

Zewail city of science and technology

Big Data Analytics CIE 427

Student Name: Ahmed Eid Abdullah

Student Name: Ahmed Kamal

Student ID: 201701264 - 201700295

Doctor Name: ElSayed Hemayed

Teaching assistant Names: Ahmed Sameh – Mohamed ElAref

Course (Code) Name: Big Data Analytics CIE 427

Mini Project 1

Analysis of Reddit website

Zewail city of science and technology

Big Data Analytics CIE 427

Table of content :-

Our approach to solve the problem.....	3
Pipeline diagram	6
Code Environment	7
Challenges	8
Results and graphs	9
Our product is efficient	17
Can our product work on any dataset?	17
Conclusion	20

Zewail city of science and technology

Big Data Analytics CIE 427

Our Approach and how we do preprocess on the data

We downloaded the Hadoop on google colab,

- 1- Firstly, we read the input json file and convert it into to python dictionary { }to access it easily and we do this job by built in library.
- 2- Then we accessed only the subreddit and body keys and values.
- 3- In Mapper.py , we did our preprocessing as the following:
 - Lower each body
 - Removing any non-alphanumeric words from each body, by using simple regex.
 - Removing any links such as <https://youtube.com> , or <https://anything.anytld> (We did this step , as we noticed that the body contains some links)
 - Tokenizing each body (to get its topics)
 - Stemming each tokenized word , then we consider each word as topic, we checked if this word is not stop words or not , not numeric , and we put a condition on the word's length (as we noticed some tokenized word may be very long, so we checked if the length is less than 15, words of length less than 15 usually be understandable) , and we take only the noun words (because the topics usually be nouns not verbs for example)

Zewail city of science and technology

Big Data Analytics CIE 427

```

21 for line in sys.stdin:
22     #read as dict
23     json_obj = json.loads(line)
24     # take only subreddit
25     subreddit_feature = json_obj["subreddit"]
26     # take only body
27     body_feature = json_obj["body"]
28     #preprocessing starts
29     # lower the body
30     lowered_body = body_feature.lower()
31     non_alphanumeric_body = re.sub(r'^\w+', '', lowered_body) # substitute in thing that is not
32     no_links_body = re.sub(r'https?:\/\/.*[\r\n]*', '', non_alphanumeric_body) # i noticed that
33     # we will tokenize body to extract topics
34     tokenization_of_body = nltk.tokenize.word_tokenize(no_links_body)

```

- Print each subreddit and its (stemmed_word which the topic , 1) separated by tab , so we count in the reducer.

- Print each stemmed_word and its upvotes value (from the json ups key)

```

35 for word in tokenization_of_body:
36     ##check if the word is not stop word or numeric
37     if word not in English_stopwords and len(word) <15 and not word.isnumeric() and word in nouns: # probably not a natural word
38         stemmed_word = stemmer.stem(word)
39         ## print the subreddit and its topics
40         print(f"{subreddit_feature}:topics", (stemmed_word, 1), sep='\t')
41         ## print the topics and its upvotes
42         print(f"{stemmed_word}:upvotes", (json_obj["ups"],), sep='\t')

```

4- In reducer, we get the most common 10 topics related with each subreddit and count the upvotes of each topic. We take only the upvotes that more than 500 , because we found that many topics count 70 , 20 , 120 , and many more in this range, so to avoid this rare counted topics, we put condition on upvotes more than 500.

Zewail city of science and technology

Big Data Analytics CIE 427

```

1 #!/usr/bin/python3
2
3 from collections import Counter
4 import sys
5 topic = None
6 word_upvotes = None
7 # here count will be initialized to zeros we will use Counter object
8 count = Counter()
9 ##initialize the counter to count the upvotes
10 count_upvotes = 0
11 word_counting = 0
12 for line in sys.stdin:
13     key, value = line.strip().split("\t")
14     # because strip function will convert it to string
15     value = eval(value)
16     *key, category = key.split(":") # as the dict contains key and value (category)
17     # first question , topics of every subreddit
18     if category == "topics":
19         #check first if the topic is none
20         if topic is None:
21             topic=key
22         # if not , then count the most common discussed topics
23         if topic != key:
24             print(f"topic:{category}", count.most_common(10), sep='\t') # top three discussed topics
25             topic = key
26             count = Counter()
27             count[value[0]] += value[1]
28     # third question , highest upvotes
29     if category == "upvotes":
30         if word_upvotes is None:
31             word_upvotes = key
32         if word_upvotes != key:
33             if word_counting > 500: ## we take 120 , as an average number , because we notice that many numbers and we want the h
34                 print(f"word_upvotes:{category}", count_upvotes, sep='\t')
35                 word_upvotes = key
36                 count_upvotes = 0
37                 word_counting = 0
38             count_upvotes += value[0]
39             word_counting += 1
40     # print the top three discussed topics
41     if category == "topics" and topic:
42         print(f"topic:{category}", count.most_common(10), sep='\t')
43     # print upvotes of the topics
44     if category == "upvotes" and word_upvotes:
45         if word_counting > 500:
46             print(f"word_upvotes:{category}", count_upvotes, sep='\t')
47
48

```

5- We open outputs files and convert them into CSV file, to be visualized later.

```

▶ f = open("part-00000", "r")
  lines = f.readlines()

[ ]

topics = []
upvotes=[]

for line in lines:
    key, value = line.strip().split("\t")
    value = eval(value)

    *key, case = key.split(":")
    key = ":".join(key) # handle the case where key has a :
    if case == "upvotes":
        if value > 500:
            upvotes.append((key, value))
    elif case == "topics":
        topics.append((key, value))

```

Zewail city of science and technology

Big Data Analytics CIE 427

```
[ ] import csv
    header = ['SUBREDDIT NAMES', 'TOP THREE DISCUSSED TOPICS']

    with open('TOP THREE DISCUSSED TOPICS FOR EVERY SUBREDDIT.csv', 'w', encoding='UTF8', newline='') as f:
        writer = csv.writer(f, delimiter=',')
        writer.writerow(header)
        # write multiple rows
        writer.writerows(topics)
```

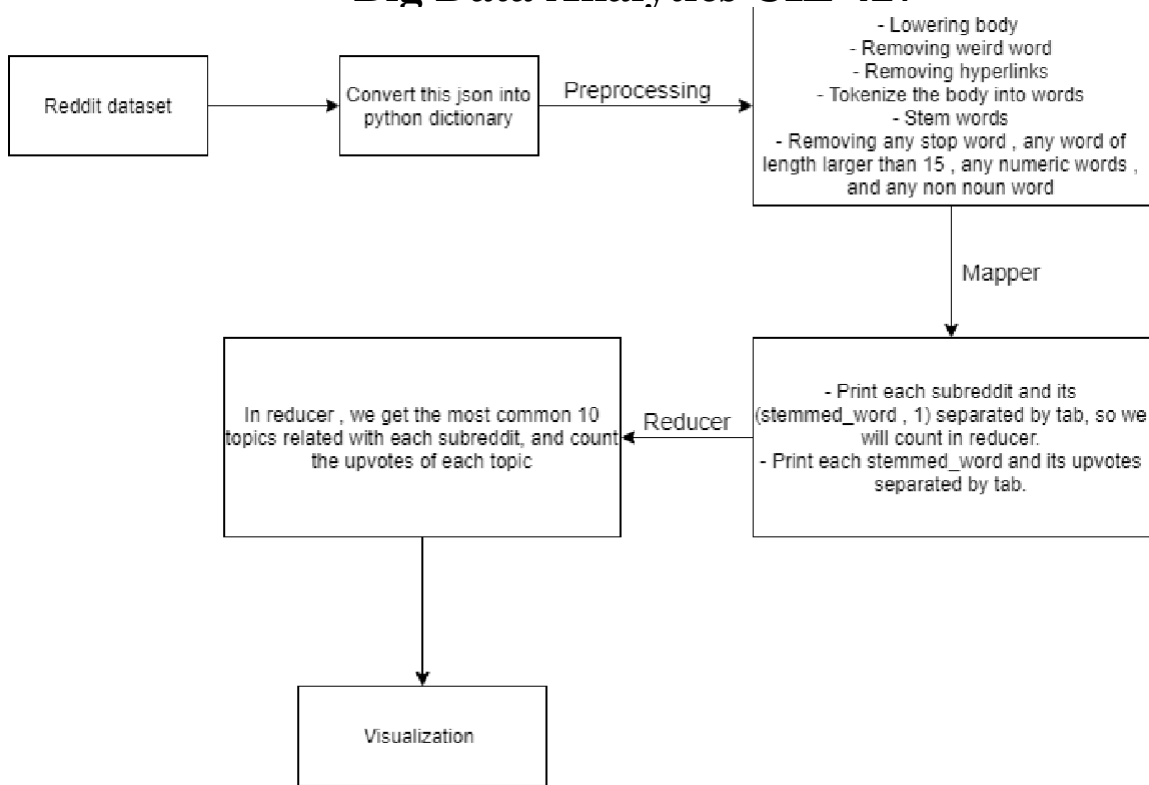
```
▶ header = ['SUBREDDIT NAMES', 'TOP THREE DISCUSSED TOPICS']

    with open('TOP UPVOTES.csv', 'w', encoding='UTF8', newline='') as f:
        writer = csv.writer(f, delimiter=',')
        writer.writerow(header)
        # write multiple rows
        writer.writerows(upvotes)
```

6- We used two websites to make most of charts which are <https://app.datawrapper.de/> and <https://online.visual-paradigm.com/>

Pipeline diagram4

Zewail city of science and technology Big Data Analytics CIE 427



Code Environment

- Programming language: Python
- Where we run the code: Google Colab <https://colab.research.google.com/>
- Submitted files:
 - 1.(mapper.py)
 - 2.(reducer.py)
 - 3.(Hadoop_Notebook_MapReduce.ipynb)
 - 4.(Convert to csv notebook.ipynb)

Zewail city of science and technology

Big Data Analytics CIE 427

- How can you run the code?

1. Open “Hadoop_Notebook_MapReduce” on google colab and follow instructions in it to download and unzip both hadoop and the dataset
2. Upload mapper and reducer file to google colab then run the codes sequentially

To get the output (part-00000) file

3. Open “Convert to csv notebook” on google colab and upload the output file (part-00000) this notebook is created to convert the output to csv file to give options of Excel function that can be applied on csv file and also tons of online websites that enabled us to make elegant charts and graphs .

Challenges we faced and how we solve them?

We faced the following challenges :

1. We cannot unzip the dataset correctly as it has special format called bz2 , none of normal compressing software such as winRare was capable of unzipping the dataset , so we figured out a way by using a library in python called bz2 is specialized to unzip this kind of compression , but first we mount colab on our drive to be able to read the dataset directly without uploading it
2. We first run the dataset locally on our machine but that was very very slow

Zewail city of science and technology

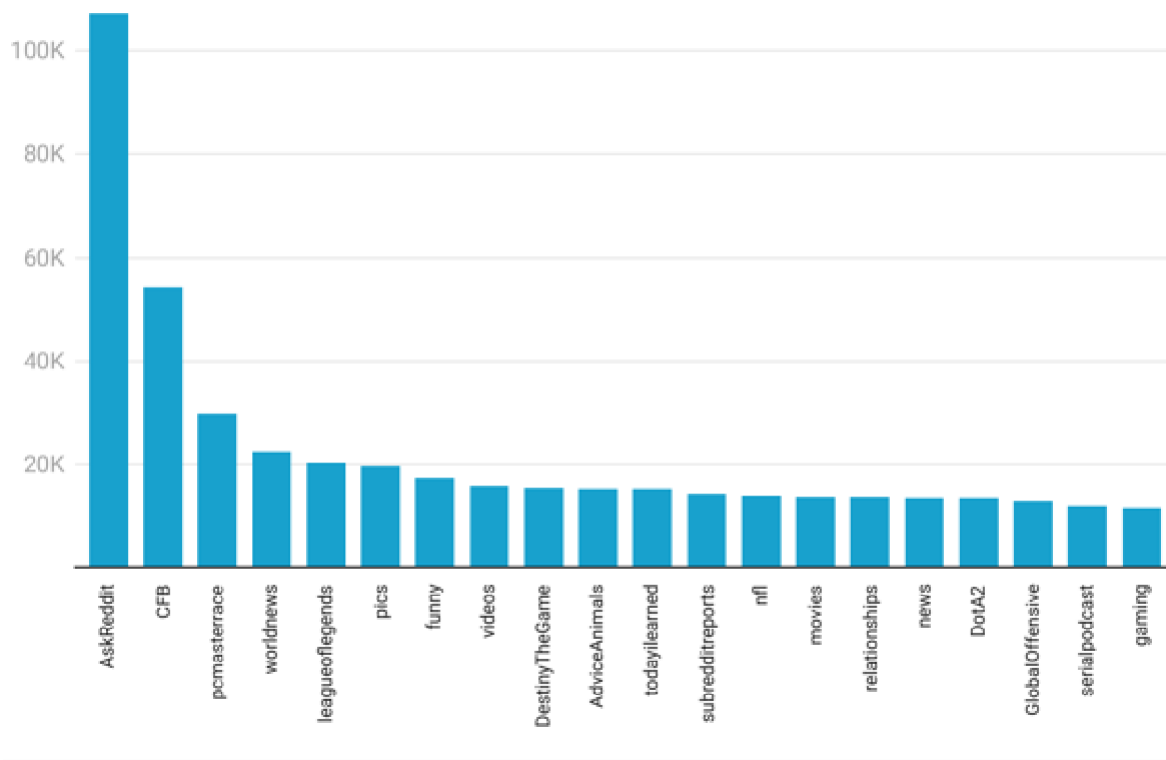
Big Data Analytics CIE 427

to processing on the whole dataset , but Aiad Assad pointed that we can use hadoop on Google colab and he tested it and it was like super power compared to our testing time , then we shifted to google colab to use hadoop online and all credit for this point goes to Aiad , this point enabled us to run on the whole dataset with suitable amount of time ,that make our life easier and was superior guide from Aiad

Results and graphs

We found that the most common subreddit is AskReddit. For simplicity we visualize the top 20 subreddits only to be well visualized.

Top 20 subreddits

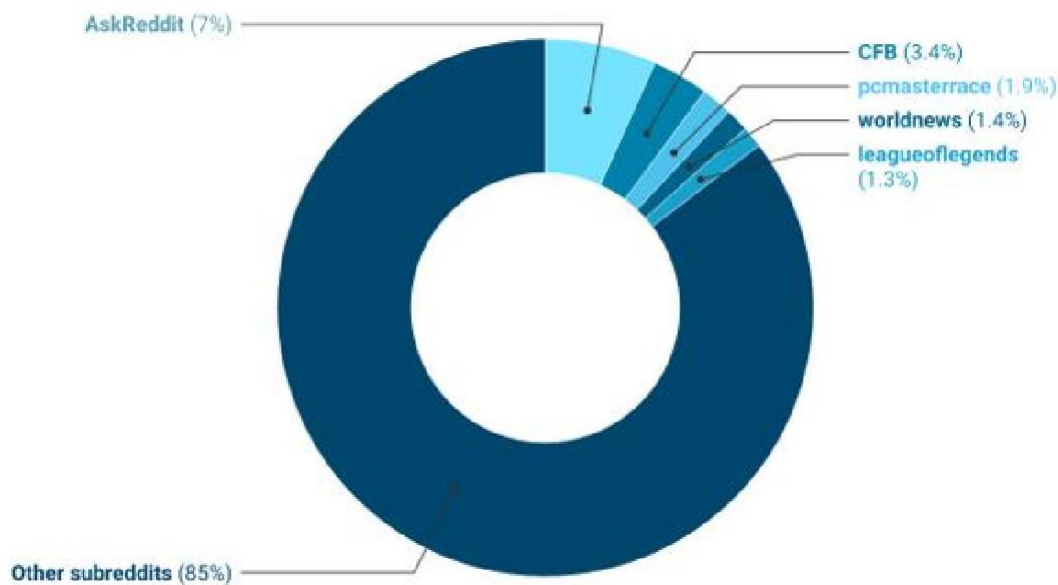


The distribution of top 1000 subreddits

Zewail city of science and technology Big Data Analytics CIE 427

The percentage of AskReddit subreddit is 7% (very high percentage), CFB is 3.4%, pcmasterrace is 1.9% , worldnews is 1.4% , leagueoflegends is 1.3% , while the other 995 topics is 85%

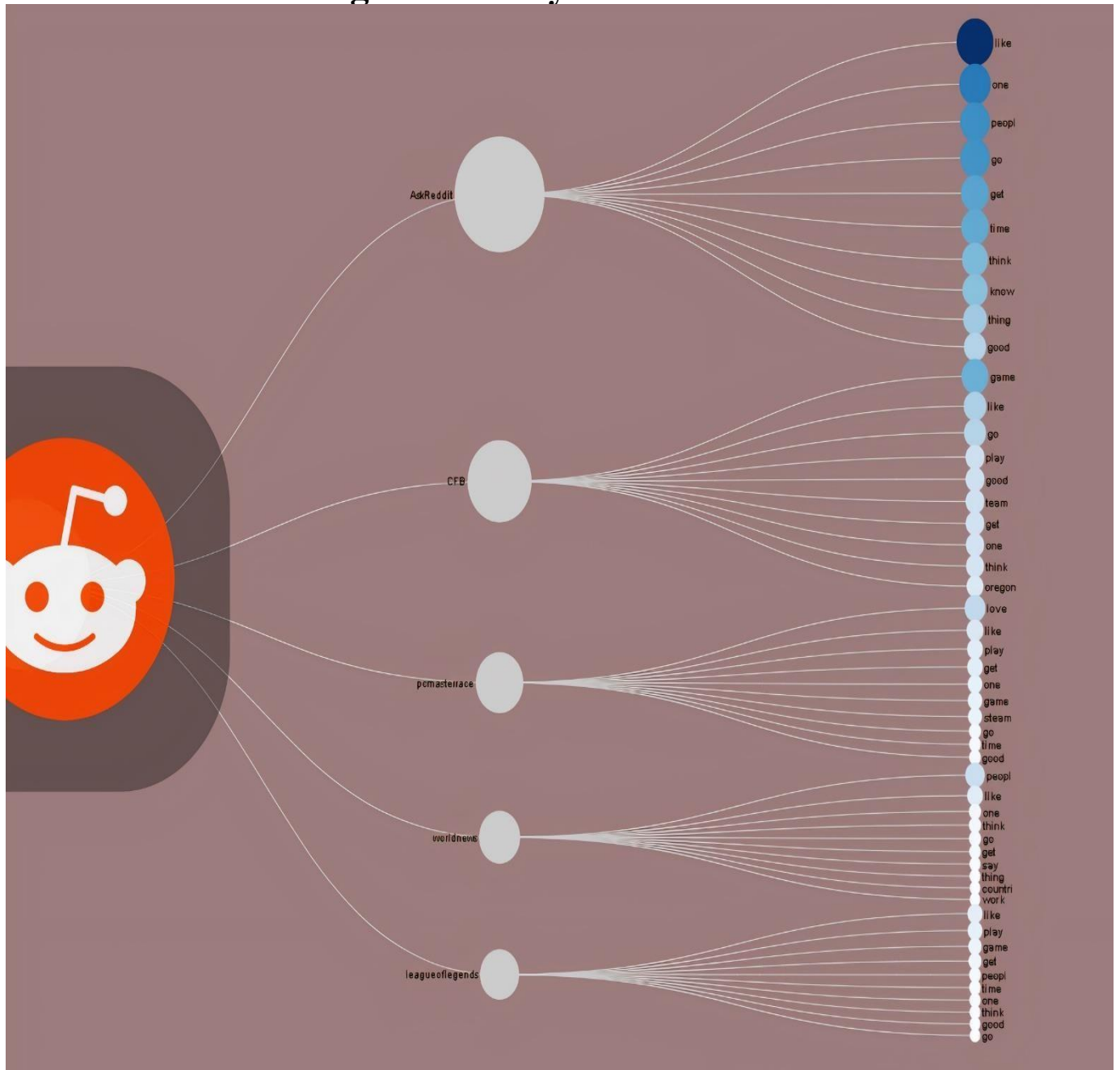
Distribution of top 1000 subreddits



Visualization of the Most 10 discussed topics associated with the top 5 subreddit in Reddit

linear Dendrogram Graph

Zewail city of science and technology Big Data Analytics CIE 427



After deletion of commas and square brackets , we use Excel functions to sort subreddits

Based on their summation on top ten topics counter , then some preprocessing of the data

To be suitable to be inserted into this [site](#) which is can create linear dendrogram with our data.

Zewail city of science and technology

Big Data Analytics CIE 427

BIG NOTE

1. Dark blue means highest frequency topic
2. White means lowest frequency topic
3. Size of the word circle means (bigger circle size higher frequency topic)
4. All topics for every subreddit is sorted descending based on their frequency in the subreddit

✓ **BIG BIG NOTE : Creative Tip**

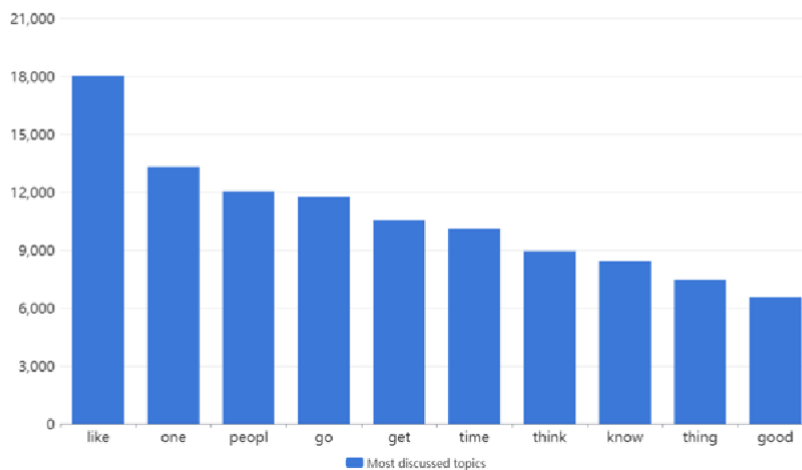
This graph is expandable which means we can use this fancy graph for all subreddits in the reddit website and for each subreddit we demonstrate every topic in that subreddit with size of circle indicating frequency of that topic (this can be a very powerful and efficient tool to know the trending topic from where it created and when and how big this topic compared to other topics) this is very similar to HASHTAG that can tackle the number of people talk about trending topic and who created this HASHTAG.

- We choose to make linear dendrogram graph on only top five subreddits with their associated top ten topics due to processing power .

Zewail city of science and technology Big Data Analytics CIE 427

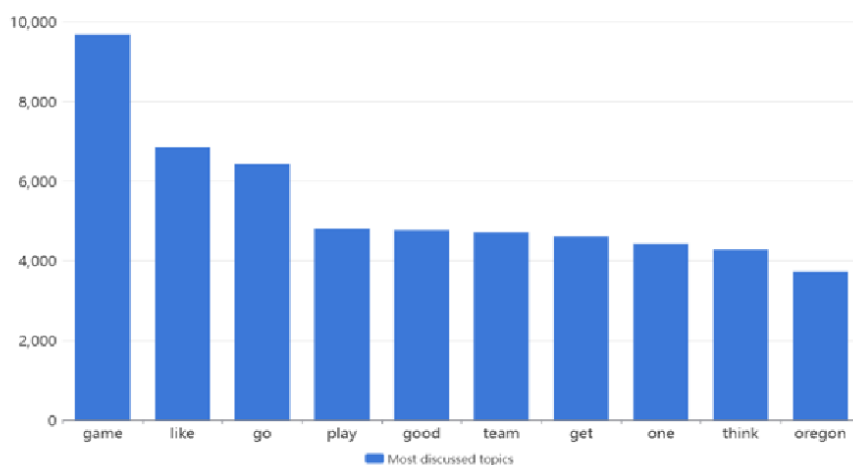
Top 10 most discussed topics in top 5 subreddits

Most 10 discussed topics in AskReddit



Top 10 most discussed topics in CFB subreddit

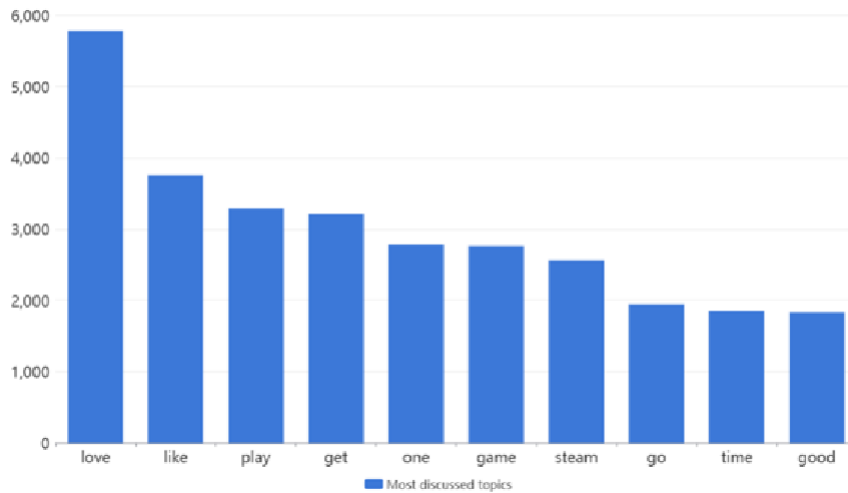
Most 10 discussed topics in CFB



Zewail city of science and technology Big Data Analytics CIE 427

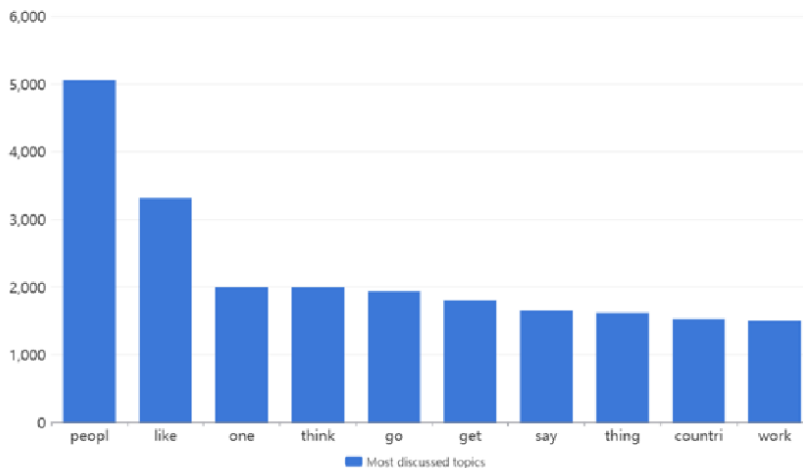
Top 10 most discussed topics in pcmasterrace subreddit

Most 10 discussed topics in pcmasterrace



Top 10 most discussed topics in worldnews subreddit

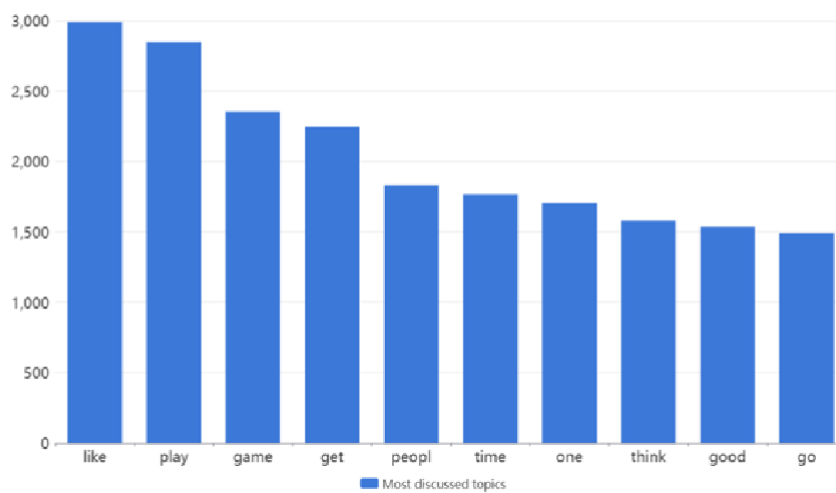
Most 10 discussed topics in worldnews



Zewail city of science and technology Big Data Analytics CIE 427

Top 10 most discussed topics in leagueoflegends subreddit

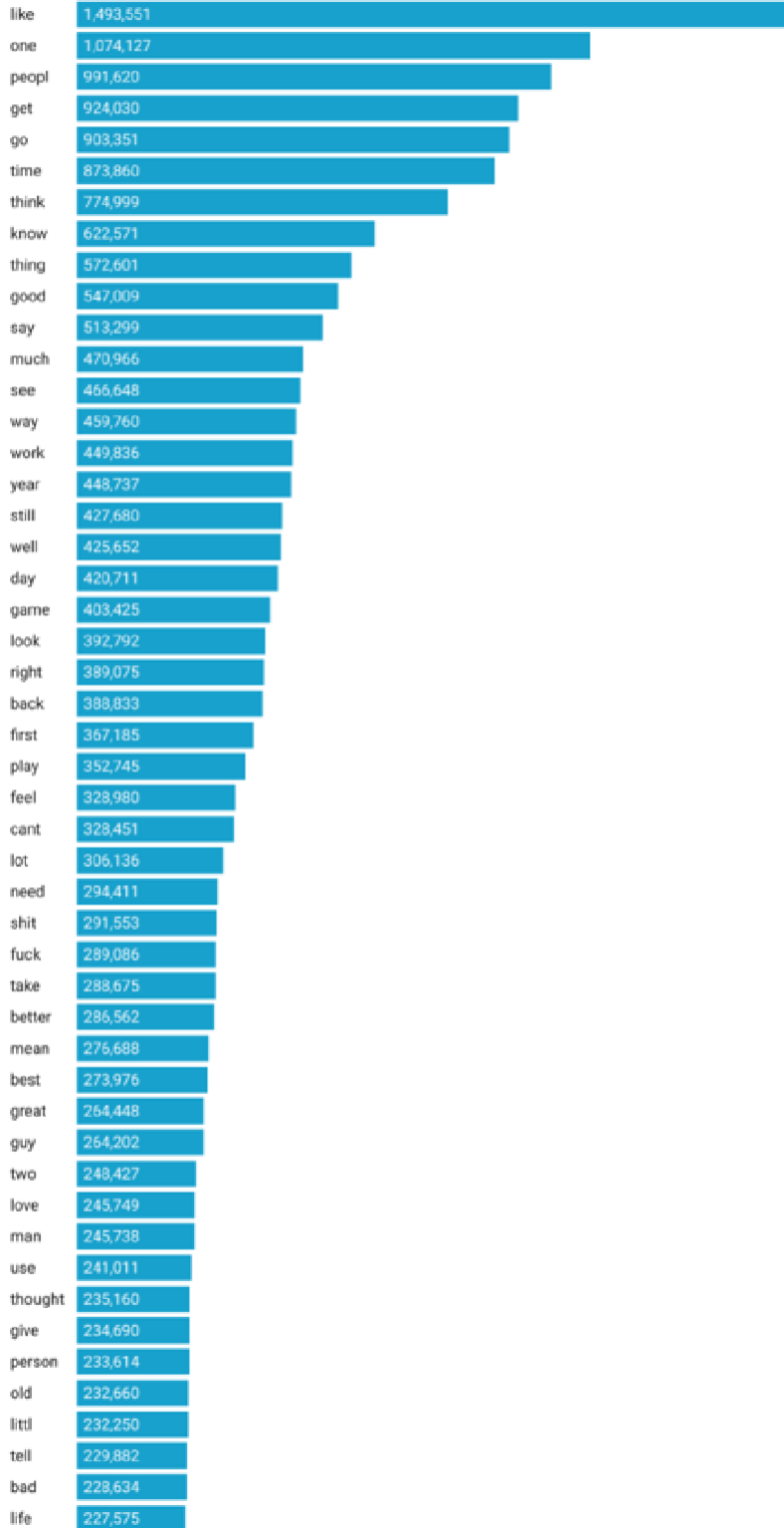
Most 10 discussed topics in leagueoflegends



Top 50 topics that yield highest upvotes

Zewail city of science and technology Big Data Analytics CIE 427

Top 50 topics that yield highest upvotes

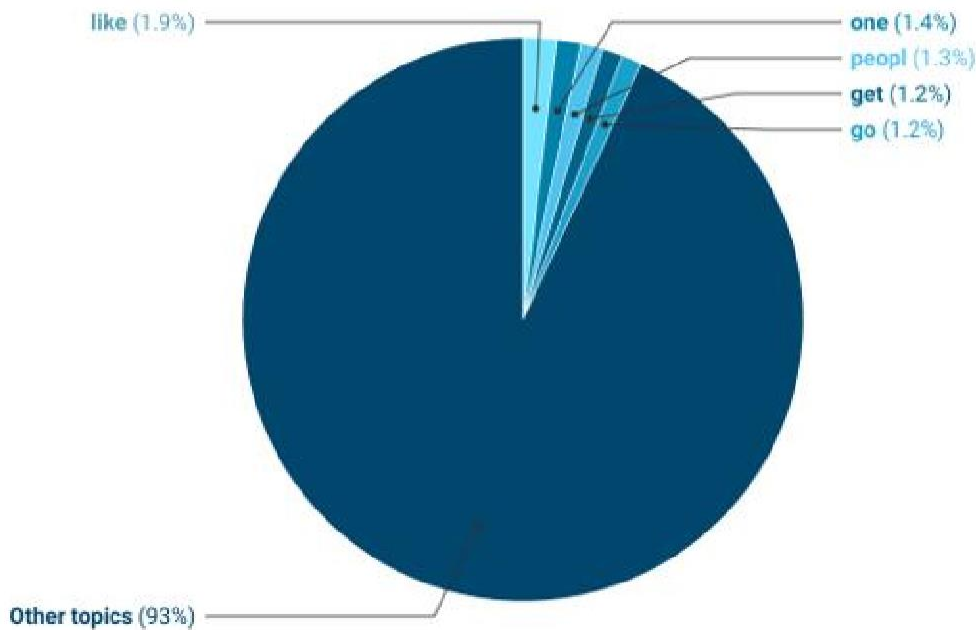


Zewail city of science and technology

Big Data Analytics CIE 427

Distribution of top 5804 topics in all subreddits

The distribution of 5804 topics



Why our product is very efficient?

1. Our product can run on any dataset of any size, because we use Hadoop system (that deals with any big data)
2. Produce good and detailed visualizations: It can generate multi charts for further analysis of Reddit (in our case). For simplicity, we visualize only the top 20 subreddits and how are they distributed, most 10 topics of the top 5 subreddits, the top 50 topics that yield highest number of upvotes (people like it), and how the top 5804 topics are distributed.

Zewail city of science and technology

Big Data Analytics CIE 427

3. Easy to run: We used google colab to run our code, which is very easy to run as all you need is to run each cell consecutively.
4. We used good data preprocessing by using different NLP libraries to remove stop words, remove not understandable characters, remove words that more than 15 characters, remove hyperlinks, we work only on nouns (as topics cannot be verbs for example), we exclude the numeric words, and we did many more preprocessing.
5. We convert all outputs into CSV files and then convert them into charts to be readable and understandable.

Creative/Innovative

- Take only Nouns from body (we download word net) to check in word is considered as noun or not wordNet use synset (presented in NLTK [here](#)) and word embedding to know that.
- Recommended creative tip that we thought of but we cannot applied it -> we observe that wordNet is consider some verb as topic also for example like wordNet consider it as noun because it can means admiration, so we think of 2 or 3 N-gram to consider the context of the sentence to be 100 % accurate in determination if this word is noun verb based on the context , this approach will

Zewail city of science and technology

Big Data Analytics CIE 427

capture only topic that 100% people talk about such as HASHTAG technology , but note this process is need powerful processing and even more with the increasing of N-gram.

- Put a threshold on up votes (>500) to be included , this threshold we see it suitable to consider topic is highly up voted or not , because if we make it small random number , the algorithm will consider any topic is highly up voted (consider case if topic is get 100 up vote but from friends of user and not relevant to the decision maker if this topic is trending or not) that is why we have high up vote threshold to capture trending topic case or even policy violation detection (higher down votes is will also be considered from the decision maker or from the data analyst)
- Before calculating Controversiality we notice that in first 10 sample its value is zero , then we checked it for the whole dataset we found it always zero , so it impossible to calculate rate of replies because Controversiality is always zero
- Linear dendrogram is very powerful tool that we created to visualize everything for each subreddit

Dendogram can describe as it is tree has all ancestors (subreddit) with its children (topics) with their frequency (circle size) this can easily expanded to include all reddit website information if processing power is available

Conclusion

Zewail city of science and technology

Big Data Analytics CIE 427

We did analysis of Reddit website, we worked on Hadoop, as our dataset is very large 30GB, we convert all outputs in CSV and then visualize all outputs. In the future we can implement GUI and connect this code with it to be a real product.