

Project: 1 How do philosopher talk?

Introduction

In this project I will be investigating how past philosopher talk. Specifically, perform a words frequency analysis to see if I can have some interesting findings.

Philosophers are thinkers of our worlds. They often bring new ideals, instigate reform, or propose hypothesis. An analyzsis of how those people talk can help us understand their ideas more clearly and bring fresh insight of how we look at the world.

Information about the dataset

The dataset I used for this project adopted from <https://www.kaggle.com/kouroshalizadeh/history-of-philosophy>. This dataset contains 360808 sentence said by famous philosophers. The time in year of when the sentence is said, and the school of which the author belongs to.

Flaw in the dataset

The data itself is mostly clean. For one author or school, it appears the same through out the dataset making it easy for later to search or group. The sentence is also very clean. There is typo or words incorrectly joined together. The creator of this database also provide a lowered case sentence as well as tokenized and lemmatized text. Unfortunately those text cannot be used for our analysis.

Another problem with this dataset is with the original publication data. Mainly the publication date of the author Plato and his respective school plato is negatvie. There are also some publication data does not correspond to an actual year. By elimination them, the dataset would loss almost 25% content. Hence, I decide to keep those data since year is not included in my analysis.

Process the data

In order to do word frequency analysis. I have to clean each sentence to remove any pronouns, preposition, punctuation ect. Then, split each word into tokens. I also lemmatized the words, this step can reduce each words to its simplist form.

For example:

" What's new, Socrates, to make you leave your usual haunts in the Lyceum and spend your time here by the king archon's court?"

Would become:

'new','socrates','make','leave','usual','haunt','lyceum','spend','time','king','archon','court'

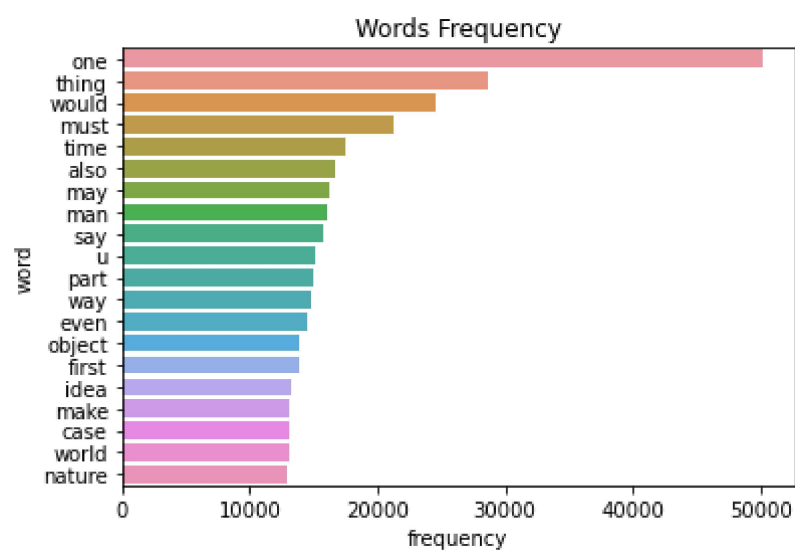
Which is the Most Famous Word Used by Philosophers

After processing the data. I can answer the first question: Which is the Most Famous Word Used by Philosophers? I present the top three words used by philosopher and their respect frequency.

Out[15]:

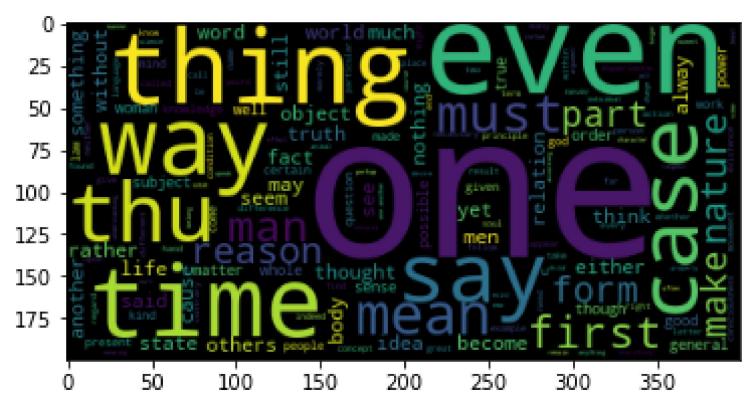
	word	frequency
48	one	50186
51	thing	28692
183	would	24614

I also make a bar chart to show the top 20 words used by philosophers.Using this result, I can demonstrate the result visually by plotting a word cloud as shown below. This would gives us some idea on the word choice of philosopher and give us insight on interpreting how philosopher talk.



Out[145]:

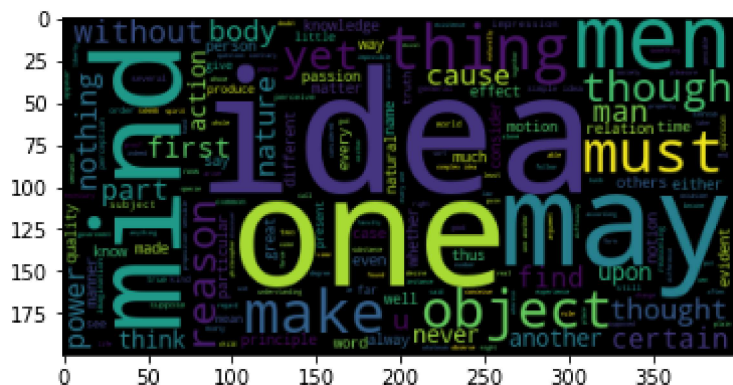
<matplotlib.image.AxesImage at 0x24bca3a8070>



obtained from the general trend. However, just by looking at the graphs and number is not accurate. In the next section, I will be performing hypothesis testing to see if each school talks differently.

Out [39]:

	word	frequency
3924	idea	5871
33	one	3375
58	may	2696



Does Each School Talk similarly?

In this section, I will use ANOVA test to see if the word choice of each school is similar.

The words that are tested is the top ten words used from all sentences (one,thing,would,must,time,also,many,say,u). The first step is to get the sum of each word in each school as shown below.

Out[42]:

word	one	thing	would	must	time	also	may	man	say	u	school
0	5556	4609	3302	2374	1333	1165	690	2162	2971	2047	plato
1	9698	7365	2027	4481	2654	4344	2444	4100	2440	1168	aristotle
2	3375	2041	1258	1488	789	530	2696	1416	677	1980	empiricism
3	3099	3184	2103	1718	807	1002	790	1319	1137	2595	rationalism
4	6437	2457	4606	2004	2083	2055	3784	852	4255	1894	analytic
5	4214	1152	1967	1529	1781	1130	743	1359	886	1020	continental
6	2863	2242	819	1535	1916	998	520	581	719	1497	phenomenology
7	7302	3072	3149	2752	2359	3332	1483	309	971	1685	german_idealism
8	1814	309	711	696	1225	556	608	334	493	333	communism
9	1610	408	2960	1257	1309	242	1613	381	310	238	capitalism
10	322	687	82	188	145	165	147	341	86	57	stoicism
11	2124	781	605	633	478	648	404	1063	437	408	nietzsche
12	1772	385	1025	582	677	597	369	1863	339	216	feminism

Now, I will perform ANOVA test on the chart we just obtained. This test will tell us does the words choice differ by each school.

Out[59]:

```
F_onewayResult(statistic=7.619955343937744, pvalue=2.8318483451450726e-10)
```

Here the p-value is $2.8318483451450726e-10$, which is less than the significant level. Hence, we can conclude that not all school have similar word choice. This is easy to understand since there are 13 different school each has their own idea. It is reasonable that some school talk differently than others. Then my question become, are there some schools have similar word choice? To investigate this, I performed pairwised ANOVA test between every school. Here is a list of which the test is significant between the two schools.

```
plato aristotle
plato empiricism
plato rationalism
plato analytic
plato continental
plato german_idealism
aristotle analytic
aristotle german_idealism
empiricism rationalism
empiricism continental
empiricism phenomenology
empiricism german_idealism
empiricism capitalism
rationalism analytic
rationalism continental
rationalism phenomenology
rationalism german_idealism
rationalism capitalism
analytic german_idealism
continental phenomenology
continental german_idealism
continental capitalism
phenomenology german_idealism
phenomenology capitalism
phenomenology nietzsche
phenomenology feminism
communism capitalism
communism nietzsche
communism feminism
capitalism nietzsche
capitalism feminism
nietzsche feminism
```

This can give us insight on which schools have similar word choice. Based on this information, further inference can be made.