data.



Given a path list, the function can also read the data from the each path from path list and concat them together to a single dataframe.

```
1  folders = ["00", "01", "02", "03", "04", "05"]
2  df_http = spark.read.json([f"data/maccdc-2012/{folder}/http.log.gz" for folder in
   folders])
3  df_dns = spark.read.json([f"data/maccdc-2012/{folder}/dns.log.gz" for folder in folders])
```

After checking the column `ts`, we can feel free to convert it to a `Timestamp` data type by `to_timestamp` with `withColumn` function. And the temp view is created for further analysis.

```
1  df_http = df_http.withColumn("ts", to_timestamp(col("ts")))
2  df_dns = df_dns.withColumn("ts", to_timestamp(col("ts")))
3  df_http.createOrReplaceTempView("http_log")
4  df_dns.createOrReplaceTempView("dns_log")
```

# (2)

In this task, we need to find the top uri that has been accessed most among the records whose `uid` is `status_code` is 200 and `method` used is `GET`.

We can use `filter` to filter the records and then group by the `uri` and count the records. Finally, we can sort the result by the count in descending order and show the top records. The code is as follows:

```
# In Spark SQL API
spark.sql("SELECT uri, COUNT(uri) AS uri_count FROM http_log WHERE status_code = 200 AND
method = 'GET' GROUP BY uri order by uri_count desc").show(5)

# In Spark DataFrame API
df_http.select(col("uri")).filter((col("status_code")==200) &
(col("method")=='GET')).groupBy(col("uri")).agg(count(col("uri")).alias("uri_count")).orde
rBy(col("uri_count").desc()).show(5)
```

And they get the same result:

```
+--------------------+---------+
|                 uri|uri_count|
+--------------------+---------+
|                   /|     9475|
|/admin/config.php...|      556|
|   /main.php?logout=1|      194|
|/top.php?stuff=15...|      191|
|            /top.php|      179|
+--------------------+---------+
only showing top 5 rows
```

# (3)

After inner joining two tables by `uid`, filtering the records in the same way as in the last task, the number of record whose `proto` is `tcp` for each `uri` can be calulated by `SUM(CASE when proto='tcp' THEN 1 ELSE 0 END) Group By uri`. Then results are sorted by the `tcp_ratio` in descending order and top records are showed below.
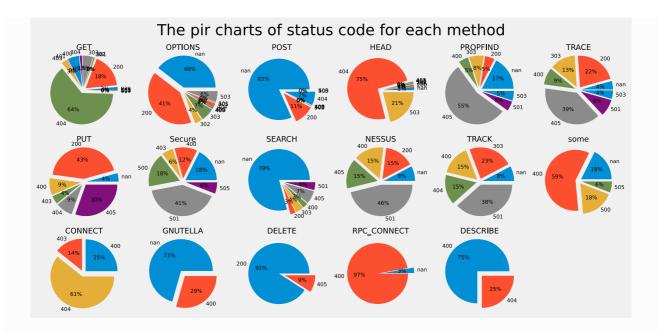
```
spark.sql('''
    SELECT h.uri, COUNT(uri) AS uri_count,
    SUM(CASE when proto='tcp' THEN 1 ELSE 0 END) as tcp_count,
    tcp_count/uri_count as tcp_ratio
    From http_log h
    Join dns_log d on h.uid = d.uid
    Where h.status_code = 200 AND h.method = 'GET'
    Group By uri
    Order By tcp_ratio desc
''').show(5)
```

```
[Stage 17:============================================>          (5 + 1) / 6]
+--------------------+---------+---------+-------------------+
|                 uri|uri_count|tcp_count|          tcp_ratio|
+--------------------+---------+---------+-------------------+
|/lib/exe/indexer....|        1|        1|                1.0|
|   /get_latest_id.php|       18|        3|0.16666666666666666|
|                   /|      133|       22|0.16541353383458646|
| /_vti_bin/shtml.dll|        9|        1| 0.1111111111111111|
|/_vti_bin/_vti_au...|        9|        1| 0.1111111111111111|
+--------------------+---------+---------+-------------------+
only showing top 5 rows
```

## (4)

According to the requirement, we need to calculate the use frequency of each method. Then `groupBy` and `agg` can be used directly. The result are showing below:

```
+-------+------+--------------------+
|method |count |freq                |
+-------+------+--------------------+
|GET    |136678|0.7354368670834992  |
|POST   |40184 |0.21622203329638517 |
|HEAD   |4805  |0.025854739945976778|
|OPTIONS|2388  |0.012849348385222174|
|CONNECT|84    |4.5198712912841813E-4|
+-------+------+--------------------+
only showing top 5 rows
```

To get the number of records in different status code for each method, two columns have been put into `GroupBy`. Before we plot the pie chart, we need to convert the result to a pandas dataframe.

| | method | status_code | count |
|---|---|---|---|
| 0 | None | NaN | 18 |
| 1 | None | 400.0 | 1299 |
| 2 | None | 414.0 | 59 |
| 3 | BXNTPG | 200.0 | 1 |
| 4 | CONNECT | 400.0 | 21 |
| ... | ... | ... | ... |
| 117 | some | NaN | 3 |
| 118 | some | 400.0 | 10 |
| 119 | some | 500.0 | 3 |
| 120 | some | 505.0 | 1 |
| 121 | user | NaN | 1 |

122 rows × 3 columns

It is notable that some values in column `method` and column `status_code` are missing. And I will drop `None` but keep `NaN` as we want to get distribution of status code regarding to a specific method. Besides, the method emerging once (only appearing with one kind of status code) will also be scrapped, and selected methods are shown below.

```
1  method_list = method_count[method_count>1].index
2  method_list

Index(['GET', 'OPTIONS', 'POST', 'HEAD', 'PROPFIND', 'TRACE', 'PUT', 'Secure',
       'SEARCH', 'NESSUS', 'TRACK', 'some', 'CONNECT', 'GNUTELLA', 'DELETE',
       'RPC_CONNECT', 'DESCRIBE'],
      dtype='object', name='method')
```

Then we can use `matplotlib` to plot the pie chart.

The pir charts of status code for each method

# Solution for Q2

## (1)

To scratch the data, I use `request` to get the data from the given website and use `BeautifulSoup` to parse the data. In the main page, there is a list of links to the detailed information of each article. I used `find_all` to get a link list and then request the page related to each link to get the article. Here is a shortcut of the website structure.



Then there is a `for` loop to get the data from each page. In each paragraph of one article in a page, texts are seperated into multiple lines under unknown resons. In this case, I first convert the soup object into a string, replace `'\n'` into `""` and split the string by `<br/><br/>`, which is a kind of paragraph seperator in html. Now each sentence has been grouped into their own paragraph. However, there are also some annoying html tags in the text. I use regular expression to remove them.

Besides, I also replace some special characters in the title of the article to avoid the error of reading files by `spark.sparkContext.textFile`.

```
1  for i in tqdm(article_list):
2      try:
3          href = i['href']
4          if "http" in href:
5              article_page = requests.get(href)
6              split_text = re.split('\n{2,}', article_page.text)
7              text = ("\n").join([i.replace("\n", " ").strip("-").strip() for i in
   split_text if i != ""]).strip("\n")
8          else:
9              article_page = requests.get(f"{suburl}{href}")
```

```
10              article_soup = BeautifulSoup(article_page.text, 'html.parser')
11              article_text = article_soup.find("font", face = "verdana")
12              para_list = str(article_text).replace("\n", " ").split("<br/><br/>")
13              text = ("\n").join([re.sub('<.*?>', '', i).strip() for i in para_list if i !=
     ""]).strip("\n")
14
15          txt_file = i.text.replace("/", "or").replace("?",
     "").replace(",","").replace(":", "")
16
17          with open(f"data/paul_articles/{txt_file}.txt", "w") as f:
18              f.write(text)
19      except Exception as e:
20          print(f"Error with {i.text}: {e}")
21          continue
```

> Two url links are not in the format of previous links, and even their related pages are not consistent with others. I use a `if` statement to deal with them separately.



> Same as the previous step, `\n` and `\n+` (more `\n` like `\n\n`, `\n\n\n`) are used simultaneously in one page. Luckily, I find that `\n+` is used to seperate paragraphs. So I use `re.split` by `\n+` to split the whole page. Then I can let every paragraph to be a single line.



## (2)

In this task, the main problem is to get key phrases in high quality. To finish it, I ask `chatgpt` for help, and my prompts as well as corresponding reponses are as follows:

> In my Pyspark task, I have a rdd that contains texts (each element is a paragraph) written by a author. Now I want to find paragraphs talking about 'career planning', could you help me find some key phrases from raw contents which I will give you to help me filter elemetns in rdd? If Ok I will give you the raw contents. And you will give me the key phrases.



The raw pages I give the gpt are "Ideas for Startups", "Web 2.0", "Why TV lost" and so on. Sometimes it can give me the key phrases directly, and sometimes it can find that the raw contents are not related to the career. Including phrases Induced by my self or by online information, I summarize all these phrases (listed in Appendix A) into the `phrases_list` and use them to filter the paragraphs in the rdd.

There are ways I think can be used:

- For each phrase in the `phrases_list`, I will tokenize it to a small phrase_tuple. If one paragraph contains all words of any phrase_tuple (even in disorder or in different morphologys), I will consider it as a paragraph talking about career planning, which will be kept in the final result.
- Let all tokens in the paragraph be the nodes of a graph, and the edges are the co-occurrence of two words in the same paragraph. Then I will use a graph algorithm to find the subgraph that contains all nodes in the phrase. If the subgraph is connected, the paragraph is considered to be talking about career planning.

> Word tenses and case should be ignored because authors may use different tenses and case to convey the same meaning. So I will convert words into lower cases as well as normalize each of them. What I used to normalize words is deleting letters. For example, If `coming` is converted to `come`, then `coming` may not be pointed out, so I use `com` to represent `coming`. But this may relate to some unwanted words like `common`. In my opinion, more is better then less.

Taking the complexity and time efficiency into account, I use the first way to filter the paragraphs. Here is a shortcut of the selected paragraphs.



To save the result into a `parquet` file, I first convert the rdd into a dataframe and then save it.

```
1  spark.createDataFrame(career_suggestion.map(lambda x:
   Row(text=x))).write.format("parquet").mode("overwrite").save("output/career_suggestion.par
   quet")
```

> To save the file successfully, we need set the configure of the spark session at its creation to avoid the error of exceeds of heap memory (shown below). Here my configuration are `("spark.executor.memory", "8g")` and `("spark.driver.memory", "8g")`
>
> 

## (3)

In order to extract noun phrases from the paragraphs, referring to the `Spacy` given by the question, I define a `UDF` function called `extract_noun_phrases`, which can extract noun phrases by model `en_core_web_sm`.

```
1  nlp = spacy.load("en_core_web_sm")
2
3  def extract_noun_phrases(text):
4      doc = nlp(text)
5      noun_phrases = [chunk.text for chunk in doc.noun_chunks]
6      return noun_phrases
7
8  spark.udf.register("extract_noun_phrases", extract_noun_phrases, ArrayType(StringType()))
```

After registering the `UDF`, I can use it in `flatMap` and get all noun phrases in one rdd. Then `CountByValue` and `sorted` by the count can be used to get top 40~50 noun phrases.

```
{'the company': 276,
 'Lisp': 266,
 'everyone': 264,
 'software': 256,
 'the way': 256,
 'themselves': 253,
 'him': 250,
 'a company': 240,
 'those': 236,
 'Google': 232,
 'anyone': 217}
```

Then the word cloud map can be plotted by `wordcloud`. As the parameter `max_words` exceeds the total number of words, the `repeat` is set to `True`.



However, the plot generated by `wordcloud` is not that satifying. So I use `stylecloud` to plot the word cloud map, which is an advanced encapsulation API of `wordcloud`. The result is shown below.

In fact, the api has some bugs, because it has not been maintained for many years, so I jump directly to his api source code and change it.

- if encountering an error that `'ImageDraw' object has no attribute 'textsize'`, that is because `ImageDraw.textsize()` called by it has been updated, and two places are needed to be changed in the source code, which is shown in the figure below (ref: **python - 'ImageDraw' object has no attribute 'textbbox' - Stack Overflow**)



- If there are not enough words in the list, you want to repeat them when plotting. Just modify the `repeat` parameter in the source code, as it call the `wordcloud` function in the `wordcloud` package.

# Appendix

## A

All key phrases I used to match the paragraphy talking about career planning are as follows:

```
1   phrases_list = [
2       "Career plan",
3       "suggestion for",
4       "suggestion of",
5       "suggestion about",
6       "Career development",
7       "Skill development",
8       "Career goal",
9       "Job market",
10      "Career mentor",
11      "Career counsel",
12      "Career transition",
13      "Self-assessment",
14      "Career prospect",
15      "Industry research",
16      "Resume and cover letter",
17      "Interview preparation",
18      "Professional networking",
19      "Continuing education",
20      "Job satisfaction",
21      "Work life balance",
22      "Career success",
23      "Professional ethic",
24      "Career growth"
25      "Find startup",
26      "Generat startup",
27      "Com up zip_with",
28      "Value of initial",
29      "Start point",
30      "Partial solution",
31      "Max future option",
32      "Work on new technolog",
33      "Conversation with friend",
34      "University research startup",
35      "Importance of collaboration and friendship",
36      "Mind wander and idea generation",
37      "Mak thing eas to use",
38      "Mak something people want",
39      "Redefin problem",
40      "Mak thing cheaper and commoditiz",
41      "Mak thing easier",
42      "Design for exit strateg",
43      "Product development on spec",
44      "Counteract monopol",
45      "Accident startup",
46      "Do what hacker enjoy"
47  ]
```