

In [1]:

```
1 # scraping the pdf files to obtain the url
```

In [36]:

```
1 import requests
2 import pandas as pd
3 import urllib.request
4 import pdftotext
5 import promptlib #pip install prompt
6 import os
7 import glob
8 #conda install poppler : conda install -c conda-forge poppler
9 # conda install pdftotext : conda install -c conda-forge pdftotext
10 # conda install nltk : conda install -c anaconda nltk
11 from nltk import *
12 from nltk.corpus import *
13 import nltk
14 import sys
15 from sklearn.metrics.pairwise import cosine_similarity
16 #import the TfidfVectorizer from Scikit-Learn.
17 from sklearn.feature_extraction.text import TfidfVectorizer
18 # Import cosine_similarity
19 from sklearn.metrics.pairwise import linear_kernel
20 import matplotlib.pyplot as plt
21 import seaborn as sns
22
23 import warnings
24 warnings.filterwarnings("ignore")
```

In [7]:

```
1 #get download.csv by going on https://tel.archives-ouvertes.fr/search/advanced-export/u
2 #then export csv
3 df = pd.read_csv("H:/Downloads/Datatasets/Tel/tel_docs.csv")
```

In [8]:

```
1 df.head()
```

Out[8]:

	halId_s	version_i	uri_s	docType_s	doId_s	nntId_s	
0	tel-03440243	1	https://hal.univ-lorraine.fr/tel-03440243	THESE	NaN	2021LORR0152	L'imp sys d'infor hosp
1	tel-03440181	1	https://tel.archives-ouvertes.fr/tel-03440181	THESE	NaN	2021UPASG065	Algorithme critères prédic
2	tel-03440058	1	https://tel.archives-ouvertes.fr/tel-03440058	THESE	NaN	2021NORMR027	La microfin a microentrepre
3	tel-03439538	1	https://pastel.archives-ouvertes.fr/tel-03439538	THESE	NaN	2020IAVF0016	Importance d domesticatic la
4	tel-03439366	1	https://hal.archives-ouvertes.fr/tel-03439366	THESE	NaN	NaN	Homor Cryptograp Privacy,Cryp

5 rows × 21 columns

In [10]:

```
1 #convert urls to list
2 list_urls=df["uri_s"].values.tolist()
```

In [11]:

```
1 #convert ids to list
2 list_ids=df["halId_s"].values.tolist()
```

In [13]:

```
1 list_ids[:10]
```

Out[13]:

```
['tel-03440243',
 'tel-03440181',
 'tel-03440058',
 'tel-03439538',
 'tel-03439366',
 'tel-03439358',
 'tel-03439354',
 'tel-03439346',
 'tel-03439261',
 'tel-03438938']
```

list_urls

In [8]:

```
1 #we need 50 urls, so we just subset the first 50
```

In [7]:

```
1 #subset 50 ids from list  
2 list_ids = list_ids[:50]
```

In [9]:

```
1 count=0  
2 #first 50 urls  
3 for i in list_urls[:50]:  
4  
5     url = i + "/document"  
6     #download pdfs  
7     urllib.request.urlretrieve(url, "{}.pdf".format(list_ids[count]))  
8     print(count,end="\r")  
9     count +=1
```

Convert pdf to text files

In [15]:

```
1 extension = "pdf"  
2 #Prompt directory in which all files were downloaded,  
3 prompter = promptlib.Files()  
4 #Set working directory to chosen directory  
5 dir = prompter.dir()  
6 os.chdir(dir)  
7 #Get list of all pdf files in directory  
8 all_filenames = [i for i in glob.glob('*.{}'.format(extension))]
```

In [34]:

```
1 #downloaded excess files so subset 50 out of the list  
2 all_filenames=all_filenames[5:55]
```

In [41]:

```
1 for file in all_filenames:  
2     with open(file, "rb") as f:  
3         pdf = pdftotext.PDF(f)  
4         # Save all text to a txt file.  
5         split_string = file.split(".", 1)  
6         substring = split_string[0]  
7         with open('{}'.format(substring+".txt"), 'w',encoding="utf-8") as f:  
8             f.write("\n\n".join(pdf))
```

In [41]:

```
1 for file in all_filenames:
2     with open(file, "rb") as f:
3         pdf = pdftotext.PDF(f)
4         # Save all text to a txt file.
5         split_string = file.split(".", 1)
6         substring = split_string[0]
7         with open('{}'.format(substring+".txt"), 'w', encoding="utf-8") as f:
8             f.write("\n\n".join(pdf))
```

LANGUAGE DETECTION

In [46]:

```
1 #download stopwords from nltk
2 nltk.download('stopwords')
3
```

```
[nltk_data] Downloading package stopwords to
[nltk_data] C:\Users\Administrator\AppData\Roaming\nltk_data...
[nltk_data] Unzipping corpora\stopwords.zip.
```

Out[46]:

True

In [47]:

```
1 stopwords.fileids()
2
```

Out[47]:

```
['arabic',
 'azerbaijani',
 'bengali',
 'danish',
 'dutch',
 'english',
 'finnish',
 'french',
 'german',
 'greek',
 'hungarian',
 'indonesian',
 'italian',
 'kazakh',
 'nepali',
 'norwegian',
 'portuguese',
 'romanian',
 'russian',
 'slovene',
 'spanish',
 'swedish',
 'tajik',
 'turkish']
```

In [13]:

```

1  # reference : https://blog.alejandronolla.com/2013/05/15/detecting-text-language-with-
2
3  try:
4      from nltk import wordpunct_tokenize
5      from nltk.corpus import stopwords
6  except ImportError:
7      print ('[!] You need to install nltk (http://nltk.org/index.html)')
8
9
10
11  #-----
12  def _calculate_languages_ratios(text):
13      """
14      Calculate probability of given text to be written in several languages and
15      return a dictionary that looks like {'french': 2, 'spanish': 4, 'english': 0}
16
17      @param text: Text whose language want to be detected
18      @type text: str
19
20      @return: Dictionary with languages and unique stopwords seen in analyzed text
21      @rtype: dict
22      """
23
24      languages_ratios = {}
25
26      ...
27      nltk.wordpunct_tokenize() splits all punctuations into separate tokens
28
29      >>> wordpunct_tokenize("That's thirty minutes away. I'll be there in ten.")
30      ['That', "'", 's', 'thirty', 'minutes', 'away', '.', 'I', "'", 'll', 'be', 'there',
31      '']
32
33      tokens = wordpunct_tokenize(text)
34      words = [word.lower() for word in tokens]
35
36      # Compute per language included in nltk number of unique stopwords appearing in analyzed
37      for language in stopwords.fileids():
38          stopwords_set = set(stopwords.words(language))
39          words_set = set(words)
40          common_elements = words_set.intersection(stopwords_set)
41
42          languages_ratios[language] = len(common_elements) # language "score"
43
44      return languages_ratios
45
46
47  #-----
48  def detect_language(text):
49      """
50      Calculate probability of given text to be written in several languages and
51      return the highest scored.
52
53      It uses a stopwords based approach, counting how many unique stopwords
54      are seen in analyzed text.
55
56      @param text: Text whose language want to be detected
57      @type text: str
58
59      @return: Most scored language guessed

```

```
60     @rtype: str
61     """
62
63     ratios = _calculate_languages_ratios(text)
64
65     most Rated language = max(ratios, key=ratios.get)
66
67     return most Rated language
68
69
70
71
```

In [63]:

```
1 # Apply to our text files
```

In [14]:

```
1 extension = "txt"
2 #Prompt directory in which all files were downloaded,
3 prompter = promptlib.Files()
4 #Set working directory to chosen directory
5 dir = prompter.dir()
6 os.chdir(dir)
7 #Get list of all pdf files
8 all_filenames = [i for i in glob.glob('*.{}'.format(extension))]
```

In [15]:

```
1 all_filenames[:10]
```

Out[15]:

```
['tel-03435883.txt',
'tel-03435884.txt',
'tel-03435885.txt',
'tel-03435936.txt',
'tel-03436011.txt',
'tel-03436023.txt',
'tel-03436024.txt',
'tel-03436025.txt',
'tel-03436087.txt',
'tel-03436137.txt']
```

In [15]:

```
1 #detect languages in 50 theses
2 for file in all_filenames:
3     with open(file, "r",encoding="utf-8") as f:
4         #read file as text and store in variable `text`
5         text = f.read()
6         split_string = file.split(".", 1)
7         substring = split_string[0]
8         language = detect_language(text)
9         print("Thesis ID :{}".format(substring) +" Language : {}".format(language))
10
```

```
Thesis ID :tel-03435883 Language : english
Thesis ID :tel-03435884 Language : english
Thesis ID :tel-03435885 Language : english
Thesis ID :tel-03435936 Language : english
Thesis ID :tel-03436011 Language : french
Thesis ID :tel-03436023 Language : english
Thesis ID :tel-03436024 Language : english
Thesis ID :tel-03436025 Language : french
Thesis ID :tel-03436087 Language : english
Thesis ID :tel-03436137 Language : english
Thesis ID :tel-03436157 Language : english
Thesis ID :tel-03436173 Language : english
Thesis ID :tel-03436335 Language : french
Thesis ID :tel-03436364 Language : french
Thesis ID :tel-03436368 Language : french
Thesis ID :tel-03436372 Language : french
Thesis ID :tel-03436394 Language : romanian
Thesis ID :tel-03436405 Language : french
Thesis ID :tel-03436409 Language : english
Thesis ID :tel-03436501 Language : english
Thesis ID :tel-03436527 Language : english
Thesis ID :tel-03436530 Language : english
Thesis ID :tel-03436542 Language : french
Thesis ID :tel-03436545 Language : french
Thesis ID :tel-03436548 Language : french
Thesis ID :tel-03436551 Language : french
Thesis ID :tel-03437063 Language : english
Thesis ID :tel-03437096 Language : english
Thesis ID :tel-03437282 Language : english
Thesis ID :tel-03437572 Language : english
Thesis ID :tel-03437573 Language : english
Thesis ID :tel-03437616 Language : french
Thesis ID :tel-03438100 Language : portuguese
Thesis ID :tel-03438101 Language : french
Thesis ID :tel-03438102 Language : french
Thesis ID :tel-03438103 Language : english
Thesis ID :tel-03438104 Language : english
Thesis ID :tel-03438105 Language : french
Thesis ID :tel-03438755 Language : english
Thesis ID :tel-03438811 Language : portuguese
Thesis ID :tel-03438828 Language : english
Thesis ID :tel-03438829 Language : portuguese
Thesis ID :tel-03438863 Language : english
Thesis ID :tel-03438921 Language : english
Thesis ID :tel-03438923 Language : french
Thesis ID :tel-03438925 Language : english
Thesis ID :tel-03438938 Language : french
Thesis ID :tel-03439261 Language : english
```

Thesis ID :tel-03439346 Language : french
Thesis ID :tel-03439354 Language : english

Process the data and use TF-IDF and cosine to assess similarity between the documents

In [16]:

```
1 from sklearn.feature_extraction.text import TfidfVectorizer
2 from sklearn.feature_extraction.text import CountVectorizer
3 import pandas as pd
4 pd.set_option("max_rows", 600)
5 from pathlib import Path
6 import glob
```

In [17]:

```
1 directory_path = "H:\\Downloads\\Datatsets\\Tel\\Tel_text"
2 text_files = glob.glob(f"{directory_path}/*.txt")
3
```

In [18]:

```
1 text_files [:10]
```

Out[18]:

```
['H:\\Downloads\\Datatsets\\Tel\\Tel_text\\tel-03435883.txt',
'H:\\Downloads\\Datatsets\\Tel\\Tel_text\\tel-03435884.txt',
'H:\\Downloads\\Datatsets\\Tel\\Tel_text\\tel-03435885.txt',
'H:\\Downloads\\Datatsets\\Tel\\Tel_text\\tel-03435936.txt',
'H:\\Downloads\\Datatsets\\Tel\\Tel_text\\tel-03436011.txt',
'H:\\Downloads\\Datatsets\\Tel\\Tel_text\\tel-03436023.txt',
'H:\\Downloads\\Datatsets\\Tel\\Tel_text\\tel-03436024.txt',
'H:\\Downloads\\Datatsets\\Tel\\Tel_text\\tel-03436025.txt',
'H:\\Downloads\\Datatsets\\Tel\\Tel_text\\tel-03436087.txt',
'H:\\Downloads\\Datatsets\\Tel\\Tel_text\\tel-03436137.txt']
```

In [19]:

```
1 text_titles = [Path(text).stem for text in text_files]
2
```

In [20]:

```
1 from spacy.lang.fr.stop_words import STOP_WORDS as fr_stop
2 from spacy.lang.en.stop_words import STOP_WORDS as en_stop
3 from spacy.lang.ro.stop_words import STOP_WORDS as ro_stop
4 from spacy.lang.pt.stop_words import STOP_WORDS as pt_stop
5
6 #use stopwords from spacy
7 final_stopwords_list = list(fr_stop) + list(en_stop) + list(ro_stop) + list(pt_stop)
8 # from nltk.corpus import stopwords
9
10 # final_stopwords_list = stopwords.words('english') + stopwords.words('french') + stopw
```


In [21]:

```
1 #instantiate vectorizer with stopwords as paramter
2 tfidf_vectorizer = TfidfVectorizer(input='filename', stop_words=final_stopwords_list)
3
```

In [22]:

```
1 #fit vectorizer on all text files
2 tfidf_vector = tfidf_vectorizer.fit_transform(text_files)
3
```

In [23]:

```
1 #Make a DataFrame out of the resulting tf-idf vector, row names as index
2 tfidf_df = pd.DataFrame(tfidf_vector.toarray(), index=text_titles, columns=tfidf_vector
3
```

In [24]:

```
1 #Add column for document frequency aka number of times word appears in all documents
2
3 tfidf_df.loc['00_Document Frequency'] = (tfidf_df > 0).sum()
4
```

In [25]:

```
1 # drop "00_Document Frequency" since we were just using it for illustration purposes.
2 tfidf_df = tfidf_df.drop('00_Document Frequency', errors='ignore')
3
```

In [26]:

```
1 #reorganize the DataFrame so that the words are in rows rather than columns.
2
3 tfidf_df.stack().reset_index()
4 tfidf_df = tfidf_df.stack().reset_index()
5
```

In [27]:

```
1 #rename columns with tfidf, document, term
2 tfidf_df = tfidf_df.rename(columns={0:'tfidf', 'level_0': 'document', 'level_1': 'term',
3
```

In [28]:

```
1 #sort values such that top 5 tf-idf of every document is obtained
2 tfidf_df=tfidf_df.sort_values(by=['document', 'tfidf'], ascending=[True, False]).groupby(
3
```

In [29]:

```
1 tfidf_df.head(10)
```

Out[29]:

	document	term	tfidf
162861	tel-03435883	toxine	0.348011
108520	tel-03435883	meliloti	0.240458
168878	tel-03435883	vapc10	0.223557
162860	tel-03435883	toxin	0.198698
168876	tel-03435883	vapc	0.196669
242882	tel-03435884	ehv	0.755768
263575	tel-03435884	herpesvirus	0.208099
247930	tel-03435884	equine	0.178468
350571	tel-03435884	virus	0.126986
264757	tel-03435884	horse	0.124838

In [30]:

```

1 import altair as alt
2 import numpy as np
3 term_list=[]
4
5 # adding a little randomness to break ties in term ranking
6 top_tfidf_plusRand = tfidf_df.copy()
7 top_tfidf_plusRand['tfidf'] = top_tfidf_plusRand['tfidf'] + np.random.rand(tfidf_df.sh
8
9 # base for all visualizations, with rank calculation
10 base = alt.Chart(top_tfidf_plusRand).encode(
11     x = 'rank:O',
12     y = 'document:N'
13 ).transform_window(
14     rank = "rank()",
15     sort = [alt.SortField("tfidf", order="descending")],
16     groupby = ["document"],
17 )
18
19 # heatmap specification
20 heatmap = base.mark_rect().encode(
21     color = 'tfidf:Q'
22 )
23
24 # red circle over terms in above list
25 circle = base.mark_circle(size=100).encode(
26     color = alt.condition(
27         alt.FieldOneOfPredicate(field='term', oneOf=term_list),
28         alt.value('red'),
29         alt.value('#FFFFFF00')
30     )
31 )
32
33 # text labels, white for darker heatmap colors
34 text = base.mark_text(baseline='middle').encode(
35     text = 'term:N',
36     color = alt.condition(alt.datum.tfidf >= 0.23, alt.value('white'), alt.value('black
37 )
38
39 # display the three superimposed visualizations
40 (heatmap + circle + text).properties(width = 600)

```

Out[30]:

document	tel-03435883	toxine	meliloti	vapc10	toxin	
	tel-03435884	ehv	herpesvirus	equine	virus	
	tel-03435885	protein	nmr	proteins	viral	
	tel-03435936	mosaïque	rhoa	patients	anomalies	
	tel-03436011	ceser	assemblée	brizio	conseillers	
	tel-03436023	temperature	illumination	optical	figure	
	tel-03436024	collisions	pb	tev	rapidity	
	tel-03436025	verre	zno	érosion	films	
	tel-03436087	plessz	pratiques	sociale	marie	
	tel-03436137	height	t0	westcott	limit	
	tel-03436157	wave	video	beach	cbathy	
	tel-03436173	patent	trade	citations	news	
	tel-03436335	sénescence	vulgaire	senescence	reproduction	
	tel-03436364	mode	appariement	utilite	galite	
	tel-03436368	sûreté	culture	conduite	manager	
	tel-03436372	pénal	prison	réinsertion	foucault	
	tel-03436394	rl	iqu	efu	td	
	tel-03436405	odin	ugtg	militants	travailleurs	
	tel-03436409	msa	speech	psp	pd	
	tel-03436501	thermal	energies	electronic	correction	
	tel-03436527	contrats	renégociation	concession	contrat	re
	tel-03436530	minority	minorities	manager	managers	
	tel-03436542	andra	déchets	sûreté	stockage	
	tel-03436545	consommation	biens	castes	préférences	
	tel-03436548	mongolie	mongoles	développement	tolgoi	
	tel-03436551	taire	mone	politique	macroe	
	tel-03437063	p53	chikv	infection	viral	
	tel-03437096	patient	real	virtual	dt	
	tel-03437282	framing	sanctions	symbolic	subjects	
	tel-03437572	cpz	liposomes	cd	drv	
	tel-03437573	ch2	fatty	aldehydes	catalyst	
	tel-03437616	sepsis	lactate	vegf	patients	
	tel-03438100	amazonie	brésil	coudreau	verissimo	
	tel-03438101	patients	alcohol	tual	alcool	
	tel-03438102	vo2	température	vanadium	minces	
	tel-03438103	postnatal	pregnancy	exposome	rownames	
	tel-03438104	bones	shape	distal	rhinos	
	tel-03438105	entartrage	carbonate	caco3	zinc	
	tel-03438755	control	quadrotor	vehicle	drone	
	tel-03438811	agricultores	agricultura	agricultor	ecológicos	
	tel-03438828	toy	learning	tool	goal	
	tel-03438829	tieta	novela	beata	beatas	
	tel-03438863	causal	xt	pq	time	
	tel-03438921	stirling	cooling	heat	regenerator	
	tel-03438923	lisa	10	bruit	sgwb	c
	tel-03438925	lungs	tree	hfcwo	mucus	
	tel-03438938	suspensions	particules	viscosité	pdms	
	tel-03439261	pendulum	zaouali	rott	energy	
	tel-03439346	droit	art	artiste	œuvre	
	tel-03439354	credit	bank	firm	branch	
						rank



N-GRAM

In [31]:

```
1 #method to generate n-grams:  
2 #params:  
3 #text-the text for which we have to generate n-grams  
4 #ngram-number of grams to be generated from the text(1,2,3,4 etc., default value=1)
```

In [32]:

```
1 from spacy.lang.fr.stop_words import STOP_WORDS as fr_stop
2 from spacy.lang.en.stop_words import STOP_WORDS as en_stop
3 from spacy.lang.ro.stop_words import STOP_WORDS as ro_stop
4 from spacy.lang.pt.stop_words import STOP_WORDS as pt_stop
5
6 # use stopwords from spacy
7 final_stopwords_list = list(fr_stop) + list(en_stop) + list(ro_stop) + list(pt_stop)
8 # Calculate bigram by using ngram_range=(2,2)
9
10 #instance vectorizer with stopwords
11 tfidf_vectorizer = TfidfVectorizer(input='filename', stop_words=final_stopwords_list, r
12
13 #fit vectorizer with all the text files
14 tfidf_vector = tfidf_vectorizer.fit_transform(text_files)
15
```

In [33]:

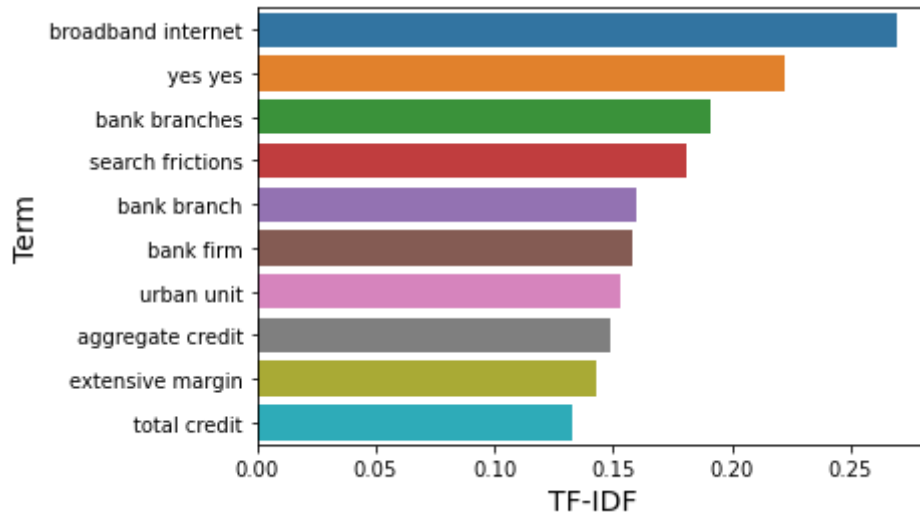
```
1 #create dataframe from vector
2 tfidf_df = pd.DataFrame(tfidf_vector.toarray(), index=text_titles, columns=tfidf_vector
3
4 tfidf_df.loc['00_Document Frequency'] = (tfidf_df > 0).sum()
5 # drop "00_Document Frequency" since we were just using it for illustration purposes.
6 tfidf_df = tfidf_df.drop('00_Document Frequency', errors='ignore')
7 #sort values such that top 5 tf-idf of every document is obtained
8 tfidf_df = tfidf_df.stack().reset_index()
9
10
11 tfidf_df = tfidf_df.rename(columns={0:'tfidf', 'level_0': 'document', 'level_1': 'term',
12
13
14
```

In [34]:

```
1 test_df = (tfidf_df.sort_values(by=['document', 'tfidf'], ascending=[True, False]).groupb
2
```

In [37]:

```
1 # Construct plot
2 p=sns.barplot(x = "tfidf", y = "term", data = test_df)
3 p.set_xlabel("TF-IDF", fontsize = 14)
4 p.set_ylabel("Term",fontsize=14)
5
6 plt.show()
```



In [40]:

```

1 from spacy.lang.fr.stop_words import STOP_WORDS as fr_stop
2 from spacy.lang.en.stop_words import STOP_WORDS as en_stop
3 from spacy.lang.ro.stop_words import STOP_WORDS as ro_stop
4 from spacy.lang.pt.stop_words import STOP_WORDS as pt_stop
5
6
7 final_stopwords_list = list(fr_stop) + list(en_stop) + list(ro_stop) + list(pt_stop)
8
9 tfidf_vectorizer = TfidfVectorizer(input='filename', stop_words=final_stopwords_list, r
10
11
12
13 #testing on 5 files only
14
15 tfidf_vector = tfidf_vectorizer.fit_transform(text_files[:5])
16
17
18 tfidf_df = pd.DataFrame(tfidf_vector.toarray(), index=text_titles[:5], columns=tfidf_ve
19
20
21 tfidf_df.loc['00_Document Frequency'] = (tfidf_df > 0).sum()
22
23
24 tfidf_df = tfidf_df.drop('00_Document Frequency', errors='ignore')
25
26
27 tfidf_df.stack().reset_index()
28
29
30 tfidf_df = tfidf_df.stack().reset_index()
31
32
33 tfidf_df = tfidf_df.rename(columns={0:'tfidf', 'level_0': 'document', 'level_1': 'term',
34
35
36 tfidf_df=tfidf_df.sort_values(by=['document', 'tfidf'], ascending=[True, False]).groupby(
37

```

In [44]:

```
1 tfidf_df.head()
```

Out[44]:

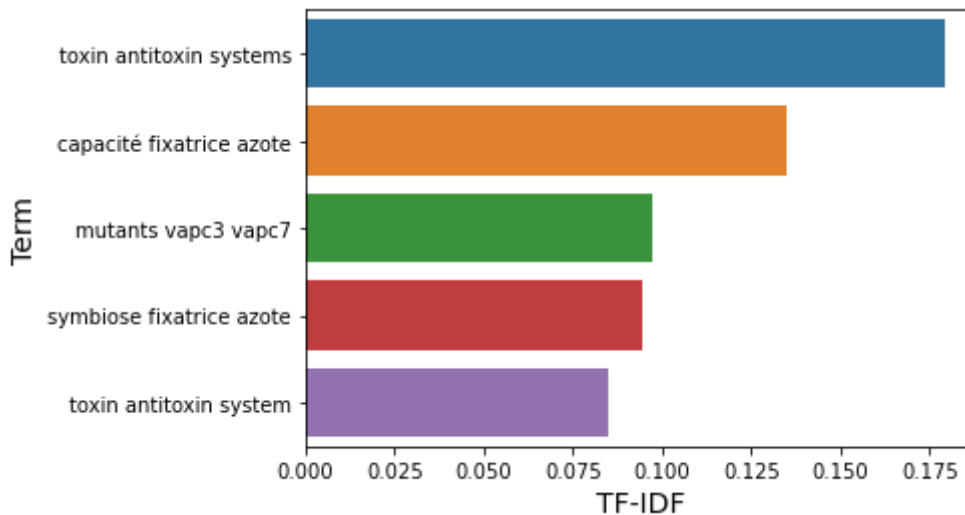
	document	term	tfidf
251726	tel-03435883	toxin antitoxin systems	0.179170
64365	tel-03435883	capacité fixatrice azote	0.135163
167348	tel-03435883	mutants vapc3 vapc7	0.097443
243732	tel-03435883	symbiose fixatrice azote	0.094300
251725	tel-03435883	toxin antitoxin system	0.084870

In [43]:

```
1 test_df2 = tfidf_df.head()
```

In [45]:

```
1 # Construct plot
2 p=sns.barplot(x = "tfidf", y = "term", data = test_df2)
3 p.set_xlabel("TF-IDF", fontsize = 14)
4 p.set_ylabel("Term",fontsize=14)
5
6 plt.show()
```



In [28]:

```
1 #https://studymachinelearning.com/cosine-similarity-text-similarity-metric/
```

check similarity between documents

In [47]:

```
1 #get names of the first 10 docs
2 docs_names=all_filenames[:10]
```

In [48]:

```
1 #read all the 10 text files and append them to a list
2 data=[]
3 for doc in docs_names:
4     with open(doc, "r",encoding="utf-8") as f:
5         #read file as text and store in variable `text`
6         text = f.read()
7         data.append(text)
8
```

In [49]:

```
1 #get the IDS of all the 10 docs
2 doc_names=[]
3 for doc in docs_names:
4
5     split_string = doc.split(".", 1)
6     substring = split_string[0]
7     doc_names.append(substring)
```

In [50]:

```
1 doc_names
```

Out[50]:

```
['tel-03435883',
'tel-03435884',
'tel-03435885',
'tel-03435936',
'tel-03436011',
'tel-03436023',
'tel-03436024',
'tel-03436025',
'tel-03436087',
'tel-03436137']
```

In [51]:

```
1 #convert the IDS to a dataframe
2 df_docnames = pd.DataFrame(doc_names)
```

In [52]:

```
1 df_docnames
```

Out[52]:

	0
0	tel-03435883
1	tel-03435884
2	tel-03435885
3	tel-03435936
4	tel-03436011
5	tel-03436023
6	tel-03436024
7	tel-03436025
8	tel-03436087
9	tel-03436137

In [53]:

```
1 #rename column o to ID
2 df_docnames=df_docnames.rename(columns={0:"halId_s"})
```

In [54]:

```
1 #merge with original dataset
2 df_merged=pd.merge(df_docnames,df,on="halId_s")
```

In [55]:

```
1 #store the names of the theses in a list
2 doc_names=df_merged["title_s"].tolist()
```

In [75]:

```
1 from sklearn.feature_extraction.text import CountVectorizer
2 #use Countvectorizer
3 count_vectorizer = CountVectorizer()
4 vector_matrix = count_vectorizer.fit_transform(data)
5 vector_matrix
```

Out[75]:

<10x53455 sparse matrix of type '<class 'numpy.int64'>' with 93787 stored elements in Compressed Sparse Row format>

In [76]:

```
1 #get names of the features
2 tokens = count_vectorizer.get_feature_names()
3
```

In [77]:

```
1 #convert matrix to a 2d numpy array
2 vector_matrix.toarray()
3
```

Out[77]:

```
array([[ 1,  7,  0, ...,  0,  0,  0],
       [60,  5,  0, ...,  0,  0,  0],
       [ 0, 11,  0, ...,  0,  0,  0],
       ...,
       [43,  1,  0, ...,  0,  0,  0],
       [37, 42,  0, ...,  0,  0,  0],
       [ 0,  0,  0, ...,  0,  0,  0]], dtype=int64)
```

In [78]:

```
1 def create_dataframe(matrix, tokens):
2
3     #doc_names = [f'doc_{i+1}' for i, _ in enumerate(matrix)]
4     df = pd.DataFrame(data=matrix, index=doc_names, columns=tokens)
5     return(df)
```

In [79]:

```
1 create_dataframe(vector_matrix.toarray(),tokens).head()
2
```

Out[79]:

	00	000	0000	0000034996	00000407	0000042954	00000457	00000569
Les systèmes Toxine-Antitoxine VapBC : des régulateurs de la symbiose fixatrice d'azote Rhizobium-Légumineuse,The VapBC Toxin-Antitoxin systems : regulators of the nitrogen-fixing symbiosis Rhizobium-Legume	1	7	0	0	0	0	0	0
Étude phylogénique de souches d'alphaherpèsvirus isolées chez les équidés français et développement d'un outil innovant pour la mesure des anticorps neutralisants après infection ou vaccination,Phylogenic study of equid alphaherpesvirus strains isolated in France and development of an innovative assay for the measurement of neutralising antibodies after infection or vaccination	60	5	0	1	0	0	0	0
Multi-scale studies of Measeles virus nucleocapsid assembly,Etudes multi-échelles de l'assemblage de la nucléocapside du virus de la rougeole	0	11	0	0	0	1	0	0
Caractérisation génomique des anomalies de la pigmentation cutanée en mosaïque,Genomic characterization of mosaic cutaneous pigmentary disorders	8	48	1	0	1	0	1	1

00 000 0000 0000034996 00000407 0000042954 00000457 00000569

Le Conseil
économique, social et
environnemental
régional : assemblée
du dialogue des
intérêts organisés
dans la région, The
Conseil économique,
social et
environnemental
régional : an assembly
of dialogue between
organised interests in
the Région

1 13 0 0 0 0 0 0

5 rows × 53455 columns



In [80]:

```
1 from sklearn.metrics.pairwise import cosine_similarity
2 #calculate the cosine similarity of matrix
3 cosine_similarity_matrix = cosine_similarity(vector_matrix)
```

In [81]:

```
1 #create dataframe from matrix with theses name
2 distance_matrix=create_dataframe(cosine_similarity_matrix,doc_names)
3
```

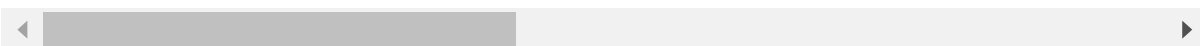
In [82]:

1 distance_matrix

Out[82]:

Les systèmes Toxine-Antitoxine VapBC : des régulateurs de la symbiose fixatrice d'azote Rhizobium- Légumineuse,The VapBC Toxin- Antitoxin systems : regulators of the nitrogen-fixing symbiosis Rhizobium- Legume	Étude phylogénique de souches d'alphaherpèsvirus isolées chez les équidés français et développement d'un outil innovant pour la mesure des anticorps neutralisants après infection ou vaccination,Phylogenic study of equid alphaherpesvirus strains isolated in France and development of an innovative assay for the measurement of neutralising antibodies after infection or vaccination	Multi-scale studies of Measeles virus nucleocapsid assembly,Etudes multi-échelles de l'assemblage de la nucléocapside du virus de la rougeole	C g ar mosa chari mos
Les systèmes Toxine-Antitoxine VapBC : des régulateurs de la symbiose fixatrice d'azote Rhizobium-Légumineuse,The VapBC Toxin-Antitoxin systems : regulators of the nitrogen-fixing symbiosis Rhizobium-Legume	1.000000	0.869968	0.436238
Étude phylogénique de souches d'alphaherpèsvirus isolées chez les équidés français et développement d'un outil innovant pour la mesure des anticorps neutralisants après infection ou vaccination,Phylogenic study of equid alphaherpesvirus strains isolated in France and development of an innovative assay for the measurement of neutralising antibodies after infection or vaccination	0.869968	1.000000	0.520814
Multi-scale studies of Measeles virus nucleocapsid assembly,Etudes multi-échelles de l'assemblage de la nucléocapside du virus de la rougeole	0.436238	0.520814	1.000000
Caractérisation génomique des anomalies de la pigmentation cutanée en mosaïque,Genomic characterization of mosaic cutaneous pigmentary disorders	0.845360	0.839129	0.624248

Le Conseil économique\, social et environnemental régional : assemblée du dialogue des intérêts organisés dans la région,The Conseil économique\, social et environnemental régional : an assembly of dialogue between organised interests in the Région	0.826060	0.714074	0.054021
Optical Developments for Microscale Measurement and Control of Temperature in Optogenetics,Développements optiques pour la mesure et le contrôle micrométrique de la température en optogénétique	0.390573	0.459029	0.902244
Z-boson and double charm production with ALICE at the LHC,Production des bosons Z et du double charme avec ALICE auprès du LHC	0.408820	0.473688	0.893786
Caractérisation opto-mécanique du verre traité par des méthodes thermo-chimiques,Opto-mechanical characterization of glass treated by thermochemical methods	0.856415	0.745178	0.137838
La Dynamique sociale des pratiques : stratification sociale\, changement social et consommation alimentaire,The Social Dynamics of Practices : social stratification\, social change and food consumption	0.858910	0.765129	0.164582
Random surface growth models : hydrodynamic limits and fluctuations,Modèles de croissance de surfaces aléatoires : limites hydrodynamiques et fluctuations	0.425379	0.487731	0.896544



In [83]:

```
1 distance_matrix.to_csv(r'H:\Downloads\Datatsets\Tel\dist_m.csv',index=True)
2
```

In [84]:

```
1 distance_matrix.reset_index().to_csv(r'H:\Downloads\Datatsets\Tel\dist_m.csv',index=True)
```

Similarity with tf-idf vectorizer instead of count vectorizer

In [85]:

```

1 from sklearn.feature_extraction.text import TfidfVectorizer
2 #same procedure as count vectorizer but we use tfidf-vectorizer which is better
3 Tfidf_vect = TfidfVectorizer()
4 vector_matrix = Tfidf_vect.fit_transform(data)
5
6 tokens = Tfidf_vect.get_feature_names()
7 create_dataframe(vector_matrix.toarray(),tokens).head()

```

Out[85]:

	00	000	0000	0000034996	00000407	0000042954	00000
Les systèmes Toxine-Antitoxine VapBC : des régulateurs de la symbiose fixatrice d'azote Rhizobium-Légumineuse,The VapBC Toxin-Antitoxin systems : regulators of the nitrogen-fixing symbiosis Rhizobium-Legume	0.000175	0.001225	0.000000	0.000000	0.000000	0.000000	0.000
Étude phylogénique de souches d'alphaherpèsvirus isolées chez les équidés français et développement d'un outil innovant pour la mesure des anticorps neutralisants après infection ou vaccination,Phylogenic study of equid alphaherpesvirus strains isolated in France and development of an innovative assay for the measurement of neutralising antibodies after infection or vaccination	0.007519	0.000627	0.000000	0.000257	0.000000	0.000000	0.000
Multi-scale studies of Measeles virus nucleocapsid assembly,Etudes multi-échelles de l'assemblage de la nucléocapside du virus de la rougeole	0.000000	0.002206	0.000000	0.000000	0.000000	0.000411	0.000
Caractérisation génomique des anomalies de la pigmentation cutanée en mosaïque,Genomic characterization of mosaic cutaneous pigmentary disorders	0.002571	0.015426	0.000659	0.000000	0.000659	0.000000	0.000

	00	000	0000	0000034996	00000407	0000042954	00000
Le Conseil économique, social et environnemental régional : assemblée du dialogue des intérêts organisés dans la région,The Conseil économique, social et environnemental régional : an assembly of dialogue between organised interests in the Région	0.000089	0.001156	0.000000	0.000000	0.000000	0.000000	0.000

5 rows × 53455 columns



In [86]:

```

1 #get cosine similarity from applying on sparse matrix + create dataframe from matrix
2 cosine_similarity_matrix = cosine_similarity(vector_matrix)
3 create_dataframe(cosine_similarity_matrix,doc_names)

```

Out[86]:

Les systèmes Toxine-Antitoxine VapBC : des régulateurs de la symbiose fixatrice d'azote Rhizobium- Légumineuse,The VapBC Toxin- Antitoxin systems : regulators of the nitrogen-fixing symbiosis Rhizobium- Legume	Étude phylogénique de souches d'alphaherpèsvirus isolées chez les équidés français et développement d'un outil innovant pour la mesure des anticorps neutralisants après infection ou vaccination,Phylogenic study of equid alphaherpesvirus strains isolated in France and development of an innovative assay for the measurement of neutralising antibodies after infection or vaccination	Multi-scale studies of Measeles virus nucleocapsid assembly,Etudes multi-échelles de l'assemblage de la nucléocapside du virus de la rougeole	C g ar mosa chari mos
Les systèmes Toxine-Antitoxine VapBC : des régulateurs de la symbiose fixatrice d'azote Rhizobium-Légumineuse,The VapBC Toxin-Antitoxin systems : regulators of the nitrogen-fixing symbiosis Rhizobium-Legume	1.000000	0.615478	0.375844
Étude phylogénique de souches d'alphaherpèsvirus isolées chez les équidés français et développement d'un outil innovant pour la mesure des anticorps neutralisants après infection ou vaccination,Phylogenic study of equid alphaherpesvirus strains isolated in France and development of an innovative assay for the measurement of neutralising antibodies after infection or vaccination	0.615478	1.000000	0.389180
Multi-scale studies of Measeles virus nucleocapsid assembly,Etudes multi-échelles de l'assemblage de la nucléocapside du virus de la rougeole	0.375844	0.389180	1.000000
Caractérisation génomique des anomalies de la pigmentation cutanée en mosaïque,Genomic characterization of mosaic cutaneous pigmentary disorders	0.686478	0.587794	0.530602

Le Conseil économique\, social et environnemental régional : assemblée du dialogue des intérêts organisés dans la région,The Conseil économique\, social et environnemental régional : an assembly of dialogue between organised interests in the Région	0.689506	0.511317	0.046973
Optical Developments for Microscale Measurement and Control of Temperature in Optogenetics,Développements optiques pour la mesure et le contrôle micrométrique de la température en optogénétique	0.333542	0.336816	0.813004
Z-boson and double charm production with ALICE at the LHC,Production des bosons Z et du double charme avec ALICE auprès du LHC	0.343989	0.341891	0.790952
Caractérisation opto-mécanique du verre traité par des méthodes thermo-chimiques,Opto-mechanical characterization of glass treated by thermochemical methods	0.716493	0.534202	0.122538
La Dynamique sociale des pratiques : stratification sociale\, changement social et consommation alimentaire,The Social Dynamics of Practices : social stratification\, social change and food consumption	0.727890	0.557179	0.145243
Random surface growth models : hydrodynamic limits and fluctuations,Modèles de croissance de surfaces aléatoires : limites hydrodynamiques et fluctuations	0.352988	0.346750	0.782971

In [87]:

```
1 #store matrix to export as csv
2 distance_matrix2 = create_dataframe(cosine_similarity_matrix,doc_names)
```

In [88]:

```
1 distance_matrix2.to_csv(r'H:\Downloads\Datatsets\Tel\dist_m2.csv',index=True)
2
```


Force-directed graph

MEENOWA Sarvesh

09/12/2021

```
library(readr)
```

```
## Warning: package 'readr' was built under R version 4.0.5
```

```
library(qgraph)
```

```
## Warning: package 'qgraph' was built under R version 4.0.5
```

```
## cosine similarity
```

```
nba <- read_csv("H:/Downloads/Datatasets/Tel/dist_m.csv", locale = readr::locale(encoding = "utf-8"))
```

```
## New names:
```

```
## * ' ' -> ...1
```

```
## Rows: 10 Columns: 12
```

```
## -- Column specification -----
```

```
## Delimiter: ","
```

```
## chr (1): index
```

```
## dbl (11): ...1, Les systèmes Toxine-Antitoxine VapBC : des régulateurs de la...
```

```
##
```

```
## i Use 'spec()' to retrieve the full column specification for this data.
```

```
## i Specify the column types or set 'show_col_types = FALSE' to quiet this message.
```

```
rownames(nba) <- nba$index
```

```
## Warning: Setting row names on a tibble is deprecated.
```

```
colnames(nba)
```

```
## [1] "...1"
```

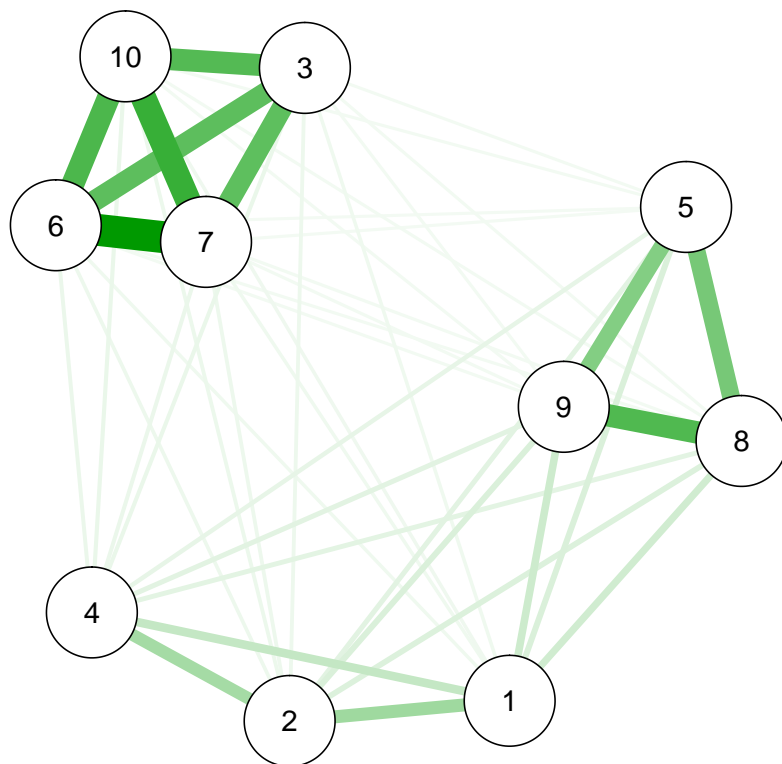
```
## [2] "index"
```

```
## [3] "Les systèmes Toxine-Antitoxine VapBC : des régulateurs de la symbiose fixatrice d'azote Rhizob"
```

```
## [4] "Étude phylogénique de souches d'alphaherpèsvirus isolées chez les équidés français et développ"
```

```
## [5] "Multi-scale studies of Measeles virus nucleocapsid assembly,Etudes multi-échelles de l'assembl
## [6] "Caractérisation génomique des anomalies de la pigmentation cutanée en mosaïque,Genomic charact
## [7] "Le Conseil économique\\, social et environnemental régional : assemblée du dialogue des intérêt
## [8] "Optical Developments for Microscale Measurement and Control of Temperature in Optogenetics,Déve
## [9] "Z-boson and double charm production with ALICE at the LHC,Production des bosons Z et du double
## [10] "Caractérisation opto-mécanique du verre traité par des méthodes thermo-chimiques,Opto-mechanica
## [11] "La Dynamique sociale des pratiques : stratification sociale\\, changement social et consommati
## [12] "Random surface growth models : hydrodynamic limits and fluctuations,Modèles de croissance de s
```

```
test= subset(nba, select = -c(...1,index))
#convert distance to coordinates and convert to matrix
dist_m <- as.matrix(dist(test))
dist_mi <- 1/dist_m # one over, as qgraph takes similarity matrices as input
#jpeg('voronoid_count.jpg', width=1000, height=1000, unit='px')
list_names = nba$index
qgraph(dist_mi, layout='spring', vsize=8,legend=TRUE,nodeNames=list_names,esize=20)
```



- 1: Les systèmes Toxir
- 2: Étude phylogénique
- 3: Multi-scale studies
- 4: Caractérisation gén
- 5: Le Conseil économ
- 6: Optical Developme
- 7: Z-boson and doubl
- 8: Caractérisation opt
- 9: La Dynamique soci
- 10: Random surface g

Database creation to generate MCQ exercises

In [1]:

```
1 import stanza
2 import spacy
3 import pandas as pd
4 import spacy_stanza
5 import os
6 import warnings
7 # import mlconjug3
8 import mlconjug3
9 warnings.filterwarnings("ignore")
```

In [2]:

```
1 #import excel corpus
2 df = pd.read_excel("7000 sentences Corpus with IDs.xlsx")
```

In [3]:

```
1 #get path
2 os.getcwd()
```

Out[3]:

'C:\\Users\\Administrator\\stanza_models'

In [4]:

```
1 #stanza.download(lang='en', model_dir='./stanza_models')
2 nlp_en = spacy_stanza.load_pipeline("en", dir = 'C:\\Users\\Administrator\\stanza_model
```

2021-12-10 00:01:30 INFO: Loading these models for language: en (English):

```
=====
| Processor | Package |
-----
| tokenize | combined |
| pos      | combined |
| lemma    | combined |
| depparse | combined |
| sentiment | sstplus  |
| ner      | ontonotes |
=====
```

```
2021-12-10 00:01:30 INFO: Use device: cpu
2021-12-10 00:01:30 INFO: Loading: tokenize
2021-12-10 00:01:30 INFO: Loading: pos
2021-12-10 00:01:32 INFO: Loading: lemma
2021-12-10 00:01:33 INFO: Loading: depparse
2021-12-10 00:01:35 INFO: Loading: sentiment
2021-12-10 00:01:37 INFO: Loading: ner
2021-12-10 00:01:41 INFO: Done loading processors!
```

In [5]:

```
1 df_eng = df["English"].dropna()
```

In [13]:

```
1 list_df_noun = []
2 list_df_verb = []
3
4 for phrase in df_eng:
5
6     list_noun = []
7     list_verb = []
8     doc = nlp_en(phrase)
9     for token in doc:
10         # Get noun and number
11         if token.pos_ == 'NOUN' and len(token.morph.get("Number")) > 0:
12             list_noun.append({token.text: token.morph.get("Number")[0]})
13
14         # Get verb
15         if token.pos_ == 'VERB' and len(token.morph.get("Tense")) > 0:
16             prefix = ''
17             if len(token.morph.get("Voice")) > 0 and token.morph.get("Voice")[0] == "Pa
18                 temp_token = token
19                 while temp_token.nbor(-1).pos_ == "AUX":
20                     prefix = temp_token.nbor(-1).text + ' ' + prefix
21                     temp_token = temp_token.nbor(-1)
22
23             list_verb.append({prefix + token.text: token.morph.get("Tense")[0]})
24
25     # Add noun
26     for noun in list_noun:
27         for key, value in noun.items():
28             list_df_noun.append([phrase, key, value])
29
30     # Add verb
31     for verb in list_verb:
32         for key, value in verb.items():
33             list_df_verb.append([phrase, key, value])
```


In [27]:

```
1 df_noun = pd.DataFrame(list_df_noun, columns=["Phrase", "Noun", "Number"])
2 df_noun
```

Out[27]:

	Phrase	Noun	Number
0	The beauty of the landscape struck the travell...	beauty	Sing
1	The beauty of the landscape struck the travell...	landscape	Sing
2	The beauty of the landscape struck the travell...	travellers	Plur
3	Nobody knows the truth about this affair.	truth	Sing
4	Nobody knows the truth about this affair.	affair	Sing
...
5734	Computer scientists find a job quickly enough.	Computer	Sing
5735	Computer scientists find a job quickly enough.	scientists	Plur
5736	Computer scientists find a job quickly enough.	job	Sing
5737	Shoemakers rapair shoes.	Shoemakers	Plur
5738	Shoemakers rapair shoes.	shoes	Plur

5739 rows × 3 columns

In [28]:

```
1 #since getting verbs take time, we convert it to csv
2 df_noun.to_csv("df_noun.csv",index=False)
```

In [8]:

```
1 df_noun = pd.read_csv("df_noun.csv")
```

In [22]:

```

1
2 df_verb = pd.DataFrame(list_df_verb, columns=["Phrase", "Verb", "Tense"])
3 df_verb

```

Out[22]:

	Phrase	Verb	Tense
0	The beauty of the landscape struck the travell...	struck	Past
1	Nobody knows the truth about this affair.	knows	Pres
2	In a dictatorship, freedom of expression is li...	is limited	Past
3	His wickedness had no limits.	had	Past
4	His elegance impressed the assembly.	impressed	Past
...
2215	Teachers teach in primary schools.	teach	Pres
2216	The plumber is going to come this afternoon.	going	Pres
2217	He quit his job because his salary was too low.	quit	Past
2218	Computer scientists find a job quickly enough.	find	Pres
2219	Shoemakers rapair shoes.	rapair	Pres

2220 rows × 3 columns

In [31]:

```

1 # getting the tense and verb takes time, so convert it to csv
2 df_verb.to_csv('df_verb.csv',index=False)

```

In [10]:

```
1 df_verb = pd.read_csv('df_verb.csv')
```

In [8]:

```

1 #create function to lemmatize the verb
2 def get_lemma(x):
3     doc=nlp_en(x)
4     for token in doc:
5         return token.lemma_
6
7

```

In [9]:

```

1 #test
2 get_lemma("struck")

```

Out[9]:

'strike'

In [11]:

```
1 #apply function to verb column to obtain lemma
2 df_verb["Lemma"]=df_verb["Verb"].apply(lambda x:get_lemma(x))
```

In [12]:

```
1 #Lemma function application takes time as well, so we conver it as csv a well
2 df_verb.to_csv('df_verb.csv',index=False)
```

In [11]:

```
1 df_verb = pd.read_csv('df_verb.csv')
```

In [12]:

```
1 df_verb
```

Out[12]:

	Phrase	Verb	Tense	Lemma
0	The beauty of the landscape struck the travell...	struck	Past	strike
1	Nobody knows the truth about this affair.	knows	Pres	know
2	In a dictatorship, freedom of expression is li...	is limited	Past	be
3	His wickedness had no limits.	had	Past	have
4	His elegance impressed the assembly.	impressed	Past	impressed
...
2215	Teachers teach in primary schools.	teach	Pres	teach
2216	The plumber is going to come this afternoon.	going	Pres	go
2217	He quit his job because his salary was too low.	quit	Past	quit
2218	Computer scientists find a job quickly enough.	find	Pres	find
2219	Shoemakers rapair shoes.	rapair	Pres	rapair

2220 rows × 4 columns

In [14]:

```
1 # choose language to set as default conjugator
2 default_conjugator = mlconjug3.Conjugator(language='en')
3
```

In [15]:

```

1 #test instantiated function on verb "know"
2 test_verb = default_conjugator.conjugate("know")
3 all_conjugated_forms = test_verb.iterate()
4 #results gives a list of tuples
5 print(all_conjugated_forms)

```

```

[('indicative', 'indicative present', '1s', 'know'), ('indicative', 'indicative present', '2s', 'know'), ('indicative', 'indicative present', '3s', 'knows'), ('indicative', 'indicative present', '1p', 'know'), ('indicative', 'indicative present', '2p', 'know'), ('indicative', 'indicative present', '3p', 'know'), ('indicative', 'indicative past tense', '1s', 'knew'), ('indicative', 'indicative past tense', '2s', 'knew'), ('indicative', 'indicative past tense', '3s', 'knew'), ('indicative', 'indicative past tense', '1p', 'knew'), ('indicative', 'indicative past tense', '2p', 'knew'), ('indicative', 'indicative past tense', '3p', 'knew'), ('indicative', 'indicative present continuous', '1s', 'knowing'), ('indicative', 'indicative present continuous', '2s', 'knowing'), ('indicative', 'indicative present continuous', '3s', 'knowing'), ('indicative', 'indicative present continuous', '1p', 'knowing'), ('indicative', 'indicative present continuous', '2p', 'knowing'), ('indicative', 'indicative present continuous', '3p', 'knowing'), ('indicative', 'indicative present perfect', '1s', 'known'), ('indicative', 'indicative present perfect', '2s', 'known'), ('indicative', 'indicative present perfect', '3s', 'known'), ('indicative', 'indicative present perfect', '1p', 'known'), ('indicative', 'indicative present perfect', '2p', 'known'), ('indicative', 'indicative present perfect', '3p', 'known'), ('infinitive', 'infinitive present', 'to know'), ('imperative', 'imperative present', '2s', 'know'), ('imperative', 'imperative present', '1p', 'know'), ('imperative', 'imperative present', '2p', 'know')]

```

In [18]:

```

1 #code to loop over a list of tuples and access individual elements of each tuple in the
2 for index, tuple in enumerate(all_conjugated_forms[:2]):
3     #most tuples are of length 4 except 1 so we verify the conditions to not have index
4     if len(tuple)==4:
5         print(tuple[0])
6         print(tuple[1])
7         print(tuple[2])
8         print(tuple[3])
9         print("-----")

```

```

indicative
indicative present
1s
know
-----
indicative
indicative present
2s
know
-----

```

We want to conjugate the lemmatized words in a single function, so to do so we will need to return multiple values

A clean way to do it is using pandas' and apply function, but the difference , we will need to output more than one column.

Below is an example

In [92]:

```
1 #https://stackoverflow.com/questions/23586510/return-multiple-columns-from-pandas-apply
```

In [16]:

```

1  #define the language of the conjugator
2  default_conjugator = mlconjug3.Conjugator(language='en')
3  """
4  The function takes parameter x, in normal pandas'- apply operation where only one output is returned.
5  But in this case, we don't use the lambda function , so the input is a series
6  To get individual values instead of series ,we extract the values from the series i.e x.values
7  We then use conjugate our lemmatize verb and we iterate over the different possibilities
8  After that, we iterate over the list of tuples and access the elements of each tuple to get the conjugated verb
9
10 For example, a tuple has the format ('indicative', 'indicative present', '1s', 'know')
11
12 The length is 4:
13
14 2nd element(1st index) : verb form
15 3rd element (2nd index): mode:    1s = 1st person singular
16                                   2s = 2nd person singular
17                                   1p = 1st person plural
18                                   2p = 2nd person plural
19
20 4th element : conjugated verb
21
22 Note : we can return more types of verbs according to the needs of the project
23
24 """
25 def conjugate_lemma_verbs(x):
26     x=x[0]
27     test_verb = default_conjugator.conjugate(x)
28     all_conjugated_forms = test_verb.iterate()
29
30     for index, tuple in enumerate(all_conjugated_forms):
31         if len(tuple)==4:
32             if tuple[1] == "indicative present" and tuple[2] == "1s":
33                 conjugated_verb_1 = tuple[3]
34
35             if tuple[1] == "indicative present" and tuple[2] == "1p":
36                 conjugated_verb_11 = tuple[3]
37
38             if tuple[1] == "indicative past tense" and tuple[2] == "1s":
39                 conjugated_verb_2 = tuple[3]
40
41             if tuple[1] == "indicative past tense" and tuple[2] == "1p":
42                 conjugated_verb_22 = tuple[3]
43
44             if tuple[1] == "indicative present continuous" and tuple[2] == "1s":
45                 conjugated_verb_3 = tuple[3]
46
47             if tuple[1] == "indicative present perfect" and tuple[2] == "1s":
48                 conjugated_verb_4 = tuple[3]
49
50             if tuple[1] == "imperative present" and tuple[2] == "2s":
51                 conjugated_verb_5 = tuple[3]
52
53     return conjugated_verb_1,conjugated_verb_11,conjugated_verb_2,conjugated_verb_22,conjugated_verb_3,conjugated_verb_4,conjugated_verb_5

```

In [23]:

```
1 #apply function by creating 6 new columns
2 df_verb[["FP singular indicative present","FP plural indicative present","FP singular i
3         "FP plural indicative past","FP indicative present continuous",\
4         "FP indicative present perfect","SP imperative present"]]=\
5 df_verb[["Lemma"]].apply(conjugate_lemma_verbs,axis=1, result_type="expand")
```

In [19]:

```
1 df_verb.head()
```

Out[19]:

	Phrase	Verb	Tense	Lemma	FP singular indicative present	FP plural indicative present	FP singular indicative past	FP plural indicative past	ind p conti
0	The beauty of the landscape struck the travell...	struck	Past	strike	strike	strike	struck	struck	:
1	Nobody knows the truth about this affair.	knows	Pres	know	know	know	knew	knew	k
2	In a dictatorship, freedom of expression is li...	is limited	Past	be	am	are	was	were	
3	His wickedness had no limits.	had	Past	have	have	have	had	had	
4	His elegance impressed the assembly.	impressed	Past	impressed	impressee	impressee	impresseed	impresseed	impre