

Winter 2024

Hao Jun Chen

February 7, 2024

1 Part 1 Understand Zipf's Law

(a) implementation of text normalization module

```
1      def is_stop_word(token):
2          stop_words = set([
3              'a', 'an', 'and', 'are', 'as', 'at', 'be', 'but', 'by',
4              'for', 'if', 'in',
5              'into', 'is', 'it', 'no', 'not', 'of', 'on', 'or', 'such', 'that', 'the',
6              'their', 'then', 'there', 'these', 'they', 'this', 'to',
7              'was', 'will', 'with'
8          ])
9
10         # Convert the word to lowercase for case-insensitive comparison
11         lowercase_word = token.lower()
12
13         return lowercase_word in stop_words
14
15     def remove_non_alphanumeric(input_string):
16         # Use regular expression to keep only English letters and numbers
17         result = re.sub(r'[^a-zA-Z0-9]', '', input_string)
18         return result
19
20     def is_integer(s):
21         try:
22             int(s)
23             return True
24         except ValueError:
25             return False
26
27     path = os.getcwd() + "/yelp"
28     tokens = {}
29     TFF_count = {}
30     DF_count = {}
31     stemmer = PorterStemmer()
32     yelp_store_list = os.listdir(path)
33     #Loop through each yelp file
34     for yelp_store in yelp_store_list:
35         store_path = path + "/" + yelp_store
36         with open(store_path, 'r') as file:
37             data = json.load(file)
38             reviews = data["Reviews"]
39             #Loop through each review in the yelp _ file
40             for review in reviews:
41                 doc = review['Content']
42                 doc_ID = review['ReviewID']
```

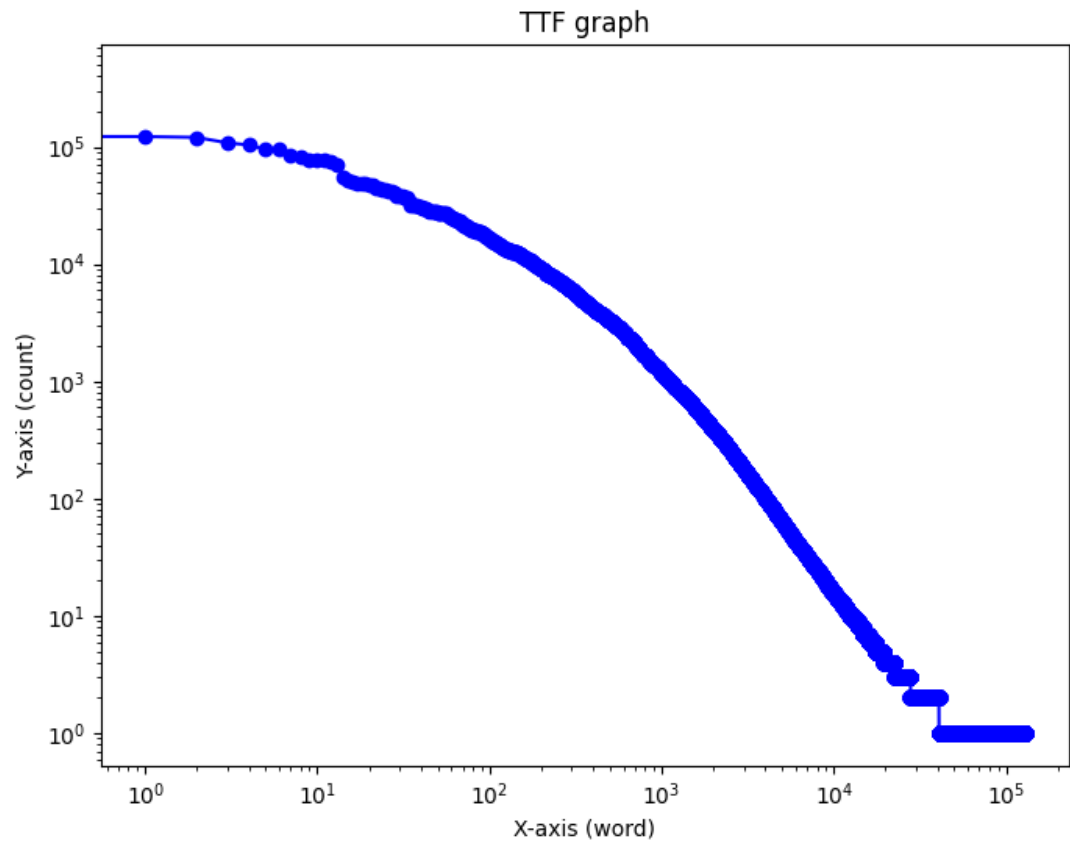
```

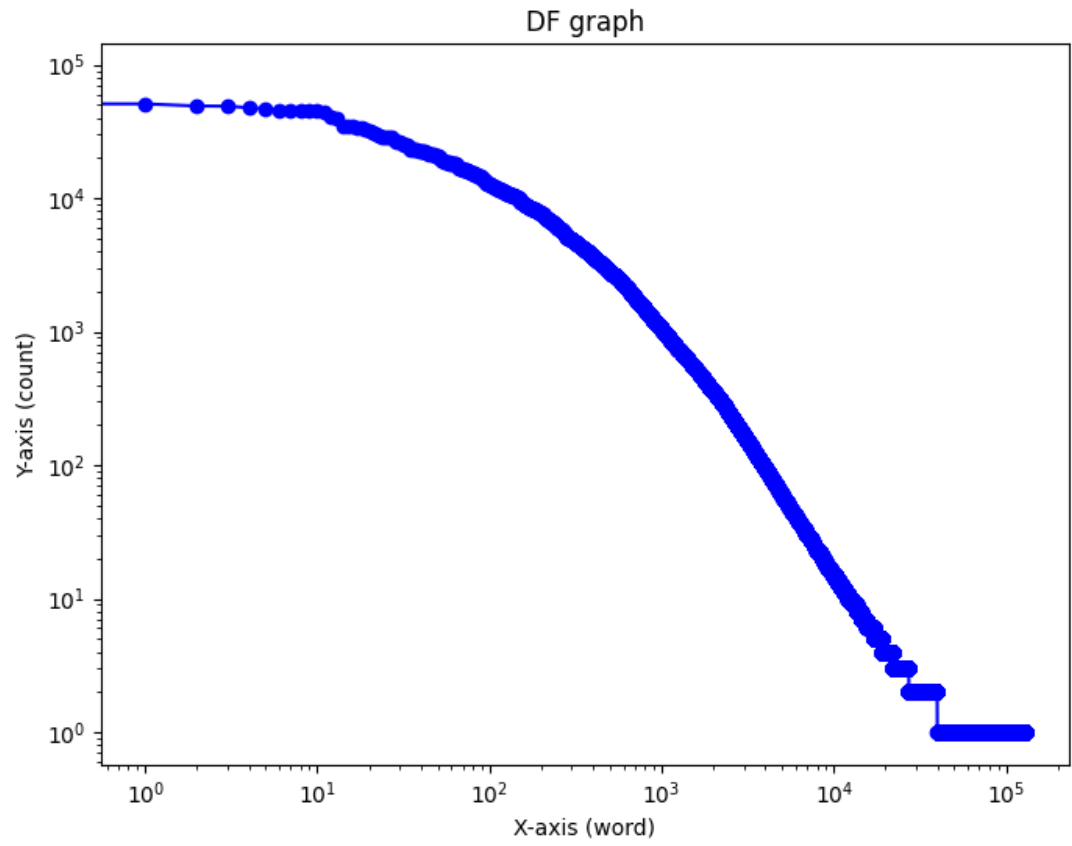
38     tokens = word_tokenize(doc) #token the review content
39     #Loop through each token and normal it
40     for token in tokens:
41         token = remove_non_alphanumeric(token)
42         if len(token) == 0:
43             continue
44         elif is_integer(token):
45             token = "NUM"
46         elif is_stop_word(token):
47             continue
48         else:
49             token = stemmer.stem(token)
50             token = token.lower()
51         TFF_count[token] = 1 + TFF_count.get(token, 0)
52         if token in DF_count and DF_count[token] is not
53             None:
54             DF_count[token].add(doc_ID)
55         else:
56             DF_count[token] = set()
57             DF_count[token].add(doc_ID)

```

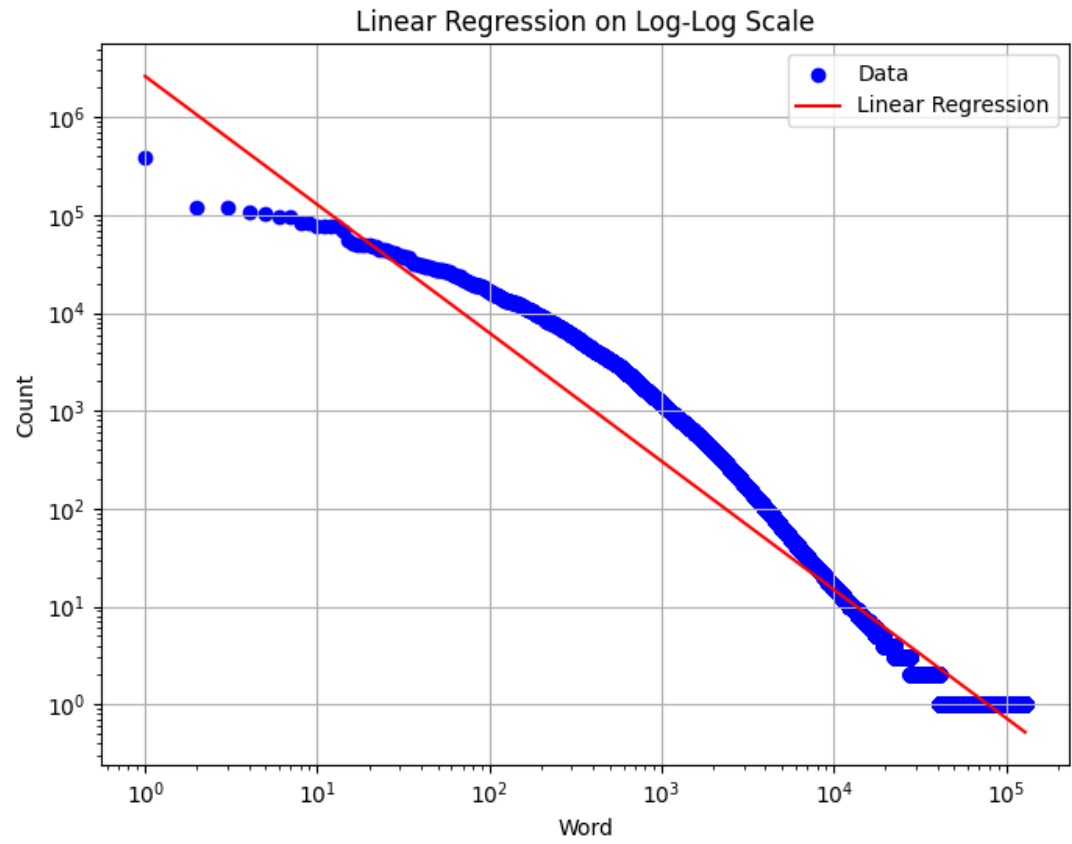
Listing 1: Implementation of text normalization

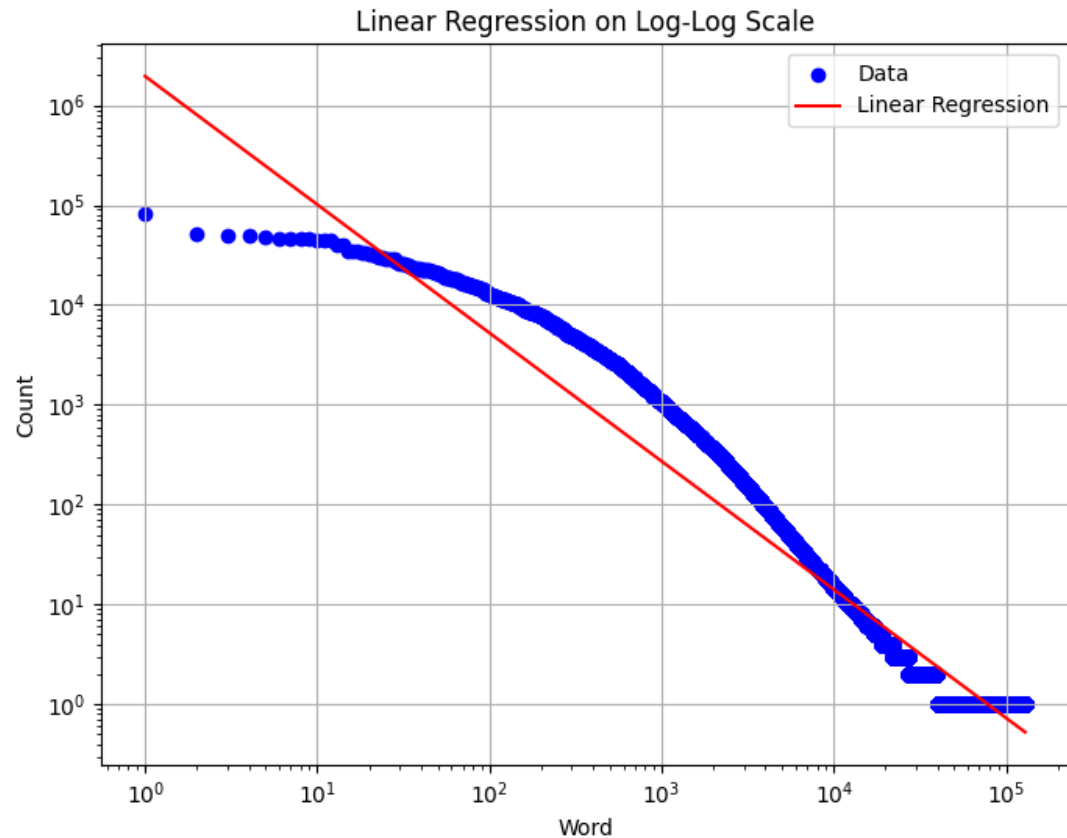
- (b) Two curves in log-log space generated above, with the corresponding slopes and intercepts of the linear interpretation results





(c) Linear Regression on TFF and DF graph





I found a fairly strong linear relationship in these two graphs. As both of their slope are around -1.3, which is very close to one.

2 Part 2 Build Inverted Index

(a) implementation of text normalization module

```

1  def is_stop_word(token):
2      stop_words = set([
3          'a', 'an', 'and', 'are', 'as', 'at', 'be', 'but', 'by'
4          , 'for', 'if', 'in',
5          'into', 'is', 'it', 'no', 'not', 'of', 'on', 'or', '
6          such', 'that', 'the',
7          'their', 'then', 'there', 'these', 'they', 'this', 'to
8          ', 'was', 'will', 'with'
9      ])
10
11     # Convert the word to lowercase for case-insensitive
12     comparison
13     lowercase_word = token.lower()
14
15     return lowercase_word in stop_words
16
17 def remove_non_alphanumeric(input_string):
18     # Use regular expression to keep only English letters and
19     numbers

```

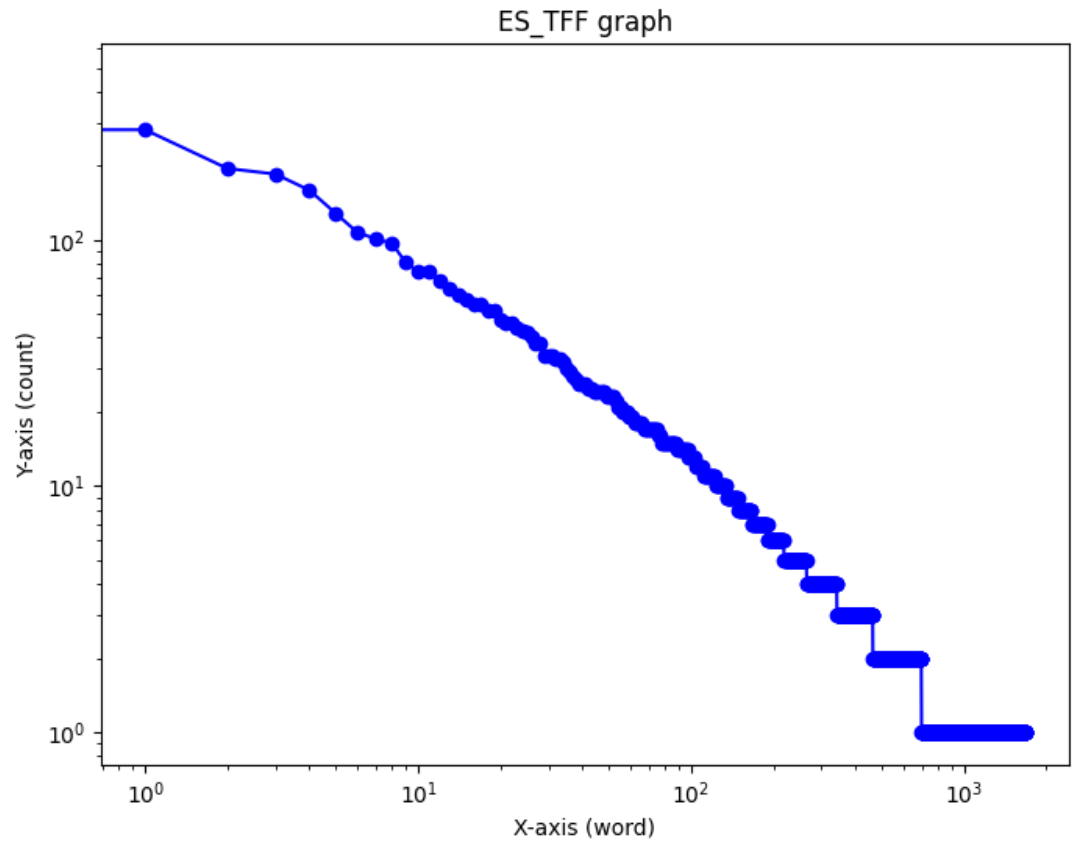
```

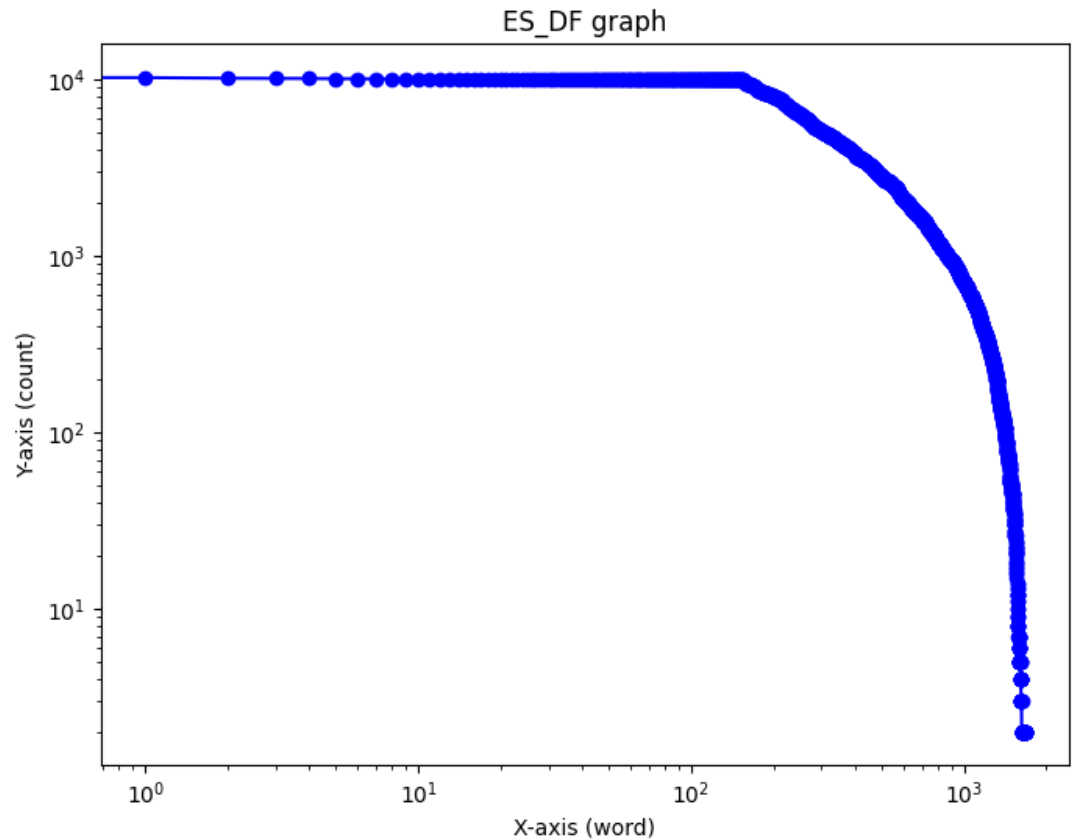
14         result = re.sub(r'[^a-zA-Z0-9]', '', input_string)
15     return result
16 def is_integer(s):
17     try:
18         int(s)
19         return True
20     except ValueError:
21         return False
22 path = os.getcwd() + "/yelp"
23 inverted_index = {}
24 stemmer = PorterStemmer()
25 yelp_store_list = os.listdir(path)
26 #loop through all the yelp data
27 for yelp_store in yelp_store_list:
28     store_path = path + "/" + yelp_store
29     with open(store_path, 'r') as file:
30         data = json.load(file)
31     reviews = data["Reviews"]
32     #loop through each review item in each data
33     for review in reviews:
34         doc = review['Content']
35         doc_ID = review['ReviewID']
36         tokens = word_tokenize(doc)
37         #Normalize each word
38         for token in tokens:
39             token = remove_non_alphanumeric(token)
40             if len(token) == 0:
41                 continue
42             elif is_integer(token):
43                 token = "NUM"
44             elif is_stop_word(token):
45                 continue
46             else:
47                 token = stemmer.stem(token)
48                 token = token.lower()
49             #For each word, record the document that contain
             it
50             if token in inverted_index and inverted_index[
                token] is not None:
51                 inverted_index[token].add(doc_ID)
52             else:
53                 inverted_index[token] = set()
54                 inverted_index[token].add(doc_ID)

```

Listing 2: Implementation of inverted index

- (b) Two curves in log-log space generated by using Elasticsearch inverted index for collecting TTF and DF statistics





Comparing to the graph from part one, We can see that both TFF and DF have similar graph shapes but have smaller amounts of frequency on both TFF and DF. This makes sense since we are looking at only 9 queries.

(c) Running time and total number of returned documents by two retrieval models

Runtime for Elastic search is 1.33 seconds and my model is 0.003 seconds. As Dr. Wang mentions in class, this is due to the system Elasticsearch is bigger. If we keep putting large data in, eventually, Elasticsearch will be faster. I got quite different results from the Elasticsearch. One major reason is how Elasticsearch handles tokenization and normalization as I put unprocessed documents into elastic search. I believe that if I put processed document, the retrieved document would be similar between my model and elastic one. Here is the table of my model vs elastic search,

	A	B	C	D
1		ES model	My model	
2	general chicken	>=10000	507	
3	fry chicken	>=10000	4767	
4	bbq sandwich	9513	442	
5	mash potato	3926	1617	
6	grill shrimp	>=10000	499	
7	lamb shank	1502	54	
8	pepperoni pizza	>=10000	669	
9	friend staff	>=10000	1572	
10	grill cheese	>=10000	1462	
11				