In [1]:	#!pip3 install numpy #!pip3 install pandas !(sys.executable) -m spacy download en_core_web_sm !(sys.executable) -m spacy download en_core_web_sm !(space re, numpy as np, pandas as pd !too pprint !eport pprint # Gensim #!pip install gensim #!pip install gensim #!pip install matplot #!pip install click #!python3 -m spacy download en #!python3 -m pip install gensim #!python3 -m pip install matplot #!python3 -m pip install nitk #!pip install python-Levenshtein !python3 -m pip install nltk.downloader stopwords !mpace gensim !eport spacy !mpace logging !mpace warnings
	<pre>import gensim.corpora as corpora # from gensim.utils import lemmatize, simple_preprocess from gensim.utils import simple_preprocess from gensim.models import CoherenceModel import matplotlib.pyplot as plt Collecting en-core-web-sm==3.1.0 Downloading https://github.com/explosion/spacy-models/releases/download/en_core_web_sm-3.1.0/en_core_web_sm-3.1.0-py3-none-any.whl (13.6 MB) Interpretation of the process of the pr</pre>
	Requirement already satisfied: jinja2 in /Users/elissaiosbaxevanis/opt/anaconda3/envs/py39-demo/lib/python3.9/s ite-packages (from spacy<3.2.0,>=3.1.0->en-core-web-sm==3.1.0) (3.0.1) Requirement already satisfied: setuptools in /Users/elissaiosbaxevanis/opt/anaconda3/envs/py39-demo/lib/python 3.9/site-packages (from spacy<3.2.0,>=3.1.0->en-core-web-sm==3.1.0) (49.6.0.post20210108) Requirement already satisfied: packaging>=20.0 in /Users/elissaiosbaxevanis/opt/anaconda3/envs/py39-demo/lib/python3.9/site-packages (from spacy<3.2.0,>=3.1.0->en-core-web-sm==3.1.0) (21.0) Requirement already satisfied: wasabi<1.1.0,>=0.8.1 in /Users/elissaiosbaxevanis/opt/anaconda3/envs/py39-demo/lib/python3.9/site-packages (from spacy<3.2.0,>=3.1.0->en-core-web-sm==3.1.0) (0.8.2) Requirement already satisfied: cymm<2.1.0,>=2.0.2 in /Users/elissaiosbaxevanis/opt/anaconda3/envs/py39-demo/lib/python3.9/site-packages (from spacy<3.2.0,>=3.1.0->en-core-web-sm==3.1.0) (2.0.5) Requirement already satisfied: catalogue<2.1.0,>=2.0.4 in /Users/elissaiosbaxevanis/opt/anaconda3/envs/py39-demo/lib/python3.9/site-packages (from spacy<3.2.0,>=3.1.0->en-core-web-sm==3.1.0) (2.0.4) Requirement already satisfied: requests<3.0.0,>=2.13.0 in /Users/elissaiosbaxevanis/opt/anaconda3/envs/py39-demo/lib/python3.9/site-packages (from spacy<3.2.0,>=3.1.0->en-core-web-sm==3.1.0) (2.25.1) Requirement already satisfied: srsly<3.0.0,>=2.4.1 in /Users/elissaiosbaxevanis/opt/anaconda3/envs/py39-demo/lib/python3.9/site-packages (from spacy<3.2.0,>=3.1.0->en-core-web-sm==3.1.0) (2.4.1) Requirement already satisfied: blis<0.8.0,>=0.4.0 in /Users/elissaiosbaxevanis/opt/anaconda3/envs/py39-demo/lib/python3.9/site-packages (from spacy<3.2.0,>=3.1.0->en-core-web-sm==3.1.0) (0.7.4) Requirement already satisfied: spacy-legacy<3.1.0,>=0.4.0 in /Users/elissaiosbaxevanis/opt/anaconda3/envs/py39-demo/lib/python3.9/site-packages (from spacy<3.2.0,>=3.1.0->en-core-web-sm==3.1.0) (0.7.4) Requirement already satisfied: spacy-legacy<3.2.0,>=3.1.0->en-core-web-sm==3.1.0) (0.7.4)
	b/python3.9/site-packages (from spacy<3.2.0,>=3.1.0->en-core-web-sm==3.1.0) (4.61.2) Requirement already satisfied: numpy>=1.15.0 in /Users/elissaiosbaxevanis/opt/anaconda3/envs/py39-demo/lib/pyth on3.9/site-packages (from spacy<3.2.0,>=3.1.0->en-core-web-sm==3.1.0) (1.21.0) Requirement already satisfied: preshed<3.1.0,>=3.0.2 in /Users/elissaiosbaxevanis/opt/anaconda3/envs/py39-demo/lib/python3.9/site-packages (from spacy<3.2.0,>=3.1.0->en-core-web-sm==3.1.0) (3.0.5) Requirement already satisfied: thinc<8.1.0,>=8.0.7 in /Users/elissaiosbaxevanis/opt/anaconda3/envs/py39-demo/lib/python3.9/site-packages (from spacy<3.2.0,>=3.1.0->en-core-web-sm==3.1.0) (8.0.7) Requirement already satisfied: pyparsing>=2.0.2 in /Users/elissaiosbaxevanis/opt/anaconda3/envs/py39-demo/lib/python3.9/site-packages (from packaging>=20.0->spacy<3.2.0,>=3.1.0->en-core-web-sm==3.1.0) (2.4.7) Requirement already satisfied: smart-open<6.0.0,>=5.0.0 in /Users/elissaiosbaxevanis/opt/anaconda3/envs/py39-demo/lib/python3.9/site-packages (from pathy>=0.3.5->spacy<3.2.0,>=3.1.0->en-core-web-sm==3.1.0) (5.1.0) Requirement already satisfied: typing-extensions>=3.7.4.3 in /Users/elissaiosbaxevanis/opt/anaconda3/envs/py39-demo/lib/python3.9/site-packages (from pydantic!=1.8,!=1.8.1,<1.9.0,>=1.7.4->spacy<3.2.0,>=3.1.0->en-core-web-sm==3.1.0) Requirement already satisfied: idna<3,>=2.5 in /Users/elissaiosbaxevanis/opt/anaconda3/envs/py39-demo/lib/python3.9/site-packages (from requests<3.0.0,>=2.13.0->spacy<3.2.0,>=3.1.0->en-core-web-sm==3.1.0) (2.10) Requirement already satisfied: urllib3<1.27,>=1.21.1 in /Users/elissaiosbaxevanis/opt/anaconda3/envs/py39-demo/lib/python3.9/site-packages (from requests<3.0.0,>=2.13.0->spacy<3.2.0,>=3.1.0->en-core-web-sm==3.1.0) (1.26.6) Requirement already satisfied: chardet<5,>=3.0.2 in /Users/elissaiosbaxevanis/opt/anaconda3/envs/py39-demo/lib/python3.9/site-packages (from requests<3.0.0,>=2.13.0->spacy<3.2.0,>=3.1.0->en-core-web-sm==3.1.0) (2.00) Requirement already satisfied: chardet<5,>=3.0.2 in /Users/elissaiosbax
In [2]:	Requirement already satisfied: click<7.2.0,>=7.1.1 in /Users/elissaiosbaxevanis/opt/anaconda3/envs/py39-demo/lib/python3.9/site-packages (from typer<0.4.0,>=0.3.0->spacy<3.2.0,>=3.1.0->en-core-web-sm==3.1.0) (7.1.2) Requirement already satisfied: MarkupSafe>=2.0 in /Users/elissaiosbaxevanis/opt/anaconda3/envs/py39-demo/lib/python3.9/site-packages (from jinja2->spacy<3.2.0,>=3.1.0->en-core-web-sm==3.1.0) (2.0.1) Download and installation successful You can now load the package via spacy.load('en_core_web_sm') ERROR: Could not find a version that satisfies the requirement nltk.downloader (from versions: none) ERROR: No matching distribution found for nltk.downloader (Users/elissaiosbaxevanis/opt/anaconda3/envs/py39-demo/lib/python3.9/site-packages/gensim/similarities/_initpy:15: UserWarning: The gensim.similarities.levenshtein submodule is disabled, because the optional Levenshte in package https://pypi.org/project/python-Levenshtein/ > is unavailable. Install Levenhstein (e.g. 'pip instal python-Levenshtein') to suppress this warning. warnings.warn(msg) #def fxn(): # warnings.warn("deprecated", DeprecationWarning) #with warnings.catch_warnings(): # warnings.simplefilter("ignore")
In [3]:	<pre>#fxn() warnings.filterwarnings("ignore", category=DeprecationWarning) Question 1 # NLTK Stop words import nltk nltk.download('stopwords') from nltk.corpus import stopwords stop_words = stopwords.words('english') stop_words.extend(['https','co','amp','from', 'subject', 're', 'edu', 'use', 'not', 'would', 'say', 'could', '_', 'be', 'know', 'good', 'go', 'get', 'do', 'done', 'try', 'many', 'some', 'nice', 'thank', 'think', 'see', 'rather', 'easy', 'easily', 'lot', 'lack', 'make', 'want', 'seem', 'run', 'need', 'even', 'right', 'line', 'even', 'also', 'may', 'take', 'come']) import nltk import ssl</pre>
Out[3]: In [4]:	
In [43]: In [44]: Out[44]:	<pre>df['word_count'].describe()</pre>
In [9]:	<pre>lof sent_to_words(sentences): for sent in sentences: sent = re.sub('\S*\S*\S*\S*', '', sent) # remove emails sent = re.sub('\S*\S*\S*', '', sent) # remove newline chars sent = re.sub("\'', "", sent) # remove single quotes sent = re.sub(r'http\S*", "", sent) # remove https sent = re.sub(r'samp','', sent) # remove amp sent = gensim.utils.simple_preprocess(str(sent), deacc= rm) yiold(sent) # Convert to list data = df.text.values.tolist() data_words = list(sent_to_words(data))</pre> bigram = gensim.models.Phrases(data_words, min_count= , threshold= 00) # higher threshold fewer
	<pre>trigram = gensim.models.Phrases:bigram[data_words], threshold=100) bigram_mod = gensim.models.phrases.Phraser(bigram) trigram_mod = gensim.models.phrases.Phraser(trigram) itrigram_mod = gensim.models.phrases.Phraser(trigram) trigram_mod = gensim.models.phrases.Phraser(trigram) itrigram_mod = gensim.models.phrases.Phraser(trigram) itrigram_mod = gensim.models.phrases.Phraser(trigram) trigram_mod = gensim.models.phrases.Phraser(trigram) texts = [word for word in simple_preprocess(str(doc)) if word not in stop_words] for doc in texts texts = [bigram_mod(doc] for doc in texts] texts = [trigram_mod(bigram_mod(doc)] for doc in texts] for sent in texts: doc = nlp = spacy.load("en', disable=['parser', 'ner']) nlp = spa</pre>
In [11]: In [12]:	<pre>id2word = corpora.Dictionary(data_ready) #create corpus corpus = [id2word.doc2bow(text) for text in data_ready] from gensim.matutils import kullback_leibler def kl_div_func(model): x = model.get_topics() score = 0 counter = 0 score_list = [] for i in x: if all(j!=i):</pre>
In [108	<pre>score += kullback_leibler(i,j)</pre>
	<pre>model_list : List of LDA topic models coherence_values : Coherence values corresponding to the LDA model with respective number of topics num coherence_values = [] model_list = [] perplexity_values = [] kl_div_values = [] kl_div_values = [] model = gensim.models.ldamulticore.LdaMulticore(corpus=corpus, num_topics=num_topics, id2word=id2word, random_state=100, workers=1, chunksize=1000, passes=10, alpha='symmetric',</pre>
In [14]:	<pre>iterations=50,</pre>
In [15]:	<pre>Num topics: 2 Num topics: 6 Num topics: 10 Num topics: 15 Num topics: 20 Num topics: 25 Num topics: 25 Num topics: 30 Num topics: 30 Num topics: 40 Num topics: 50 Num topics: 75 Num topics: 100 Num topics: 125 Num topics: 125 Num topics: 125 Num topics: 120 Num topics: 200</pre> <pre>plt.plot(x, coherence_values) plt.xlabel("Num Topics")</pre>
	plt.ylabel("Coherence_values"), loc='best') plt.legend(("Coherence_values"), loc='best') plt.show() plt.show() plt.plot(x, perplexity_values) plt.xlabel("Num Topics") plt.ylabel("Perplexity score") plt.legend(("perplexity_values"), loc='best') plt.title("Perplexity over number of topics") plt.show() Coherence over number of topics Coherence over number of topics
	1e16 Perplexity over number of topics Topics 1e16 Perplexity over number of topics Topics Num Topics Num Topics
In [16]:	<pre>results = pd.DataFrame (data = [x,</pre>
Out[16]:	0 2 -3.508417 3.632350e+02 2.022938 1 6 -4.322848 4.191805e+02 3.082971 2 10 -4.856069 5.962021e+02 4.096701 3 15 -5.780834 9.829778e+02 4.757345 4 20 -7.355326 1.361981e+03 5.344533
Out[16]:	R1_div_values).T results Num Topics = results Num Topics .** Num Topics = results Num Topics
	Column
	Page
	Comparison Com
In [207	Constitution Company Constitution Constitut
In [207	
In [207	### 15 Comment ## 12 Comment ## 12 Comment ### 15 C
In [207 In [208	The Company of Compa
In [207 In [208	March Colored Colore
In [209	Column
In [209 In [209	March Marc
In [209 In [209	March Marc
In [209 In [209	March Marc
In [209 In [211 In [211	March Marc
In [209 In [211 In [211	March Marc
In [209 In [211 In [211	The state of the content of the cont

