

State of the Union (1947-2022): Sentiment Analysis and Topic Modeling

By: Shirsho Dasgupta (2022)

General Notes:

Some visualizations may not be compatible with viewing directly on Github. View it by copying the URL [here](#) or run the code after downloading this repo.

Notes on Sentiment Analysis:

The code reads the text of every State of the Union speech (delivered in-person on the Hill) from 1947 to 2022 and performs sentiment analysis on them using the NRC Word-Emotion Association Lexicon (EmoLex).

The NRC Emotion Lexicon is a list of English words and their associations with eight basic emotions (anger, fear, anticipation, trust, surprise, sadness, joy, and disgust) and two sentiments (negative and positive). The annotations were manually done by crowdsourcing.

The code matches the words in the speeches to that of the dictionary then adds up the factor of the emotion.

Since EmoLex is, at the end of the day, crowdsourced and finite, there might be some words which are outside its scope or some nuances which are not accounted for (the code performs word-to-word comparison). The analysis is always at least an approximation.

Notes on Topic Modelings

The code deploys machine-learning modules on the text of every State of the Union speech (delivered in-person on the Hill) from 1947 to 2022. The algorithm calculates the frequency of each word in the speech. It then uses a neural network to predict the date of the speech.

The code uses the Scikitlearn module to compute Non-Negative Matrix Factorization (NMF)/Latent Semantic Indexing (LSI) and Latent Dirichlet Allocation (LDA) methods. It also deploys

Stat. 111, M111, A111, 121, 131, M111, M121, 121M111, Stat 111

第1章

Importing libraries

```
In [1]: 1 import numpy as np  
2 import re  
3 import pandas as pd  
4 %matplotlib inline
```

Creating dataframe of speeches

Importing SOTU speeches

```
In [2]: 1 import glob  
2  
3 filenames = glob.glob("notebooks/*.ipynb")
```

```
In [3]: 1 speeches = [open(filename).read() for filename in filenames]
2 len(speeches)
```

Out[3]: 75

Storing SOTU speeches in one dataframe

```
In [4]: 1 ### stored entire text of speech under one column and the associated filename under another
2 speeches_df = pd.DataFrame({"text": speeches, "filename": filenames})
3 speeches_df.head(3)
```

Out[4]:	text	filename
0	Mr. President, Mr. Speaker, Members of the Con...	sotu/01-06-1955.txt
1	Mr. President, Mr. Speaker, Members of the Eig...	sotu/01-07-1955.txt
2	Mr. Speaker, Mr. President, my colleagues in t...	sotu/01-20-1972.txt

```
In [5]: 1  ### creates two new columns to store year and name of president
2  speeches_df["year"] = " "
3  speeches_df["pres"] = " "
4  speeches_df["id"] = " "
5
6  ### loop runs through dataframe
7  for i in range(0, len(speeches_df)):
8
9    ### extracts (from filename) and stores year of speech
10   speeches_df["year"][i] = int(speeches_df["filename"][i][-8:-4])
11   x = speeches_df["year"][i]
12
13  ### condition checks what year the speech was delivered and stores the name of the president accordingly
14  if (x >= 1947) & (x <= 1952):
15    speeches_df["pres"][i] = "truman"
16  elif (x >= 1953) & (x <= 1960):
17    speeches_df["pres"][i] = "eisenhower"
18  elif (x >= 1961) & (x <= 1963):
19    speeches_df["pres"][i] = "kennedy"
20  elif (x >= 1964) & (x <= 1969):
21    speeches_df["pres"][i] = "johnson"
22  elif (x >= 1970) & (x <= 1974):
23    speeches_df["pres"][i] = "nixon"
24  elif (x >= 1975) & (x <= 1977):
25    speeches_df["pres"][i] = "ford"
26  elif (x >= 1978) & (x <= 1980):
27    speeches_df["pres"][i] = "carter"
28  elif (x >= 1981) & (x <= 1988):
29    speeches_df["pres"][i] = "reagan"
30  elif (x >= 1989) & (x <= 1992):
31    speeches_df["pres"][i] = "bush sr."
32  elif (x >= 1993) & (x <= 2000):
33    speeches_df["pres"][i] = "clinton"
```

```

24     if (x == 2001) & (x == 2002):
25         speeches_df["pres"][i] = "bush jr."
26     elif (x >= 2009) & (x <= 2016):
27         speeches_df["pres"][i] = "obama"
28     elif (x >= 2017) & (x <= 2020):
29         speeches_df["pres"][i] = "trump"
30     elif (x > 2020):
31         speeches_df["pres"][i] = "biden"
32
33 speeches_df["id"][i] = speeches_df["pres"][i] + "-" + str(speeches_df["year"][i])

```

In [6]:

```

1 ##### displays final master dataframe
2 speeches_df

```

Out[6]:

	text	filename	year	pres	id
0	Mr. President, Mr. Speaker, Members of the Con...	sotu/01-06-1955.txt	1955	eisenhower	eisenhower-1955
1	Mr. President, Mr. Speaker, Members of the Eig...	sotu/01-07-1954.txt	1954	eisenhower	eisenhower-1954
2	Mr. Speaker, Mr. President, my colleagues in t...	sotu/01-20-1972.txt	1972	nixon	nixon-1972
3	Mr. President, Mr. Speaker, Members of the Hou...	sotu/02-17-1993.txt	1993	clinton	clinton-1993
4	Madam Speaker, Madam Vice President—no Preside...	sotu/04-28-2021.txt	2021	biden	biden-2021
...
70	Mr. Speaker, Mr. Vice President, Members of Co...	sotu/01-30-2018.txt	2018	trump	trump-2018
71	Mr. Speaker, Vice President Cheney, Members of...	sotu/01-28-2003.txt	2003	bush jr.	bush jr-2003
72	Thank you very much. Mr. Speaker, Vice Preside...	sotu/01-29-2002.txt	2002	bush jr.	bush jr-2002
73	Mr. President, Mr. Speaker, Members of the Con...	sotu/01-08-1951.txt	1951	truman	truman-1951
74	Mr. Vice President, my old colleague from Mass...	sotu/01-11-1962.txt	1962	kennedy	kennedy-1962

75 rows x 5 columns

Creating separate dataframes for each president

In [7]:

```

1 ##### filters by name of president and stores in separate dataframes
2 speeches_df_obama = speeches_df[speeches_df["pres"] == "obama"]
3 speeches_df_trump = speeches_df[speeches_df["pres"] == "trump"]
4 speeches_df_biden = speeches_df[speeches_df["pres"] == "biden"]

```

Sentiment analysis

Importing NRC Emotion Lexicon

Download compressed file from <http://saifmohammad.com/WebPages/NRC-Emotion-Lexicon.htm>

Move the relevant text-file to the directory

In [8]:

```

1 filepath = "NRC-Emotion-Lexicon-Wordlevel-v0.92.txt"
2 emolex_df = pd.read_csv(filepath, names = ["word", "emotion", "association"], sep = "\t", keep_default_na = False)
3 emolex_df = emolex_df.pivot(index = "word", columns = "emotion", values = "association").reset_index()
4 emolex_df.head()

```

Out[8]:

emotion	word	anger	anticipation	disgust	fear	joy	negative	positive	sadness	surprise	trust
0	aback	0	0	0	0	0	0	0	0	0	0
1	abacus	0	0	0	0	0	0	0	0	0	1
2	abandon	0	0	0	1	0	1	0	1	0	0
3	abandoned	1	0	0	1	0	1	0	1	0	0
4	abandonment	1	0	0	1	0	1	0	1	1	0

Calculating share of words using TfidfVectorizer

In [9]:

```

1 from sklearn.feature_extraction.text import TfidfVectorizer
2
3 vec = TfidfVectorizer(vocabulary = emolex_df.word, use_idf = False, norm = "l1")
4 matrix = vec.fit_transform(speeches_df.text)
5 vocab = vec.get_feature_names()
6 wordcount_df = pd.DataFrame(matrix.toarray(), columns = vocab)
7 wordcount_df.head()

```

/Users/shirshdasgupta/Documents/Python/anaconda3/lib/python3.9/site-packages/sklearn/utils/deprecation.py:87: FutureWarning: Function get_feature_names is deprecated; get_feature_names is deprecated in 1.0 and will be removed in 1.2. Please use get_feature_names_out instead.
warnings.warn(msg, category=FutureWarning)

Out[9]:

	aback	abacus	abandon	abandoned	abandonment	abate	abatement	abba	abbot	abbreviate	...	zephyr	zeppelin	zest	zip	zodiac	zone	zoo	zoological	zoology	zoom
0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	...	0.0	0.0	0.0	0.0	0.00000	0.0	0.0	0.0	0.0	0.0
1	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	...	0.0	0.0	0.0	0.0	0.00000	0.0	0.0	0.0	0.0	0.0
2	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	...	0.0	0.0	0.0	0.0	0.00000	0.0	0.0	0.0	0.0	0.0
3	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	...	0.0	0.0	0.0	0.0	0.00000	0.0	0.0	0.0	0.0	0.0
4	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	...	0.0	0.0	0.0	0.0	0.00105	0.0	0.0	0.0	0.0	0.0

5 rows x 14154 columns

Applying EmoLex

The words which appear per category in the NRC Emotion Lexicon dictionary are identified and their shares are added up. The sum gives the total share of the document that relates to a particular emotion or sentiment. For instance say a document has two "angry" words each making up 0.01 of the text, then the anger-factor as it were of the text is 0.02 (= 0.01 + 0.01).

In [10]:

```

1 ##### identifies all negative words
2 neg_words = emolex_df[emolex_df.negative == 1]["word"]
3 ##### adds up shares of negative words
4 speeches_df["negative"] = wordcount_df[neg_words].sum(axis = 1)
5
6 ##### the above process is repeated per sentiment
7
8 pos_words = emolex_df[emolex_df.positive == 1]["word"]
9 speeches_df["positive"] = wordcount_df[pos_words].sum(axis = 1)
10
11 angry_words = emolex_df[emolex_df.anger == 1]["word"]
12 speeches_df["anger"] = wordcount_df[angry_words].sum(axis = 1)
13
14 anticip_words = emolex_df[emolex_df.anticipation == 1]["word"]
15 speeches_df["anticipation"] = wordcount_df[anticip_words].sum(axis = 1)

```

```

16
17 disgust_words = emolex_df[emolex_df.disgust == 1]["word"]
18 speeches_df["disgust"] = wordcount_df[disgust_words].sum(axis = 1)
19
20 fear_words = emolex_df[emolex_df.fear == 1]["word"]
21 speeches_df["fear"] = wordcount_df[fear_words].sum(axis = 1)
22
23 joy_words = emolex_df[emolex_df.joy == 1]["word"]
24 speeches_df["joy"] = wordcount_df[joy_words].sum(axis = 1)
25
26 sad_words = emolex_df[emolex_df.sadness == 1]["word"]
27 speeches_df["sadness"] = wordcount_df[sad_words].sum(axis = 1)
28
29 surprise_words = emolex_df[emolex_df.surprise == 1]["word"]
30 speeches_df["surprise"] = wordcount_df[surprise_words].sum(axis = 1)
31
32 trust_words = emolex_df[emolex_df.trust == 1]["word"]
33 speeches_df["trust"] = wordcount_df[trust_words].sum(axis = 1)

```

In [11]:
1 ### displays final dataframe
2 speeches_df.head()

Out[11]:

	text	filename	year	pres	id	negative	positive	anger	anticipation	disgust	fear	joy	sadness	surprise	trust
0	Mr. President, Mr. Speaker, Members of the Con...	sotu/01-06-1955.txt	1955	eisenhower	eisenhower-1955	0.109145	0.328276	0.045512	0.128108	0.028234	0.093974	0.096924	0.040455	0.030763	0.226717
1	Mr. President, Mr. Speaker, Members of the Eig...	sotu/01-07-1954.txt	1954	eisenhower	eisenhower-1954	0.114540	0.304916	0.044979	0.116632	0.022490	0.091004	0.080544	0.041318	0.031381	0.211297
2	Mr. Speaker, Mr. President, my colleagues in t...	sotu/01-20-1972.txt	1972	nixon	nixon-1972	0.112086	0.292398	0.049708	0.115010	0.039961	0.082846	0.092593	0.036062	0.034113	0.195906
3	Mr. President, Mr. Speaker, Members of the Hou...	sotu/02-17-1993.txt	1993	clinton	clinton-1993	0.126192	0.243971	0.042625	0.123948	0.021873	0.058329	0.076276	0.053842	0.035895	0.162086
4	Madam Speaker, Madam Vice President—no Preside...	sotu/04-28-2021.txt	2021	biden	biden-2021	0.132283	0.266142	0.058268	0.114961	0.034646	0.087664	0.080315	0.061942	0.039370	0.167454

Exporting the dataset

In [12]:
1 speeches_df.to_csv("sotu_analysis.csv", index = False)

Re-importing created dataset

In [13]:
1 speeches = pd.read_csv("sotu_analysis.csv")

In [14]:
1 speeches.shape

Out[14]: (75, 15)

Plotting sentiments

Importing Altair libraries

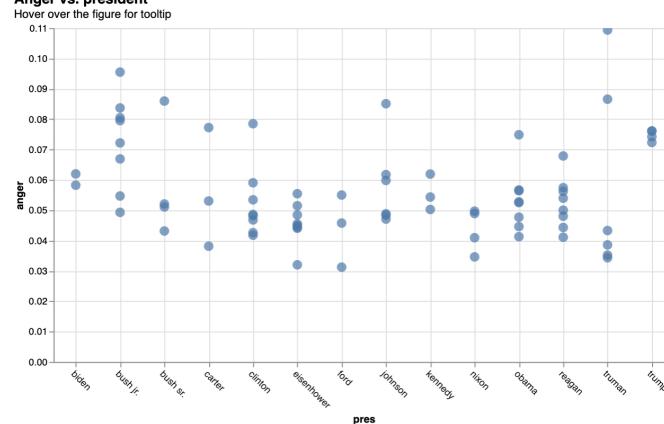
In [15]:
1 import altair as alt
2
3 alt.renderers.enable("default")

Out[15]: RenderRegistry.enable('default')

Plotting anger-factor by president

In [16]:
1 alt.Chart(speeches).mark_circle(size = 100).encode(
2 x = "pres",
3 y = "anger",
4 tooltip = ["pres", "year", "negative"]
5).configure_axis(
6 grid = True
7).properties(
8 width = 650,
9 height = 350,title = {
10 "text": "Anger vs. president",
11 "subtitle": "Hover over the figure for tooltip",
12 }
13).configure_title(
14 fontSize = 15,
15 anchor = "start",
16).configure_axisX(
17 labelAngle = 45
18).interactive()

Out[16]: Anger vs. president



Plotting negativity-factor by year

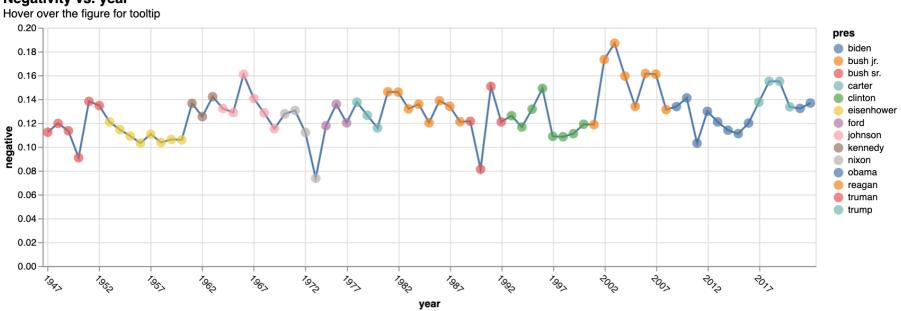
In [17]:
1 ### creates chart with point-markings
2 points = alt.Chart(speeches).mark_circle(size = 100).encode(
3 alt.X("year:", axis = alt.Axis(values = list(range(1947, 2022, 5)))),
4 alt.Y("negative"),
5 alt.Color("pres"),
6 tooltip = ["pres", "year", "negative"] ### adds tooltip on mouse-hover
7)

```

8     ## creates line chart
9     line = alt.Chart(speeches).mark_line().encode(
10        alt.X("year:O", axis = alt.Axis(values = list(range(1947, 2022, 5)))),
11        alt.Y("negative"),
12    )
13 )
14
15     ## layers both charts and adds some other properties
16     (line + points).configure_axis(
17        grid = True
18    ).configure_axisX(
19        labelAngle = 45
20    ).properties(
21        width = 810,
22        height = 250,
23        title = {
24            "text": "Negativity vs. year",
25            "subtitle": "Hover over the figure for tooltip",
26        }
27    ).configure_title(
28        fontSize = 15,
29        anchor = "start",
30    ).interactive()

```

Out[17]: Negativity vs. year



Topic modeling

Setting dataframe

```

In [18]: 1 sotu = speeches[["text", "filename", "year", "pres", "id"]].copy()
2 sotu.head()

```

Out[18]:

		text	filename	year	pres	id
0	Mr. President, Mr. Speaker, Members of the Con...	sotu/01-06-1955.txt	1955	eisenhower	eisenhower-1955	
1	Mr. President, Mr. Speaker, Members of the Ele...	sotu/01-07-1954.txt	1954	eisenhower	eisenhower-1954	
2	Mr. Speaker, Mr. President, my colleagues in ...	sotu/01-20-1972.txt	1972	nixon	nixon-1972	
3	Mr. President, Mr. Speaker, Members of the Hou...	sotu/02-17-1993.txt	1993	clinton	clinton-1993	
4	Madam Speaker, Madam Vice President—no Preside...	sotu/04-28-2021.txt	2021	biden	biden-2021	

Non-Negative Matrix Factorization (NMF)/Latent Semantic Indexing (LSI) topic modeling

Vectorizing all words

```

In [19]: 1 from sklearn.feature_extraction.text import TfidfVectorizer
2 from nltk.stem import SnowballStemmer
3
4 stemmer = SnowballStemmer("english")
5
6 class StemmedTfidfVectorizer(TfidfVectorizer):
7     def build_analyzer(self):
8         analyzer = super(StemmedTfidfVectorizer, self).build_analyzer()
9         return lambda doc:(stemmer.stem(word) for word in analyzer(doc))
10
11 tfidf_vectorizer = StemmedTfidfVectorizer(stop_words = "english", min_df = 5, max_df = 0.40)
12 x = tfidf_vectorizer.fit_transform(sotu.text.astype(str))

```

Deploying NMF module

```

In [20]: 1 %%time
2
3     ## imports modules
4     from sklearn.decomposition import NMF
5
6     ## selects number of topics
7     model = NMF(n_components = 10)
8     model.fit(x)
9
10    ## selects number of words per model
11    n_words = 7
12    ## stores each word
13    feature_names = tfidf_vectorizer.get_feature_names()
14
15    ## creates new arrays to store list of topics and topic numbers
16    topic_list_nmf = []
17    topic_list_index_nmf = []
18
19    ## applies NMF modeling
20    for topic_idx, topic in enumerate(model.components_):
21        top_n = [feature_names[i]
22                 for i in topic.argsort()
23                 [-n_words:][::-1]]
24        top_features = ' '.join(top_n)
25        topic_list_nmf.append(f'{topic_idx}_{top_features}')
26        topic_list_index_nmf.append("Topic " + str(topic_idx + 1))
27
28    ## displays topics
29    print("Topic " + str(topic_idx + 1) + ": " + top_features)

```

Topic 1: kid class manufactur oil doesn innov gas
Topic 2: agricultur expenditur adequ atom adjust 1946 communist
Topic 3: iraq iraqi terrorist al qaida saddam regim
Topic 4: ryan isi pillar gang hero incred veteran

```

Topic 5: nicaragua revolut excel veto gramm rudman governor
Topic 6: oil 1974 1980 barrel 1976 capabl branch
Topic 7: vietnam south expeditur communist abund desir wish
Topic 8: 21st ought gun got read class grade
Topic 9: vaccin pandem gun folk 19 childcar cancer
Topic 10: communist session atlant disarma nam viet lack
CPU times: user 480 ms, sys: 20.8 ms, total: 500 ms
Wall time: 147 ms

/Users/shirshodasgupta/Documents/Python/anaconda3/lib/python3.9/site-packages/sklearn/decomposition/_nmf.py:289: FutureWarning: The 'init=None' and 'n_components' is less than 'n_samples' and 'n_features', will be changed from 'nndsvd' to 'nndsvda' in 1.1 (renaming of 0.26).
    warnings.warn(
/Users/shirshodasgupta/Documents/Python/anaconda3/lib/python3.9/site-packages/sklearn/utils/deprecation.py:87: FutureWarning: Function get_feature_names is deprecated; get_feature_names is deprecated in 1.0 and will be removed in 1.2. Please use get_feature_names_out instead.
    warnings.warn(msg, category=FutureWarning)

```

Converting computed factors into dataset

```

In [21]: 1  ### converts counts into numbers
2  amounts = model.transform(x) * 100
3
4  ### sets up dataframe with corresponding topic-numbers and amount numbers
5  topics_nmf = pd.DataFrame(amounts, columns = topic_list_index_nmf)
6
7  ### sets up dataframe with corresponding topics and subjects
8  topics_list_nmf = pd.DataFrame(amounts, columns = topic_list_nmf)
9
10 topics_nmf

```

```

Out[21]:
   Topic 1   Topic 2   Topic 3   Topic 4   Topic 5   Topic 6   Topic 7   Topic 8   Topic 9   Topic 10
0  0.613844  54.268987  0.000000  0.000000  0.200546  0.000000  1.904521  0.173698  0.000000  1.615704
1  0.330403  51.035306  0.000000  0.000000  0.000000  3.638860  2.363694  0.000000  0.000000  0.189744
2  0.000000  0.000000  0.000000  1.137020  6.012578  12.888952  7.798940  0.000000  0.000000  25.867264
3  8.488645  1.810723  0.000000  0.000000  0.000000  7.050203  2.624654  24.156173  0.000000  2.904309
4  0.000000  0.000000  0.000000  0.000000  0.000000  1.639435  0.063283  0.000000  90.373168  0.000000
...
...
...
...
...
...
70 0.000000  0.000000  0.000000  77.829350  0.000000  0.000000  0.000000  0.000000  0.000000  0.000000
71 0.000000  1.374149  52.437751  0.000000  0.000000  0.000000  0.000000  2.176445  0.000000  0.000000
72 0.000000  0.000000  41.596349  9.516587  1.327411  0.274926  0.000000  0.000000  6.703602  0.000000
73 0.000000  8.065589  0.000000  1.125444  0.020000  0.000000  0.000000  0.000000  5.305285  39.098590
74 0.000000  4.553163  0.000000  0.000000  1.240334  2.385833  0.000000  0.000000  2.058022  63.850685

```

75 rows × 10 columns

Storing topic list as a dataframe

```

In [22]: 1  topic_indices = list(topics_nmf)
2  topic_subjects = list(topics_list_nmf)
3  topic_desc = pd.DataFrame(topic_indices, topic_subjects).reset_index()
4  topic_desc_nmf = topic_desc.rename(columns = {"index": "topic_desc", 0: "topic"})
5  topic_desc_nmf

```

```

Out[22]:
   topic_desc      topic
0  kid_class_manufactur  Topic 1
1  agricultur_expenditur_adequ  Topic 2
2  iraq_iraqi_terrorist  Topic 3
3  ryan_isi_pilliar  Topic 4
4  nicaragua_revolut_excel  Topic 5
5  oil_1974_1980  Topic 6
6  vietnam_south_expeditur  Topic 7
7  21st_ought_gun  Topic 8
8  vaccin_pandem_gun  Topic 9
9  communist_session_atlant  Topic 10

```

Merging speeches database with topics dataframe

```

In [23]: 1  nmf_merged = topics_nmf.join(sotu)
2  nmf_merged.head(3)

```

```

Out[23]:
   Topic 1   Topic 2   Topic 3   Topic 4   Topic 5   Topic 6   Topic 7   Topic 8   Topic 9   Topic 10          text      filename  year     pres      id
0  0.613844  54.268987  0.0  0.00000  0.200546  0.000000  1.904521  0.173698  0.0  1.615704  Mr. President, Mr. Speaker, Members of the Con...  sotu/01-06-1955.txt  1955  eisenhower  eisenhower-1955
1  0.330403  51.035306  0.0  0.00000  0.000000  3.638860  2.363694  0.000000  0.0  0.189744  Mr. President, Mr. Speaker, Members of the Eig...  sotu/01-07-1954.txt  1954  eisenhower  eisenhower-1954
2  0.000000  0.000000  0.0  1.13702  6.012578  12.888952  7.798940  0.000000  0.0  25.867264  Mr. Speaker, Mr. President, my colleagues in t...  sotu/01-20-1972.txt  1972  nixon  nixon-1972

```

Rearranging dataframe for plotting

```

In [24]: 1  nmf_1 = nmf_merged[["text", "filename", "year", "pres", "id", "Topic 1"]].copy()
2  nmf_1["topic"] = "Topic 1"
3  nmf_1 = nmf_1.rename(columns = {"Topic 1": "factor"})
4
5  nmf_2 = nmf_merged[["text", "filename", "year", "pres", "id", "Topic 2"]].copy()
6  nmf_2["topic"] = "Topic 2"
7  nmf_2 = nmf_2.rename(columns = {"Topic 2": "factor"})
8
9  nmf_3 = nmf_merged[["text", "filename", "year", "pres", "id", "Topic 3"]].copy()
10 nmf_3["topic"] = "Topic 3"
11 nmf_3 = nmf_3.rename(columns = {"Topic 3": "factor"})
12
13 nmf_4 = nmf_merged[["text", "filename", "year", "pres", "id", "Topic 4"]].copy()
14 nmf_4["topic"] = "Topic 4"
15 nmf_4 = nmf_4.rename(columns = {"Topic 4": "factor"})
16
17 nmf_5 = nmf_merged[["text", "filename", "year", "pres", "id", "Topic 5"]].copy()
18 nmf_5["topic"] = "Topic 5"
19 nmf_5 = nmf_5.rename(columns = {"Topic 5": "factor"})
20
21 nmf_6 = nmf_merged[["text", "filename", "year", "pres", "id", "Topic 6"]].copy()
22 nmf_6["topic"] = "Topic 6"
23 nmf_6 = nmf_6.rename(columns = {"Topic 6": "factor"})
24

```

```

25 nmf_7 = nmf_merged[["text", "filename", "year", "pres", "id", "Topic 7"]].copy()
26 nmf_7["topic"] = "Topic 7"
27 nmf_7 = nmf_7.rename(columns = {"Topic 7": "factor"})
28
29 nmf_8 = nmf_merged[["text", "filename", "year", "pres", "id", "Topic 8"]].copy()
30 nmf_8["topic"] = "Topic 8"
31 nmf_8 = nmf_8.rename(columns = {"Topic 8": "factor"})
32
33 nmf_9 = nmf_merged[["text", "filename", "year", "pres", "id", "Topic 9"]].copy()
34 nmf_9["topic"] = "Topic 9"
35 nmf_9 = nmf_9.rename(columns = {"Topic 9": "factor"})
36
37 nmf_10 = nmf_merged[["text", "filename", "year", "pres", "id", "Topic 10"]].copy()
38 nmf_10["topic"] = "Topic 10"
39 nmf_10 = nmf_10.rename(columns = {"Topic 10": "factor"})
40
41 nmf_master = pd.concat([nmf_1, nmf_2, nmf_3, nmf_4, nmf_5, nmf_6, nmf_7, nmf_8, nmf_9, nmf_10], ignore_index = True)
42 nmf_master = pd.merge(nmf_master, topic_desc_nmf, on = "topic", how = "left")
43
44 nmf_master

```

Out[24]:

	text	filename	year	pres	id	factor	topic	topic_desc
0	Mr. President, Mr. Speaker, Members of the Con...	sotu/01-06-1955.txt	1955	eisenhower	eisenhower-1955	0.613844	Topic 1	kid_class_manufactur
1	Mr. President, Mr. Speaker, Members of the Eig...	sotu/01-07-1954.txt	1954	eisenhower	eisenhower-1954	0.330403	Topic 1	kid_class_manufactur
2	Mr. Speaker, Mr. President, my colleagues in t...	sotu/01-20-1972.txt	1972	nixon	nixon-1972	0.000000	Topic 1	kid_class_manufactur
3	Mr. President, Mr. Speaker, Members of the Hou...	sotu/02-17-1993.txt	1993	clinton	clinton-1993	8.488645	Topic 1	kid_class_manufactur
4	Madam Speaker, Madam Vice President—no Preside...	sotu/04-28-2021.txt	2021	biden	biden-2021	0.000000	Topic 1	kid_class_manufactur
...
745	Mr. Speaker, Mr. Vice President, Members of Co...	sotu/01-30-2018.txt	2018	trump	trump-2018	0.000000	Topic 10	communist_session_atlant
746	Mr. Speaker, Vice President Cheney, Members of...	sotu/01-28-2003.txt	2003	bush jr.	bush jr-2003	0.000000	Topic 10	communist_session_atlant
747	Thank you very much. Mr. Speaker, Vice Preside...	sotu/01-29-2002.txt	2002	bush jr.	bush jr-2002	0.000000	Topic 10	communist_session_atlant
748	Mr. President, Mr. Speaker, Members of the Con...	sotu/01-08-1951.txt	1951	truman	truman-1951	39.098590	Topic 10	communist_session_atlant
749	Mr. Vice President, my old colleague from Mass...	sotu/01-11-1962.txt	1962	kennedy	kennedy-1962	63.850685	Topic 10	communist_session_atlant

750 rows × 8 columns

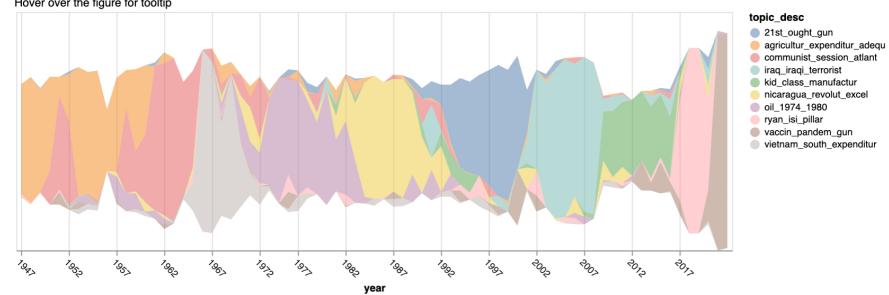
Plotting steamgraph

```

In [25]: 1 alt.Chart(nmf_master).mark_area(opacity = 0.5).encode(
2     alt.X("year:0", axis = alt.Axis(values = list(range(1947, 2022, 5))),),
3     alt.Y("sum(factor):0", stack = "center", axis = None),
4     alt.Color("topic_desc", scale = alt.Scale(scheme = "tableau10")),
5     tooltip = ["pres", "year", "topic_desc"]
6 ).configure_axis(
7     grid = True
8 ).configure_axisX(
9     labelAngle = 45
10 ).properties(
11     width = 750,
12     height = 250,
13     title = {
14         "text": "Topics by year",
15         "subtitle": "Hover over the figure for tooltip",
16     }
17 ).configure_title(
18     fontSize = 15,
19     anchor = "start",
20 ).interactive()

```

Out[25]: Topics by year



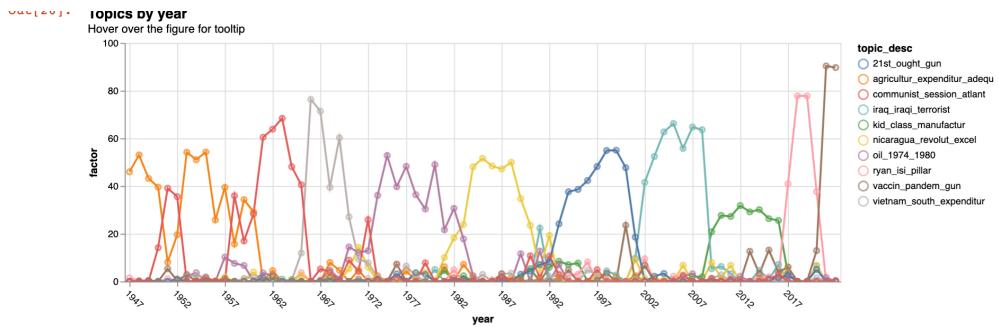
Plotting line chart

```

In [26]: 1 lines = alt.Chart(nmf_master).mark_line().encode(
2     alt.Color("topic_desc", scale = alt.Scale(scheme = "tableau10")),
3     alt.X("year:0", axis = alt.Axis(values = list(range(1947, 2022, 5))),),
4     y = "factor",
5 )
6
7 points = alt.Chart(nmf_master).mark_point().encode(
8     alt.Color("topic_desc", scale = alt.Scale(scheme = "tableau10")),
9     alt.X("year:0", axis = alt.Axis(values = list(range(1947, 2022, 5))),),
10    y = "factor",
11    tooltip = ["pres", "year", "topic_desc"]
12 )
13
14 chart = lines + points
15
16 chart.configure_axis(
17     grid = True
18 ).configure_axisX(
19     labelAngle = 45
20 ).properties(
21     width = 750,
22     height = 250,
23     title = {
24         "text": "Topics by year",
25         "subtitle": "Hover over the figure for tooltip",
26     }
27 ).configure_title(
28     fontSize = 15,
29     anchor = "start",
30 ).interactive()

```

Out[26]:



Note: This same figure can also be drawn using selected topics.

Latent Dirichlet Allocation (LDA) topic modeling

Displaying dataframe

```
In [27]: 1 sotu
```

```
Out[27]:
```

		text	filename	year	pres	id
0	Mr. President, Mr. Speaker, Members of the Con...	sotu/01-06-1955.txt	1955	eisenhower	eisenhower-1955	
1	Mr. President, Mr. Speaker, Members of the Eig...	sotu/01-07-1954.txt	1954	eisenhower	eisenhower-1954	
2	Mr. Speaker, Mr. President, my colleagues in t...	sotu/01-20-1972.txt	1972	nixon	nixon-1972	
3	Mr. President, Mr. Speaker, Members of the Hou...	sotu/02-17-1993.txt	1993	clinton	clinton-1993	
4	Madam Speaker, Madam Vice President—no Preside...	sotu/04-28-2021.txt	2021	biden	biden-2021	
...
70	Mr. Speaker, Mr. Vice President, Members of Co...	sotu/01-30-2018.txt	2018	trump	trump-2018	
71	Mr. Speaker, Vice President Cheney, Members of...	sotu/01-28-2003.txt	2003	bush jr.	bush jr.-2003	
72	Thank you very much. Mr. Speaker, Vice Preside...	sotu/01-29-2002.txt	2002	bush jr.	bush jr.-2002	
73	Mr. President, Mr. Speaker, Members of the Con...	sotu/01-08-1951.txt	1951	truman	truman-1951	
74	Mr. Vice President, my old colleague from Mass...	sotu/01-11-1962.txt	1962	kennedy	kennedy-1962	

75 rows × 5 columns

Vectorizing all words

```
In [28]: 1 from sklearn.feature_extraction.text import CountVectorizer
2 from nltk.stem import SnowballStemmer
3
4 stemmer = SnowballStemmer("english")
5
6 analyzer = CountVectorizer().build_analyzer()
7
8 class StemmedCountVectorizer(CountVectorizer):
9     def build_analyzer(self):
10         analyzer = super(CountVectorizer, self).build_analyzer()
11         return lambda doc: (stemmer.stem(word) for word in analyzer(doc))
12
13 vectorizer = StemmedCountVectorizer(stop_words = "english", min_df = 5, max_df = 0.4)
14 x = vectorizer.fit_transform(sotu.text)
```

Note: LDA modeling has in-built TF-IDF unlike NMF, so only a CountVectorizer is required.

Deploying LDA module

```
In [29]: 1 %%time
2
3 ### imports modules
4 from sklearn.decomposition import LatentDirichletAllocation
5
6 ### assigns number of topics
7 n_topics = 10
8 model = LatentDirichletAllocation(n_components = n_topics)
9 model.fit(x)
10
11 ### assigns number of words per topic
12 n_words = 7
13 feature_names = vectorizer.get_feature_names()
14
15 ### applies LDA modeling
16 topic_list_lda = []
17 topic_list_index_lda = []
18 for topic_idx, topic in enumerate(model.components_):
19     top_n = [feature_names[i]
20             for i in topic.argsort()
21             [-n_words:][::-1]]
22     top_features = " ".join(top_n)
23     topic_list_lda.append(f"\t{topic_idx}\t{top_features}")
24     topic_list_index_lda.append("Topic " + str(topic_idx + 1))
25
26 print(f"Topic {topic_idx + 1}: {top_features}")

Topic 1: ought class kid governor revolut got tough
Topic 2: terrorist iraq iraqi afghanistan al qaida regim
Topic 3: vietnam oil communist latin expenditur south needi
Topic 4: hussein saddam inspector intellig chemic biolog agent
Topic 5: vietnam south voic door knowledg desir planet
Topic 6: session surplus read invit john senior mayor
Topic 7: communist agricultur expenditur adequ republ branch affair
Topic 8: veteran disput bargain 1946 collect atom jurisdict
Topic 9: vietnam januari abund battl travel 1969 seventi
Topic 10: 21st veteran class terrorist manufactur gun kid
CPU times: user 3.96 s, sys: 45.8 ms, total: 4.01 s
Wall time: 1.01 s
```

```
/Users/shirshodasgupta/Documents/Python/anaconda3/lib/python3.9/site-packages/sklearn/utils/deprecation.py:87: FutureWarning: Function get_feature_names is deprecated; get_feature_names is deprecated in 1.0 and will be removed in 1.2. Please use get_feature_names_out instead.
  warnings.warn(msg, category=FutureWarning)
```

Converting computed factors into dataset

```
In [30]: 1 ##### converts counts into numbers
2 amounts = model.transform(x) * 100
3
4 ##### sets up dataframe with corresponding topic-numbers and amount numbers
5 topics_lda = pd.DataFrame(amounts, columns = topic_list_index_lda)
6
7 ##### sets up dataframe with corresponding topics and subjects
8 topics_list_lda = pd.DataFrame(amounts, columns = topic_list_lda)
9
10 topics_lda
```

Out[30]:

	Topic 1	Topic 2	Topic 3	Topic 4	Topic 5	Topic 6	Topic 7	Topic 8	Topic 9	Topic 10
0	0.012183	0.012182	0.012183	0.012182	0.012182	99.890359	0.012183	0.012183	0.012183	0.012183
1	0.014840	0.014839	0.014840	0.014839	0.014839	99.866447	0.014839	0.014839	0.014839	0.014839
2	23.976947	0.030310	14.822932	0.030308	0.030310	0.030309	49.192545	0.030308	11.825722	0.030311
3	63.740839	0.019497	19.967020	0.019496	0.019498	0.019497	9.034511	0.019498	0.019497	7.140647
4	0.011754	0.011753	0.011754	0.011752	0.011753	0.011752	0.011754	0.011752	0.011753	99.894222
...
70	0.013282	0.013284	0.013282	0.013281	0.013282	0.013282	0.013281	0.013282	0.013282	99.880461
71	0.013407	46.632949	0.013407	53.259794	0.013407	0.013407	0.013408	0.013406	0.013407	0.013407
72	0.019965	99.820326	0.019963	0.019964	0.019963	0.019963	0.019964	0.019963	0.019963	0.019967
73	0.028498	0.028497	0.028497	0.028494	1.266543	0.028493	93.560961	0.028494	0.028493	4.973030
74	0.013390	0.013389	55.967734	0.013388	0.013390	0.013389	43.925150	0.013389	0.013390	0.013390

75 rows × 10 columns

Storing topics list as a dataframe

```
In [31]: 1 topic_indices = list(topics_lda)
2 topic_subjects = list(topics_list_lda)
3 topic_desc = pd.DataFrame(topic_indices, topic_subjects).reset_index()
4 topic_desc_lda = topic_desc.rename(columns = {"index": "topic_desc", 0: "topic"})
5 topic_desc_lda
```

Out[31]:

	topic_desc	topic
0	ought_class_kid	Topic 1
1	terrorist_iraq_iraqi	Topic 2
2	vietnam_oil_communist	Topic 3
3	hussein_saddam_inspector	Topic 4
4	vietnam_south_voic	Topic 5
5	session_surplus_read	Topic 6
6	communist_agricultur_expenditur	Topic 7
7	veteran_disput_bargain	Topic 8
8	vietnam_januari_abund	Topic 9
9	21st_veteran_class	Topic 10

Merging speeches database with topics dataframe

```
In [32]: 1 lda_merged = topics_lda.join(sotu)
2 lda_merged.head(3)
```

Out[32]:

	Topic 1	Topic 2	Topic 3	Topic 4	Topic 5	Topic 6	Topic 7	Topic 8	Topic 9	Topic 10	text	filename	year	pres	id
0	0.012183	0.012182	0.012183	0.012182	0.012182	99.890359	0.012183	0.012183	0.012183	0.012183	Mr. President, Mr. Speaker, Members of the Con...	sotu/01-06-1955.txt	1955	eisenhower	eisenhower-1955
1	0.014840	0.014839	0.014840	0.014839	0.014839	99.866447	0.014839	0.014839	0.014839	0.014839	Mr. President, Mr. Speaker, Members of the Eig...	sotu/01-07-1954.txt	1954	eisenhower	eisenhower-1954
2	23.976947	0.030310	14.822932	0.030308	0.030310	49.192545	0.030308	11.825722	0.030311	0.030311	Mr. Speaker, Mr. President, my colleagues in t...	sotu/01-20-1972.txt	1972	nixon	nixon-1972

Rearranging dataframe for plotting

```
In [33]: 1 lda_1 = lda_merged[["text", "filename", "year", "pres", "id", "Topic 1"]].copy()
2 lda_1["topic"] = "Topic 1"
3 lda_1 = lda_1.rename(columns = {"Topic 1": "factor"})
4
5 lda_2 = lda_merged[["text", "filename", "year", "pres", "id", "Topic 2"]].copy()
6 lda_2["topic"] = "Topic 2"
7 lda_2 = lda_2.rename(columns = {"Topic 2": "factor"})
8
9 lda_3 = lda_merged[["text", "filename", "year", "pres", "id", "Topic 3"]].copy()
10 lda_3["topic"] = "Topic 3"
11 lda_3 = lda_3.rename(columns = {"Topic 3": "factor"})
12
13 lda_4 = lda_merged[["text", "filename", "year", "pres", "id", "Topic 4"]].copy()
14 lda_4["topic"] = "Topic 4"
15 lda_4 = lda_4.rename(columns = {"Topic 4": "factor"})
16
17 lda_5 = lda_merged[["text", "filename", "year", "pres", "id", "Topic 5"]].copy()
18 lda_5["topic"] = "Topic 5"
19 lda_5 = lda_5.rename(columns = {"Topic 5": "factor"})
20
21 lda_6 = lda_merged[["text", "filename", "year", "pres", "id", "Topic 6"]].copy()
22 lda_6["topic"] = "Topic 6"
23 lda_6 = lda_6.rename(columns = {"Topic 6": "factor"})
24
25 lda_7 = lda_merged[["text", "filename", "year", "pres", "id", "Topic 7"]].copy()
26 lda_7["topic"] = "Topic 7"
27 lda_7 = lda_7.rename(columns = {"Topic 7": "factor"})
28
29 lda_8 = lda_merged[["text", "filename", "year", "pres", "id", "Topic 8"]].copy()
30 lda_8["topic"] = "Topic 8"
31 lda_8 = lda_8.rename(columns = {"Topic 8": "factor"})
32
33 lda_9 = lda_merged[["text", "filename", "year", "pres", "id", "Topic 9"]].copy()
34 lda_9["topic"] = "Topic 9"
35 lda_9 = lda_9.rename(columns = {"Topic 9": "factor"})
36
37 lda_10 = lda_merged[["text", "filename", "year", "pres", "id", "Topic 10"]].copy()
38 lda_10["topic"] = "Topic 10"
39 lda_10 = lda_10.rename(columns = {"Topic 10": "factor"})
```

```

41 lda_master = pd.concat([lda_1, lda_2, lda_3, lda_4, lda_5, lda_6, lda_7, lda_8, lda_9, lda_10], ignore_index = True)
42 lda_master = pd.merge(lda_master, topic_desc_lda, on = "topic", how = "left")
43
44 lda_master

```

Out[33]:

	text	filename	year	pres	id	factor	topic	topic_desc
0	Mr. President, Mr. Speaker, Members of the Con...	sotu/01-06-1955.txt	1955	eisenhower	eisenhower-1955	0.012183	Topic 1	ought_class_kid
1	Mr. President, Mr. Speaker, Members of the Eig...	sotu/01-07-1954.txt	1954	eisenhower	eisenhower-1954	0.014840	Topic 1	ought_class_kid
2	Mr. Speaker, Mr. President, my colleagues in t...	sotu/01-20-1972.txt	1972	nixon	nixon-1972	23.976947	Topic 1	ought_class_kid
3	Mr. President, Mr. Speaker, Members of the Hou...	sotu/02-17-1993.txt	1993	clinton	clinton-1993	63.740839	Topic 1	ought_class_kid
4	Madam Speaker, Madam Vice President—no Preside...	sotu/04-28-2021.txt	2021	biden	biden-2021	0.011754	Topic 1	ought_class_kid
...
745	Mr. Speaker, Mr. Vice President, Members of Co...	sotu/01-30-2018.txt	2018	trump	trump-2018	99.880461	Topic 10	21st_veteran_class
746	Mr. Speaker, Vice President Cheney, Members of...	sotu/01-28-2003.txt	2003	bush jr.	bush jr-2003	0.013407	Topic 10	21st_veteran_class
747	Thank you very much. Mr. Speaker, Vice Preside...	sotu/01-29-2002.txt	2002	bush jr.	bush jr-2002	0.019967	Topic 10	21st_veteran_class
748	Mr. President, Mr. Speaker, Members of the Con...	sotu/01-08-1951.txt	1951	truman	truman-1951	4.973030	Topic 10	21st_veteran_class
749	Mr. Vice President, my old colleague from Mass...	sotu/01-11-1962.txt	1962	kennedy	kennedy-1962	0.013390	Topic 10	21st_veteran_class

750 rows × 8 columns

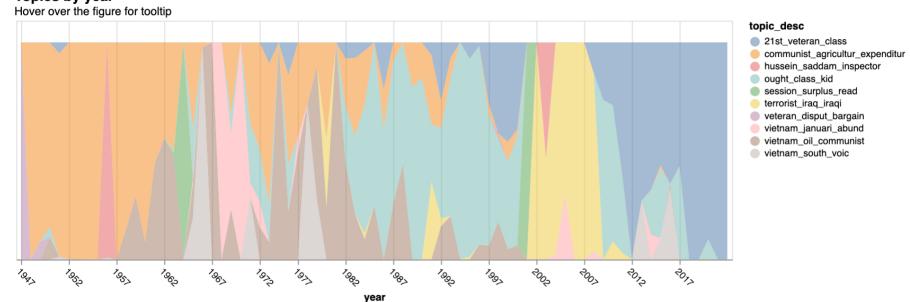
Plotting steamgraph

```

1 alt.Chart(lda_master).mark_area(opacity = 0.5).encode(
2     alt.X("year:O", axis = alt.Axis(values = list(range(1947, 2022, 5)))),
3     alt.Y("sum(factor):Q", stack = "center", axis = None),
4     alt.Color("topic_desc", scale = alt.Scale(scheme = "tableau10")),
5     tooltip = ["pres", "year", "topic_desc"]
6 ).configure_axis(
7     grid = True
8 ).configure_axisX(
9     labelAngle = 45
10 ).properties(
11     width = 750,
12     height = 250,
13     title = {
14         "text": "Topics by year",
15         "subtitle": "Hover over the figure for tooltip",
16     }
17 ).configure_title(
18     fontSize = 15,
19     anchor = "start",
20 ).interactive()

```

Out[34]: Topics by year



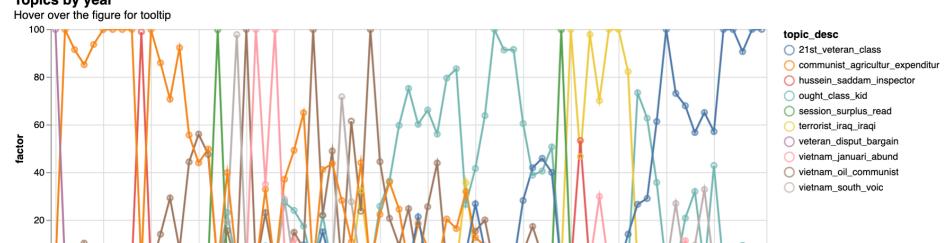
Plotting line chart

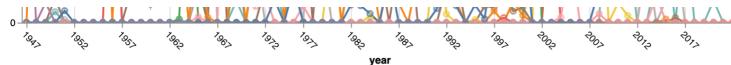
```

1 lines = alt.Chart(lda_master).mark_line().encode(
2     alt.Color("topic_desc", scale = alt.Scale(scheme = "tableau10")),
3     alt.X("year:O", axis = alt.Axis(values = list(range(1947, 2022, 5)))),
4     y = "factor",
5 )
6
7 points = alt.Chart(lda_master).mark_point().encode(
8     alt.Color("topic_desc", scale = alt.Scale(scheme = "tableau10")),
9     alt.X("year:O", axis = alt.Axis(values = list(range(1947, 2022, 5)))),
10    y = "factor",
11    tooltip = ["pres", "year", "topic_desc"]
12 )
13
14 chart = lines + points
15
16 chart.configure_axis(
17     grid = True
18 ).configure_axisX(
19     labelAngle = 45
20 ).properties(
21     width = 750,
22     height = 250,
23     title = {
24         "text": "Topics by year",
25         "subtitle": "Hover over the figure for tooltip",
26     }
27 ).configure_title(
28     fontSize = 15,
29     anchor = "start",
30 ).interactive()

```

Out[35]: Topics by year





Comparing NMF and LDA models:

The topics generated by both models are nearly identical but the ones from NMF seem to make a bit more sense.

Note: When applying topic modeling, sometimes which method is chosen boils down to a subjective decision. GridSearchCV can be used to deploy an algorithm to suggest the best parameters but that too is not always foolproof. The primary question should always be "What makes sense?"

Gensim modeling

Setting dataframe

```
In [36]: 1 sotu.text = sotu.text.str.replace("[^A-Za-z ]", " ")
2
3 sotu.head(3)
/var/folders/8/_vf0g3np57bdb6_x59hshmtcr0000gn/T/ipykernel_49503/1114454686.py:1: FutureWarning: The default value of regex will change from True to False in a
future version.
    sotu.text = sotu.text.str.replace("[^A-Za-z ]", " ")

Out[36]:   text      filename  year  pres          id
0  Mr President Mr Speaker Members of the Con...  sotu/01-06-1955.txt  1955  eisenhower  eisenhower-1955
1  Mr President Mr Speaker Members of the Eig...  sotu/01-07-1954.txt  1954  eisenhower  eisenhower-1954
2  Mr Speaker Mr President my colleagues in t...  sotu/01-20-1972.txt  1972    nixon    nixon-1972
```

Deploying LSI/NMF modeling with Gensim

```
In [37]: 1 ### imports modules
2 from gensim.utils import simple_preprocess
3
4 ### stores text of speeches and pre-processes
5 texts = sotu.text.apply(simple_preprocess)
6
7 ### imports modules
8 from gensim import corpora
9
10 dictionary = corpora.Dictionary(texts)
11 dictionary.filter_extremes(no_below = 5, no_above = 0.4)
12 corpus = [dictionary.doc2bow(text) for text in texts]
13
14 ### imports modules
15 from gensim import models
16
17 tfidf = models.TfidfModel(corpus)
18 corpus_tfidf = tfidf[corpus]
19
20 ### assigns number of topics
21 n_topics = 10
22
23 ### builds a model
24 nmf_model = models.LsiModel(corpus_tfidf,
25                             id2word = dictionary,
26                             num_topics = n_topics)
27
28 ### displays topic results
29 nmf_model.print_topics()

Out[37]: [(0,
  '0.085*"iraq" + 0.067*"vietnam" + 0.065*"oil" + 0.064*"terrorists" + 0.058*"kids" + 0.057*"recommend" + 0.055*"companies" + 0.054*"ought" +
  0.053*"communist"),
 (1,
  '0.135*"iraq" + 0.104*"terrorists" + -0.096*"expenditures" + -0.093*"recommend" + 0.087*"al" + -0.082*"peoples" + -0.082*"communist" + 0.081
  *"iraqi" + 0.079*"qaida"),
 (2,
  '0.325*"iraq" + 0.262*"iraqi" + 0.178*"terrorists" + 0.176*"saddam" + 0.169*"iraqis" + 0.168*"terror" + 0.154*"hussein" + 0.134*"al" + 0.129
  *"ought"),
 (3,
  '-0.247*"ryan" + -0.209*"isis" + -0.162*"gang" + -0.119*"heroes" + -0.119*"totally" + -0.118*"immigration" + -0.115*"incredible" + 0.104*"ou
  terans" + -0.090*"david"),
 (4,
  '-0.243*"vietnam" + -0.111*"nicaragua" + 0.103*"folks" + 0.082*"oil" + 0.082*"manufacturing" + -0.078*"drugs" + -0.075*"revolution" + 0.071*
  "veterans" + 0.070*"collective"),
 (5,
  '-0.202*"ought" + 0.170*"oil" + 0.119*"vietnam" + -0.101*"millennium" + -0.097*"saddam" + -0.096*"bosnia" + -0.082*"tobacco" + -0.081*"grade
  in" + -0.073*"americorps"),
 ...]
```

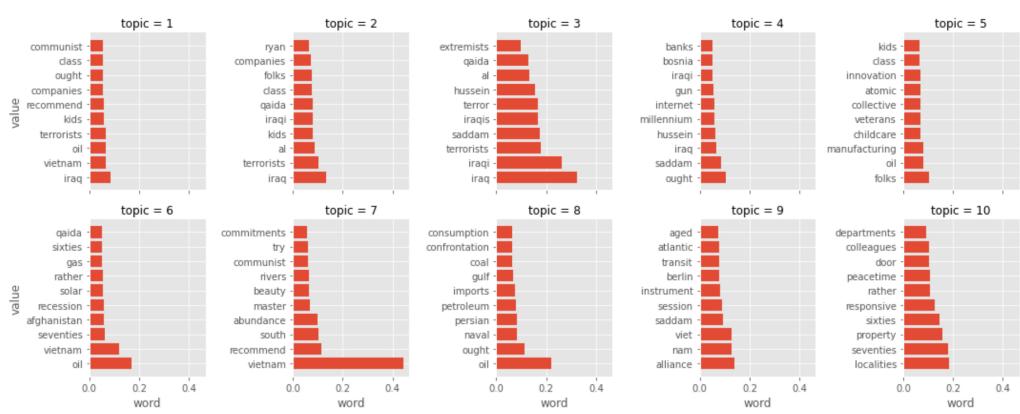
```
In [38]: 1 ### assigns number of words per topic
2 n_words = 10
3
4 topic_words = pd.DataFrame({})
5
6 ### computes dataframe with results
7 for i, topic in enumerate(nmf_model.get_topics()):
8     top_feature_ids = topic.argsort()[-n_words:][:-1]
9     feature_values = topic[top_feature_ids]
10    words = {dictionary[id] for id in top_feature_ids}
11    topic_df = pd.DataFrame({"value": feature_values, "word": words, "topic": (i + 1)})
12    topic_words = pd.concat([topic_words, topic_df], ignore_index = True)
13
14 topic_words.head()
```

```
Out[38]:   value      word  topic
0  0.085336    iraq    1
1  0.066523  vietnam    1
2  0.064880      oil    1
3  0.063994  terrorists    1
4  0.057784      kids    1
```

Visualizing Gensim-LSI/NMF model using Seaborn

```
In [39]: 1 import matplotlib.pyplot as plt
2 plt.style.use("ggplot")
3 import seaborn as sns
4
```

```
5 g = sns.FacetGrid(topic_words, col = "topic", col_wrap = 5, sharey = False)
6 g.map(plt.barh, "word", "value")
```



Deploying LDA model with Gensim

```
In [40]: 1  ### imports modules
2  from gensim.utils import simple_preprocess
3
4  ### stores text of speeches and pre-processes
5  texts = sotu.text.apply(simple_preprocess)
6
7  ### imports modules
8  from gensim import corpora
9
10 dictionary = corpora.Dictionary(texts)
11 dictionary.filter_extremes(no_below = 5, no_above = 0.4, keep_n = 2000)
12 corpus = [dictionary.doc2bow(text) for text in texts]
13
14 ### assigns number of topics
15 from gensim import models
16
17 ### builds a model
18 n_topics = 10
19 lda_model = models.LdaModel(corpus = corpus, num_topics = n_topics)
20
21 ### displays topic results
22 lda_model.print_topics()
```

```
Out[40]: [(0,
  '0.004*\"1497" + 0.003*\"1290" + 0.002*\"1084" + 0.002*\"1404" + 0.002*\"623" + 0.002*\"1514" + 0.002*\"1334" + 0.002*\"1021" + 0.002*\"1531" + 0.002
(1,
  '0.003*\"1514" + 0.003*\"1497" + 0.003*\"980" + 0.003*\"425" + 0.002*\"1290" + 0.002*\"623" + 0.002*\"1099" + 0.002*\"189" + 0.002*\"1404" + 0.002*\"8
(2,
  '0.004*\"1497" + 0.003*\"1021" + 0.003*\"1514" + 0.003*\"1404" + 0.003*\"623" + 0.003*\"559" + 0.002*\"262" + 0.002*\"1229" + 0.002*\"1334" + 0.002*
(3,
  '0.003*\"425" + 0.003*\"1404" + 0.003*\"1084" + 0.002*\"1497" + 0.002*\"559" + 0.002*\"84" + 0.002*\"1290" + 0.002*\"79" + 0.002*\"980" + 0.002*\"189"
(4,
  '0.005*\"1497" + 0.004*\"1404" + 0.003*\"1514" + 0.002*\"1229" + 0.002*\"1721" + 0.002*\"456" + 0.002*\"425" + 0.002*\"1290" + 0.002*\"1364" + 0.002*
(5,
  '0.003*\"623" + 0.003*\"1497" + 0.002*\"1229" + 0.002*\"1514" + 0.002*\"1404" + 0.002*\"1721" + 0.002*\"1084" + 0.002*\"1334" + 0.002*\"1230" + 0.002
(6,
  '0.003*\"1290" + 0.003*\"79" + 0.002*\"1021" + 0.002*\"425" + 0.002*\"623" + 0.002*\"1045" + 0.002*\"1229" + 0.002*\"980" + 0.002*\"1084" + 0.002*\"15
(7,
  '0.003*\"1514" + 0.002*\"1404" + 0.002*\"189" + 0.002*\"1497" + 0.002*\"1721" + 0.002*\"623" + 0.002*\"7" + 0.002*\"1531" + 0.002*\"425" + 0.002*\"102
(8,
  '0.003*\"1497" + 0.003*\"1404" + 0.002*\"623" + 0.002*\"425" + 0.002*\"1290" + 0.002*\"1021" + 0.002*\"189" + 0.002*\"1084" + 0.002*\"1045" + 0.002*"
(9,
```

```
In [41]: 1  ### assigns number of words per topic
2  n_words = 10
3
4  topic_words = pd.DataFrame({})
5
6  ### computes dataframe with results
7  for i, topic in enumerate(lda_model.get_topics()):
8      top_feature_ids = topic.argsort()[-n_words:][::-1]
9      feature_values = topic[top_feature_ids]
10     words = [dictionary[id] for id in top_feature_ids]
11     topic_df = pd.DataFrame({'value': feature_values, "word": words, "topic": (i + 1)})
12     topic_words = pd.concat([topic_words, topic_df], ignore_index = True)
13
14 topic_words.head()
```

Out[41]:

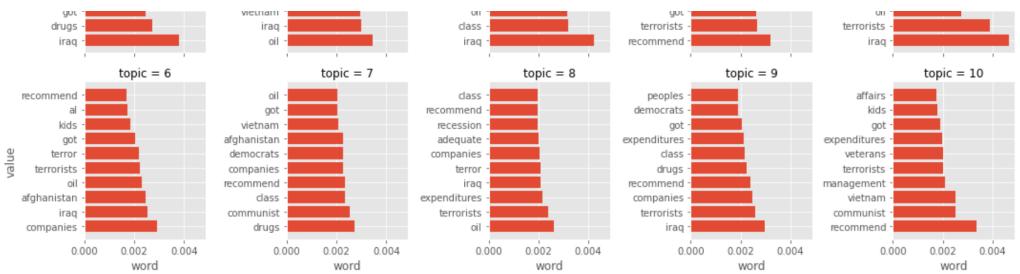
	value	word	topic
0	0.003800	iraq	1
1	0.002715	drugs	1
2	0.002465	got	1
3	0.002459	terrorists	1
4	0.002437	companies	1

Visualizing Gensim-LDA model using Seaborn

```
In [42]: 1 import matplotlib.pyplot as plt
2 plt.style.use("ggplot")
3 import seaborn as sns
4
5 g = sns.FacetGrid(topic_words, col = "topic", col_wrap = 5, sharey = False)
6 g.map(plt.barh, "word", "value")
```

A facet grid visualization showing word frequency for five different topics. The facets are labeled 'topic = 1' through 'topic = 5'. Each facet contains a horizontal bar chart where the x-axis represents word frequency and the y-axis lists words. The words are ordered by frequency within each topic.

topic	words (frequency)
topic = 1	idea, recession, class, kids, oil, companies
topic = 2	try, terrorists, expenditures, idea, companies, drugs
topic = 3	republicans, kids, afghanistan, immigration, veterans, companies
topic = 4	expenditures, vietnam, communist, drugs, concern, veterans
topic = 5	affordable, qaida, drugs, recommend, retirement, terror, veterans



Visualizing Gensim-LDA model using PyLDAvis

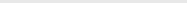
```
In [43]: 1 import pyLDAvis
2 import pyLDAvis.gensim_models as gensimvis
3
4 pyLDAvis.enable_notebook()
5 vis = gensimvis.prepare(lda_model, corpus, dictionary)
6 vis
```

```
/Users/shirshdasgupta/Documents/Python/anaconda3/lib/python3.9/site-packages/pyLDAvis/_prepare.py:246: FutureWarning: In a future version of pandas all arguments of Dataframe.drop except for the argument 'labels' will be keyword-only.
    default_term_info = default_term_info.sort_values(
/Users/shirshdasgupta/Documents/Python/anaconda3/lib/python3.9/site-packages/past/builtins/misc.py:45: DeprecationWarning: the imp module is deprecated in favor of importlib; see the module's documentation for alternative uses
    from imp import reload
/Users/shirshdasgupta/Documents/Python/anaconda3/lib/python3.9/site-packages/past/builtins/misc.py:45: DeprecationWarning: the imp module is deprecated in favor of importlib; see the module's documentation for alternative uses
    from imp import reload
/Users/shirshdasgupta/Documents/Python/anaconda3/lib/python3.9/site-packages/past/builtins/misc.py:45: DeprecationWarning: the imp module is deprecated in favor of importlib; see the module's documentation for alternative uses
    from imp import reload
/Users/shirshdasgupta/Documents/Python/anaconda3/lib/python3.9/site-packages/past/builtins/misc.py:45: DeprecationWarning: the imp module is deprecated in favor of importlib; see the module's documentation for alternative uses
    from imp import reload
/Users/shirshdasgupta/Documents/Python/anaconda3/lib/python3.9/site-packages/past/builtins/misc.py:45: DeprecationWarning: the imp module is deprecated in favor of importlib; see the module's documentation for alternative uses
    from imp import reload
/Users/shirshdasgupta/Documents/Python/anaconda3/lib/python3.9/site-packages/past/builtins/misc.py:45: DeprecationWarning: the imp module is deprecated in favor of importlib; see the module's documentation for alternative uses
    from imp import reload
/Users/shirshdasgupta/Documents/Python/anaconda3/lib/python3.9/site-packages/past/builtins/misc.py:45: DeprecationWarning: the imp module is deprecated in favor of importlib; see the module's documentation for alternative uses
```

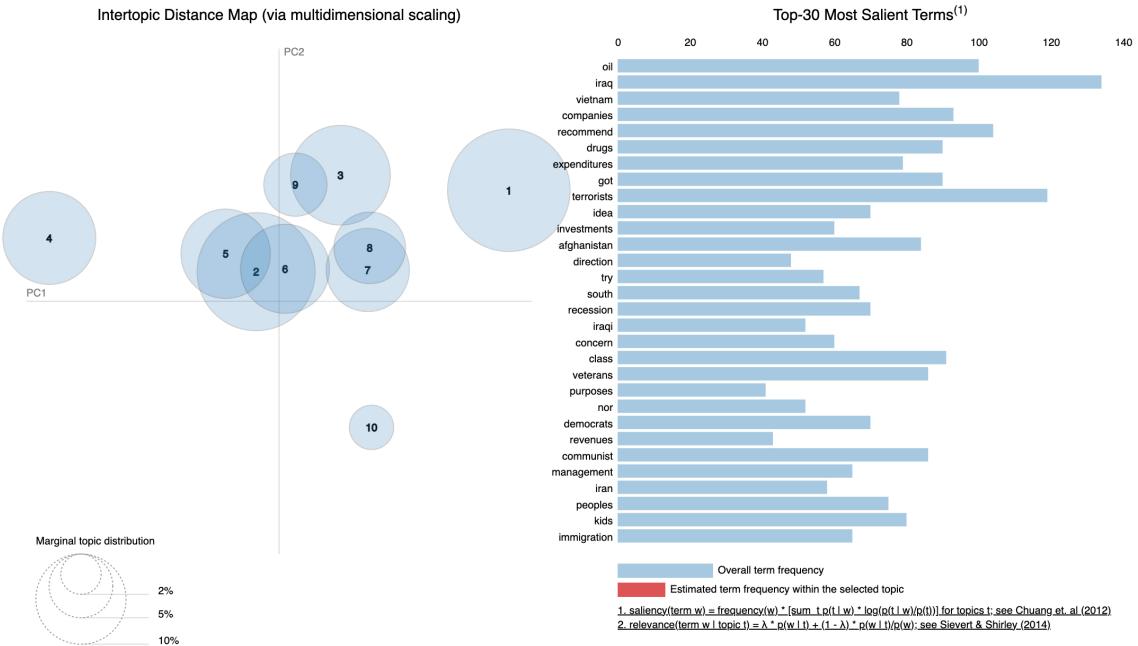
Out[43]: Selected Topic: 0 [Previous Topic](#) [Next Topic](#) [Clear Topic](#)

Slide to adjust relevance metric:(2)

$\lambda = 1$



0.0 0.2 0.4 0.6 0.8 1



In []: 1