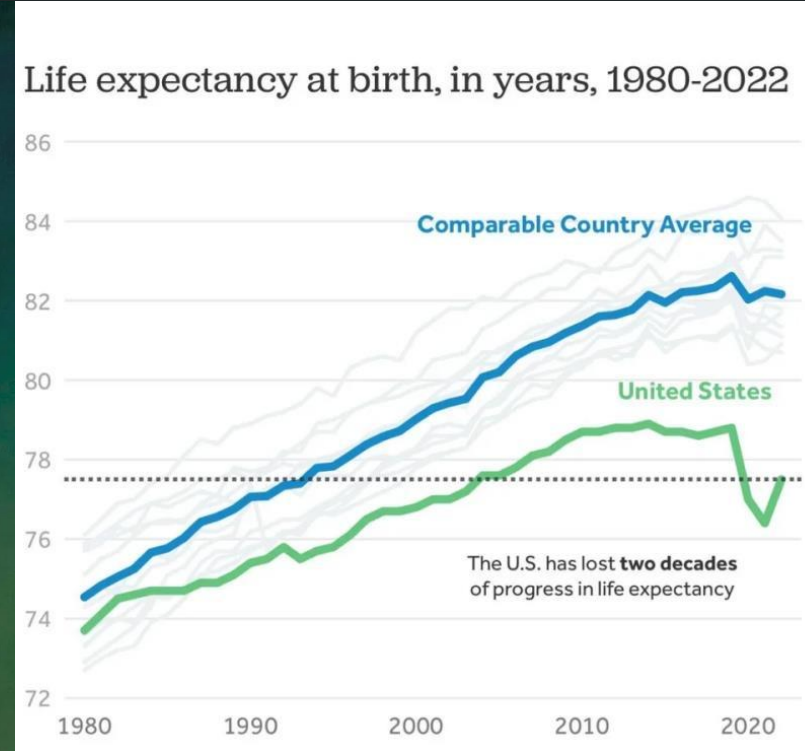
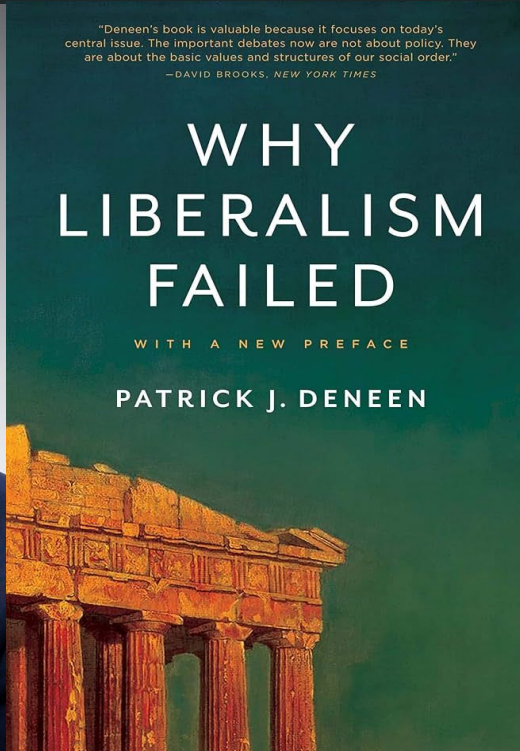


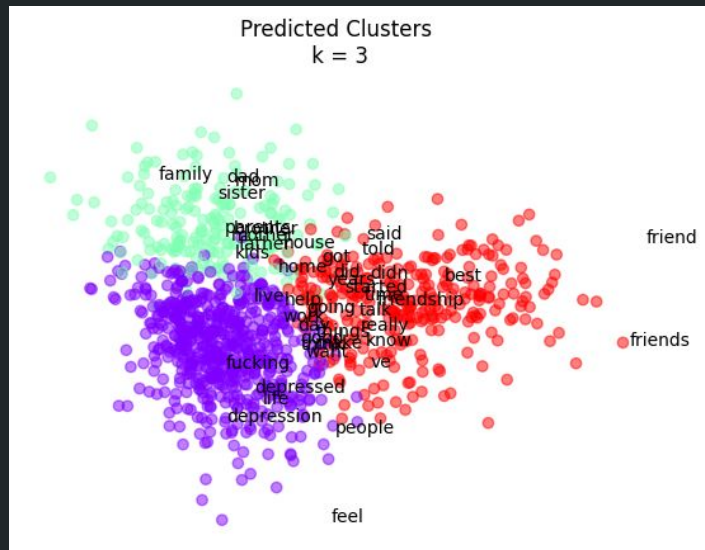
Final Project: Purpose & Meaning

Dan Gilles

Inspiration



K-Means analysis



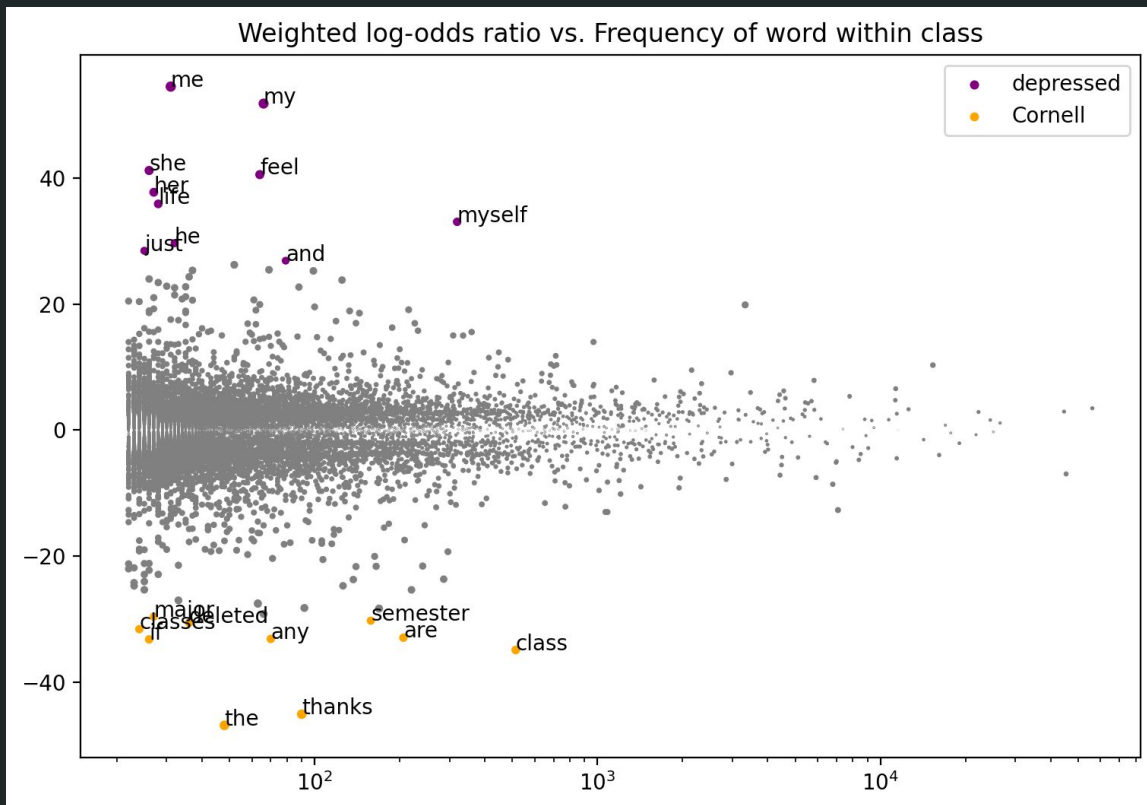
Top terms per cluster: (3 clusters)

['Cluster 0:', ' feel', ' life', ' want', ' depression', ' people', ' ve', ' know', ' think', ' really', ' time', ' depressed', ' things', ' help', ' day', ' fucking', ' make', ' work', ' good', ' going', ' years', ' happy', ' better', ' feeling', ' hate', ' love', ' wish', ' job', ' talk', ' makes', ' need', ' try', ' bad', ' right', ' sad', ' thoughts', ' having', ' anymore', ' shit', ' friends', ' year', ' today', ' tired', ' way', ' fuck', ' doing', ' school', ' long', ' getting', ' say', ' does',

'Cluster 1:', ' friend', ' friends', ' best', ' know', ' ve', ' feel', ' time', ' really', ' people', ' want', ' said', ' didn', ' told', ' friendship', ' talk', ' make', ' years', ' did', ' things', ' started', ' school', ' got', ' think', ' group', ' person', ' new', ' life', ' say', ' going', ' girl', ' guy', ' way', ' doesn', ' help', ' day', ' tell', ' hang', ' bad', ' talking', ' boyfriend', ' good', ' year', ' close', ' recently', ' months', ' sure', ' asked', ' thing', ' advice', ' need',

'Cluster 2:', ' family', ' dad', ' mom', ' sister', ' parents', ' mother', ' brother', ' father', ' house', ' time', ' want', ' home', ' know', ' years', ' kids', ' got', ' told', ' live', ' really', ' said', ' husband', ' help', ' feel', ' old', ' things', ' didn', ' year', ' wife', ' ve', ' say', ' cousin', ' doesn', ' away', ' money', ' going', ' think', ' went', ' daughter', ' does', ' child', ' work', ' day', ' ago', ' school', ' car', ' mum', ' right', ' died', ' asked', ' care']

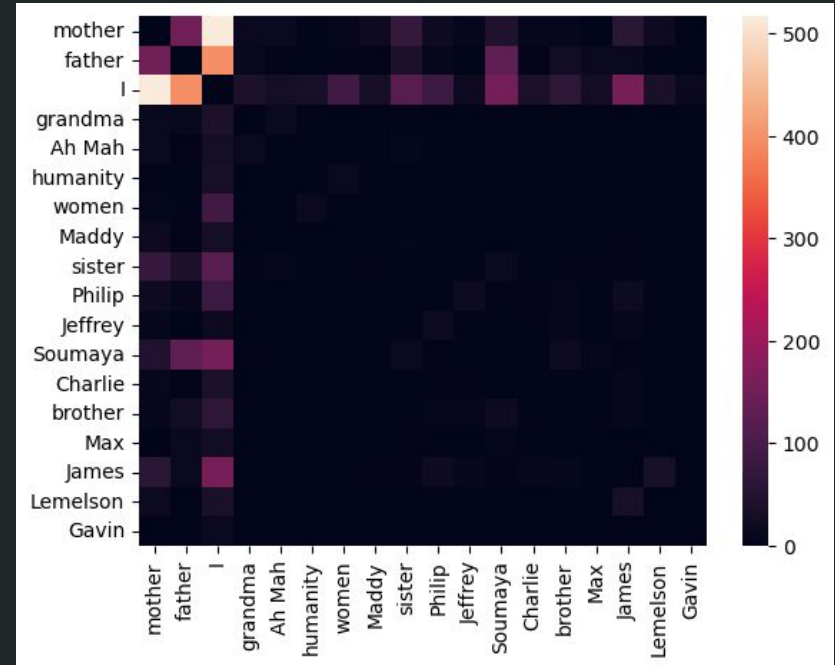
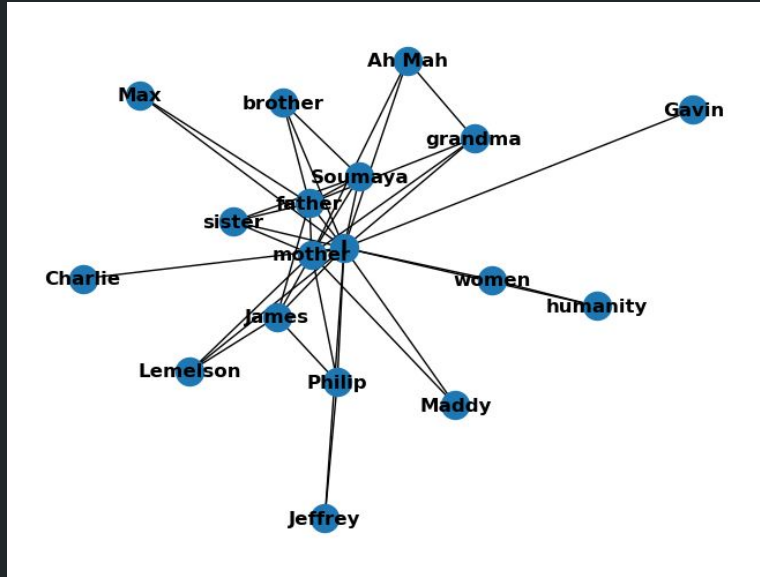
Word frequency differences



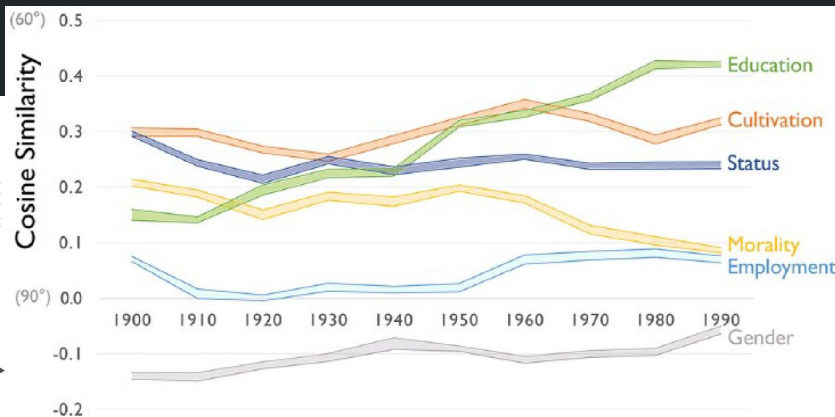
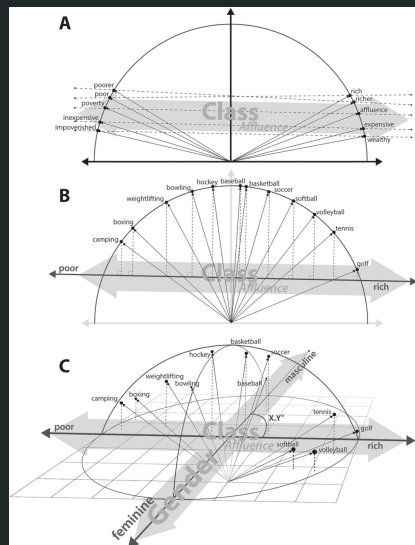
	z-score	class
ngram		
the	-46.813787	Cornell
thanks	-45.052157	Cornell
class	-34.871888	Cornell
if	-33.188343	Cornell
any	-33.114831	Cornell
...
her	37.721990	depressed
feel	40.518676	depressed
she	41.175328	depressed
my	51.782092	depressed
me	54.478234	depressed

15000 rows x 2 columns

Network analysis



Word embeddings



```
...
'-2.490995526313781738e-01 1.513628512620925903e-01 1.627609878778457642e-02 -2.899127304553985596e-01 1.994638331234455109e-02
1.047520861829624939e-01 -9.935868531465530396e-02 -2.254071988714797974e-01 3.088236153125762939e-01 3.820571303367614746e-02
6.73935562372076416e-02 -3.0886010694950378418e-01 -7.813987135887145996e-02 3.222133815288543701e-01 1.317477226257324219e-01
1.605398459299087524e-01 2.64114588499692139e-01 -2.05675289034843448e-02 1.061243712902609092e-01 1.458893120288848877e-01
3.055666014552116394e-02 -3.502352535724639893e-02 5.022452399134635925e-02 4.805423691868782043e-02 4.483234882354736328e-02
1.564268767833789717e-01 -2.481821775436401367e-01 -9.431868791580200195e-01 1.273611746728420258e-02 2.694861590862274170e-01
-8.399840444326400757e-02 -2.579927444458007812e-01 -5.300612002611160278e-02 -1.724688364131927490e-01 -2.331216633319854736e-02
-8.260920643806457520e-02 -7.970795788420383345e-02 -3.962805122137069702e-02 -1.564607769250869751e-01 2.645880877971649170e-01
-2.503050863742828369e-02 -2.256768494844436646e-01 7.921236753463745117e-02 1.294871885329484940e-02 -2.509132921695709229e-01
-3.748797997832298279e-02 8.977919071912765503e-02 -2.409224361181259155e-01 -6.001283879714012146e-02 -2.081173509359359741e-01
-2.2428885102272083369e-01 -4.367685317993164062e-01 4.219847725200653076e-01 1.48108065128326160e-01 -8.94341079387664795e-02
-1.915721036493778229e-02 -8.969138234853744507e-02 -2.136332541704177856e-01 6.727373600060103516e-02 3.982214492948663174e-03
2.2375484973434448e-01 3.088929622173309326e-01 3.113623635447780239e-02 -1.1491076648233310e-01 3.556222219602963609e-01
1.673831206813725146e-02 -2.2526720623806100e-01 -1.008317503929138104e-01 -4.0040755780580077e-02 -2.4377236675410614e-01
3.614936023958576702e-02 2.05481631758588574e-01 -1.72097970611188354e-01 -1.522434338782430176e-01 -1.006038171852937730e-01
1.344614382833242416e-02 -7.894081622362136841e-02 8.38536024036279297e-02 -1.52665256748199463e-01 6.38793334648910522e-02
2.209412604570388794e-01 -2.350811148062348366e-02 2.722778916358847754e-01 6.32160827515994604e-02 -1.082860814237594604e-01
3.690200914306640625e-01 7.115402817726135254e-02 -2.922659739851951599e-03 -2.514926493167877197e-01 1.480864042323857800e-02
1.137095987796783447e-01 -3.514844775199890137e-01 3.623995184898376465e-01 1.382033515036773682e-01 5.965615008963775635e-02
-1.012735366821289062e-01 1.3107146322272728337e-01 -1.72013590615987777e-02 4.326710104942321777e-01 3.409578083022117615e-02
-1.16186784580349222e-02 1.492834836244583130e-01 4.845631122589111328e-02 1.368832141160964966e-01 6.269634515047073364e-02
7.08634318223953247e-02 6.277590833679962158e-02 -1.750740176208066700e-01 1.022537946781049805e-01 -3.569718526924213330e-01
1.389058073982596397e-03 -2.133055831299591064e-01 -2.727279812097549438e-02 -1.6921156644822116999e-01 -1.856350749731063843e-01
-7.838120311498641968e-02 -5.291774868965148926e-01 1.172335594892501831e-01 1.879646331071853638e-01 7.884481549263000408e-02
2.23128691315509399e-01 1.292670071125030518e-01 -1.811326593160629272e-01 -1.280843615531921387e-01 -2.8077837508760833740e-01
-1.997592002153396006e-01 9.234853088855743408e-02 -4.534598812460899353e-02 1.868978515267372131e-02 -2.621622383594512939e-01
```

w2v_ngram_models	32.38 GB
ngram_w2v_readme.txt	1 KB
syn0_ngram_1900_1909_full.txt	2.7 GB
syn0_ngram_1910_1919_full.txt	2.56 GB
syn0_ngram_1920_1929_full.txt	2.66 GB
syn0_ngram_1930_1939_full.txt	2.55 GB
syn0_ngram_1940_1949_full.txt	2.44 GB
syn0_ngram_1950_1959_full.txt	2.83 GB
syn0_ngram_1960_1969_full.txt	3.66 GB
syn0_ngram_1970_1979_full.txt	3.88 GB
syn0_ngram_1980_1989_full.txt	4.11 GB
syn0_ngram_1990_1999_full.txt	4.96 GB
vocab_list_ngram_1900_1909_full.txt	3.1 MB
vocab_list_ngram_1910_1919_full.txt	2.9 MB
vocab_list_ngram_1920_1929_full.txt	3 MB
vocab_list_ngram_1930_1939_full.txt	2.9 MB
vocab_list_ngram_1940_1949_full.txt	2.7 MB
vocab_list_ngram_1950_1959_full.txt	3.2 MB
vocab_list_ngram_1960_1969_full.txt	4.2 MB
vocab_list_ngram_1970_1979_full.txt	4.4 MB
vocab_list_ngram_1980_1989_full.txt	4.6 MB
vocab_list_ngram_1990_1999_full.txt	5.7 MB

Word embeddings

✓ w2v_ngram_models	32.38 GB
ngram_w2v_readme.txt	1 KB
syn0_ngram_1900_1909_full.txt	2.7 GB
syn0_ngram_1910_1919_full.txt	2.56 GB
syn0_ngram_1920_1929_full.txt	2.66 GB
syn0_ngram_1930_1939_full.txt	2.55 GB
syn0_ngram_1940_1949_full.txt	2.44 GB
syn0_ngram_1950_1959_full.txt	2.83 GB
syn0_ngram_1960_1969_full.txt	3.66 GB
syn0_ngram_1970_1979_full.txt	3.88 GB
syn0_ngram_1980_1989_full.txt	4.11 GB
syn0_ngram_1990_1999_full.txt	4.96 GB
vocab_list_ngram_1900_1909_full.txt	3.1 MB
vocab_list_ngram_1910_1919_full.txt	2.9 MB

```
for year in start_years:
    vocab = pd.read_csv(f'/Volumes/Dan2/w2v_ngram_models/vocab_list_ngram_{year}_{year+9}_full.txt')
    df_chunk = pd.read_csv(f'/Volumes/Dan2/w2v_ngram_models/syn0_ngram_{year}_{year+9}_full.txt',
                           n = 0)

    with open(f'/Volumes/Dan2/ngrams/ngrams_{year}_{year+9}_{VOCAB_LEN}.txt', 'w') as f:
        f.write(str(VOCAB_LEN) + ' 300\n')
        for chunk in df_chunk:
            vec = chunk.iloc[0, 0]
            idx = chunk.index[0]
            # vectors.loc[n] = [idx, vocab.iloc[idx][0], vec]
            f.write(f'{vocab.iloc[idx][0]} {vec}\n')
            if n % 1000 == 0:
                print(year, n)
            n += 1
        if n == VOCAB_LEN:
            break
```

	0
0	the
1	of
2	,
3	to
4	.
...	...
334024	urticastrum

```
for year in start_years:
    print(year, vars()[f'vecs_{year}']).
    most_similar(positive=['woman',
                          'king'], negative=['man'], topn=1))
```

✓ 0.1s

Python

```
1900 [('inside', 0.44141054153442383)]
1910 [('queen', 0.7039340138435364)]
1920 [('mouth', 0.47583141922950745)]
1930 [('hydrogen', 0.41790416836738586)]
1940 [('employed', 0.4347355365753174)]
1950 [('soon', 0.5123521685600281)]
1960 [('held', 0.44984015822410583)]
1970 [('growth', 0.5318295359611511)]
1980 [('wife', 0.42465054988861084)]
1990 [('pounds', 0.4013866186141968)]
```

```
for year in start_years:
    print(year, vars()[f'vecs_{year}']).
    most_similar(positive=['woman',
                          'king'], negative=['man'], topn=1))
```

✓ 0.0s

Python

```
1900 [('queen', 0.7076151371002197)]
1910 [('queen', 0.7039340138435364)]
1920 [('queen', 0.7088885307312012)]
1930 [('queen', 0.7042350769042969)]
1940 [('queen', 0.6755873560905457)]
1950 [('queen', 0.7101660370826721)]
1960 [('queen', 0.7316992878913879)]
1970 [('queen', 0.6998559236526489)]
1980 [('queen', 0.6730139255523682)]
1990 [('queen', 0.6461030840873718)]
```

Antonym pairs

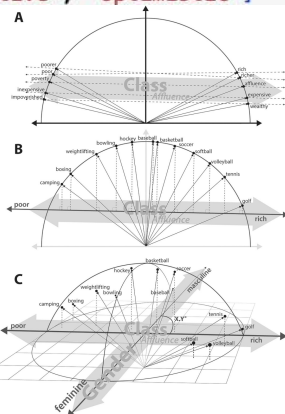
```
axes = ['purpose', 'career', 'happiness', 'satisfaction', 'anxiety', 'positivity', 'success', 'power', 'age', 'self', 'alone', 'wealth']
purpose_words_pos = ['purpose', 'significance', 'significant', 'meaning']#, 'worthwhile', 'meaningful']
purpose_words_neg = ['purposeless', 'insignificance', 'insignificant', 'meaningless']#, 'worthless', 'meaningless']
career_words_pos = ['employed', 'skilled', 'occupied', 'professional']
career_words_neg = ['unemployed', 'unskilled', 'unoccupied', 'idle']
happiness_words_pos = ['happy', 'joyful', 'glad']
happiness_words_neg = ['sad', 'sorrowful', 'grieving']
satisfaction_words_pos = ['satisfied', 'content', 'satisfaction']
satisfaction_words_neg = ['dissatisfied', 'discontent', 'dissatisfaction']
anxiety_words_pos = ['calm', 'unafraid', 'fearless']
anxiety_words_neg = ['anxious', 'afraid', 'fearful']
positivity_words_pos = ['hopeful', 'positive', 'optimistic']
positivity_words_neg = ['hopeless', 'negative', 'pessimistic']
success_words_pos = ['successful', 'thriving']
success_words_neg = ['failing', 'struggling']
power_words_pos = ['powerful', 'robust', 'strong']
power_words_neg = ['powerless', 'fragile', 'weak']
age_words_pos = ['young', 'youthful', 'energetic']
age_words_neg = ['old', 'decrepit', 'drained']
self_words_pos = ['self', 'myself', 'i', 'me', 'mine']
self_words_neg = ['family', 'friends', 'us', 'we', 'our']
alone_words_pos = ['alone', 'lonely', 'isolated']
alone_words_neg = ['together', 'connected', 'united']
wealth_words_pos = ['wealthy', 'rich', 'affluent']
wealth_words_neg = ['poor', 'impoverished', 'destitute']
```


Antonym pairs

```
axes = ['purpose', 'career', 'happiness', 'satisfaction', 'anxiety', 'positivity', 'success', 'power', 'age', 'self', 'alone', 'wealth']
purpose_words_pos = ['purpose', 'significance', 'significant', 'meaning']#, 'worthwhile', 'meaningful']
purpose_words_neg = ['purposeless', 'insignificance', 'insignificant', 'meaningless']#, 'worthless', 'meaningless']
career_words_pos = ['employed', 'skilled', 'occupied', 'professional']
career_words_neg = ['unemployed', 'unskilled', 'unoccupied', 'idle']
happiness_words_pos = ['happy', 'joyful', 'glad']
happiness_words_neg = ['sad', 'sorrowful', 'grieving']
satisfaction_words_pos = ['satisfied', 'content', 'satisfaction']
satisfaction_words_neg = ['dissatisfied', 'discontent', 'dissatisfaction']
anxiety_words_pos = ['calm', 'unafraid', 'fearless']
anxiety_words_neg = ['anxious', 'afraid', 'fearful']
positivity_words_pos = ['hopeful', 'positive', 'optimistic']
positivity_words_neg = ['hopeless', 'negative', 'pessimistic']
success_words_pos = ['successful', 'thriving']
success_words_neg = ['failing', 'struggling']
power_words_pos = ['powerful', 'robust', 'strong']
power_words_neg = ['powerless', 'fragile', 'weak']
age_words_pos = ['young', 'youthful', 'energetic']
age_words_neg = ['old', 'decrepit', 'drained']
self_words_pos = ['self', 'myself', 'i', 'me', 'mine']
self_words_neg = ['family', 'friends', 'us', 'we', 'our']
alone_words_pos = ['alone', 'lonely', 'isolated']
alone_words_neg = ['together', 'connected', 'united']
wealth_words_pos = ['wealthy', 'rich', 'affluent']
wealth_words_neg = ['poor', 'impoverished', 'destitute']
```

Axis vectors

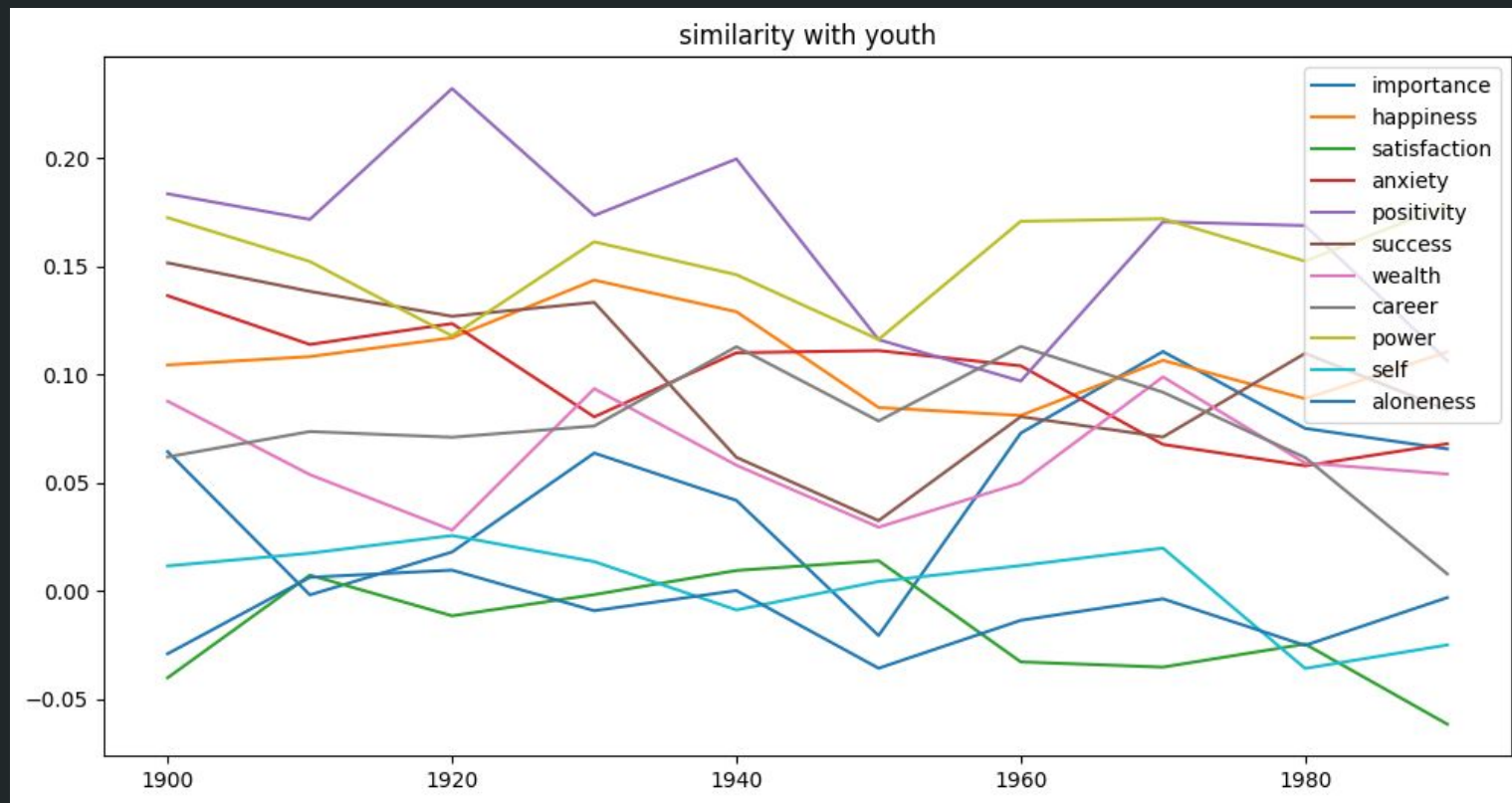
```
axes = ['purpose', 'career', 'happiness', 'satisfaction', 'anxiety', 'wealth']
purpose_words_pos = ['purpose', 'significance', 'significant', 'importance', 'meaningful', 'valuable', 'worthwhile', 'significant', 'important', 'valuable', 'worthwhile']
purpose_words_neg = ['purposeless', 'insignificance', 'insignificant', 'unimportant', 'meaningless', 'valueless', 'worthless', 'insignificant', 'unimportant', 'valueless', 'worthless']
career_words_pos = ['employed', 'skilled', 'occupied', 'professional', 'career-oriented', 'ambitious', 'productive', 'productive', 'productive', 'productive']
career_words_neg = ['unemployed', 'unskilled', 'unoccupied', 'idle', 'unproductive', 'unproductive', 'unproductive', 'unproductive']
happiness_words_pos = ['happy', 'joyful', 'glad', 'pleased', 'content', 'satisfied', 'satisfied', 'satisfied', 'satisfied']
happiness_words_neg = ['sad', 'sorrowful', 'grieving', 'grieving', 'grieving', 'grieving', 'grieving', 'grieving']
satisfaction_words_pos = ['satisfied', 'content', 'satisfaction', 'satisfaction', 'satisfaction', 'satisfaction', 'satisfaction', 'satisfaction']
satisfaction_words_neg = ['dissatisfied', 'discontent', 'dissatisfaction', 'dissatisfaction', 'dissatisfaction', 'dissatisfaction', 'dissatisfaction', 'dissatisfaction']
anxiety_words_pos = ['calm', 'unafraid', 'fearless', 'fearless', 'fearless', 'fearless', 'fearless', 'fearless']
anxiety_words_neg = ['anxious', 'afraid', 'fearful', 'fearful', 'fearful', 'fearful', 'fearful', 'fearful']
positivity_words_pos = ['hopeful', 'positive', 'optimistic', 'optimistic', 'optimistic', 'optimistic', 'optimistic', 'optimistic']
positivity_words_neg = ['hopeless', 'negative', 'pessimistic', 'pessimistic', 'pessimistic', 'pessimistic', 'pessimistic', 'pessimistic']
success_words_pos = ['successful', 'thriving', 'thriving', 'thriving', 'thriving', 'thriving', 'thriving', 'thriving']
success_words_neg = ['failing', 'struggling', 'struggling', 'struggling', 'struggling', 'struggling', 'struggling', 'struggling']
power_words_pos = ['powerful', 'robust', 'robust', 'robust', 'robust', 'robust', 'robust', 'robust']
power_words_neg = ['powerless', 'fragile', 'fragile', 'fragile', 'fragile', 'fragile', 'fragile', 'fragile']
age_words_pos = ['young', 'youthful', 'youthful', 'youthful', 'youthful', 'youthful', 'youthful', 'youthful']
age_words_neg = ['old', 'decrepit', 'decrepit', 'decrepit', 'decrepit', 'decrepit', 'decrepit', 'decrepit']
self_words_pos = ['self', 'myself', 'myself', 'myself', 'myself', 'myself', 'myself', 'myself']
self_words_neg = ['family', 'friends', 'friends', 'friends', 'friends', 'friends', 'friends', 'friends']
alone_words_pos = ['alone', 'lonely', 'lonely', 'lonely', 'lonely', 'lonely', 'lonely', 'lonely']
alone_words_neg = ['together', 'connected', 'connected', 'connected', 'connected', 'connected', 'connected', 'connected']
wealth_words_pos = ['wealthy', 'rich', 'rich', 'rich', 'rich', 'rich', 'rich', 'rich']
wealth_words_neg = ['poor', 'impoverished', 'impoverished', 'impoverished', 'impoverished', 'impoverished', 'impoverished', 'impoverished']
```



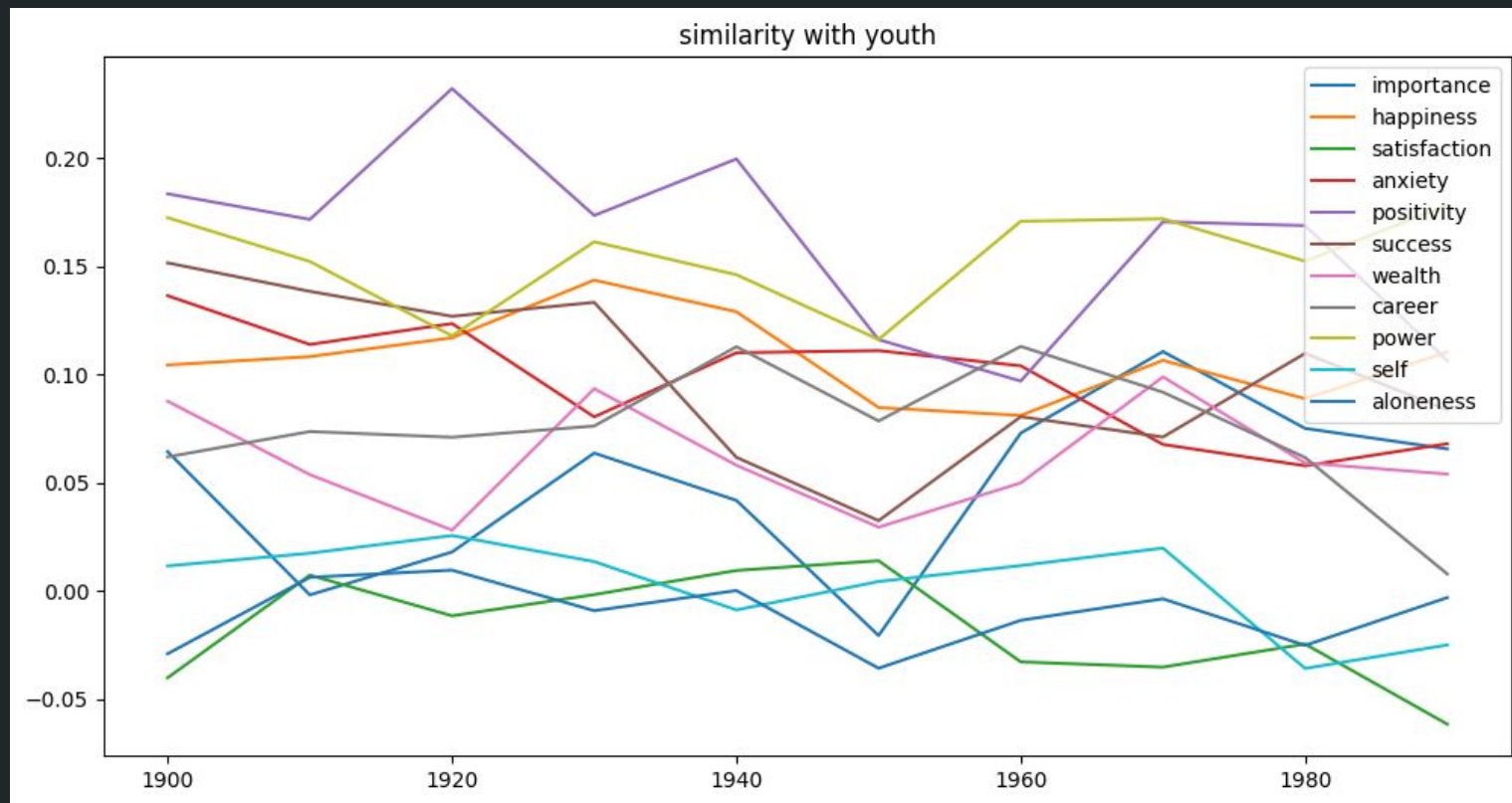
	vecs	purpose_axis	career_axis	happiness_axis	satisfaction_axis	anxiety_axis	wealth_axis
1900	[[[-0.15527616, -0.044475354, 0.24905019, 0.069...	[0.27155071, -0.05503708, -0.11242683, 0.47478...	[-0.046234615, -0.25147384, 0.20027569, -0.060...	[0.26558653, -0.0071866014, -0.042752195, -0.1...	[-0.15267624, -0.29770717, -0.28999034, -0.073...		
1910	[[[0.13818951, 0.08336567, 0.43356463, 0.008813...	[-0.09118675, 0.10703693, 0.21998593, 0.029964...	[0.18839166, -0.13435283, -0.15498304, 0.25415...	[-0.04038347, -0.29692504, -0.12958808, 0.2171...	[-0.5792359, -0.42665628, -0.16115755, -0.3567...		
1920	[[[0.08622483, 0.059710536, 0.03332708, 0.30869...	[-0.5595601, 0.016212419, 0.11945161, -0.00287...	[0.051258057, 0.0387393, 0.17393678, 0.1084501...	[-0.0063034794, 0.15272544, -0.36101255, -0.08...	[-0.15660228, 0.040637117, 0.13808829, 0.17992...		
1930	[[[0.18046615, 0.025464488, 0.057560947, -0.195...	[-0.23274247, 0.0610854, -0.4943271, -0.062550...	[0.10255523, -0.011350911, 0.009990338, 0.1768...	[0.06374652, 0.051403016, -0.2491311, 0.150112...	[0.18366778, 0.115343146, 0.2480477, -0.079701...		
1940	[[[-0.24909955, 0.15136285, 0.016276099, -0.289...				[-0.087414175, -0.0725108, 0.29694104, 0.75403...		
1950	[[[0.04905158, 0.21111171, 0.059865676, -0.3999...				[0.5251739, -0.11746201, -0.047229905, -0.6627...		
1960	[[[-0.10262189, -0.20712672, -0.25004148, -0.36...	0.35192...	-0.25...	-0.5432...	0.27791998, -0.1405...		

$$\sum_p \frac{|P| \vec{p}_1 - \vec{p}_2}{|P|}$$

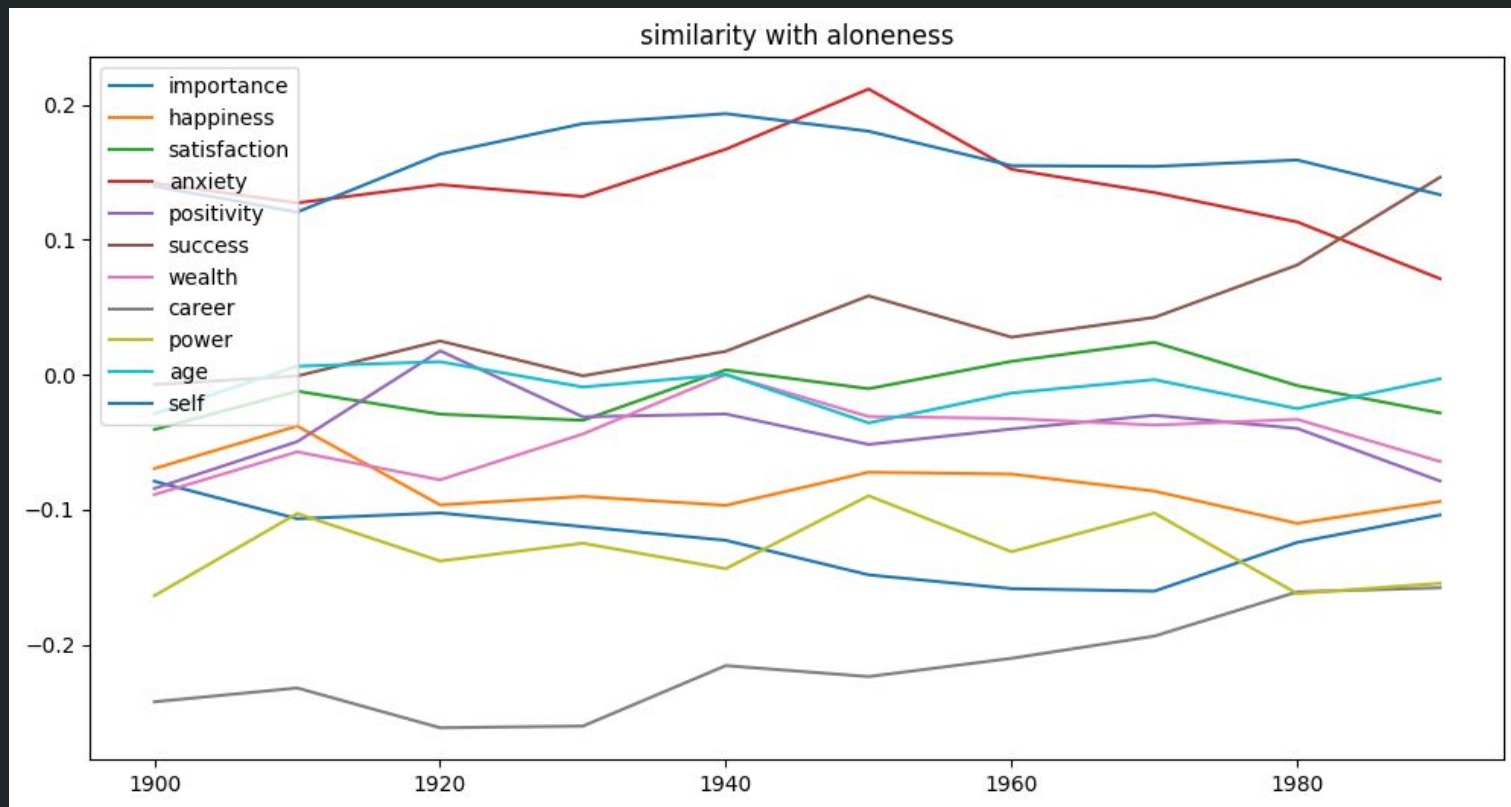
Axis similarity: youth



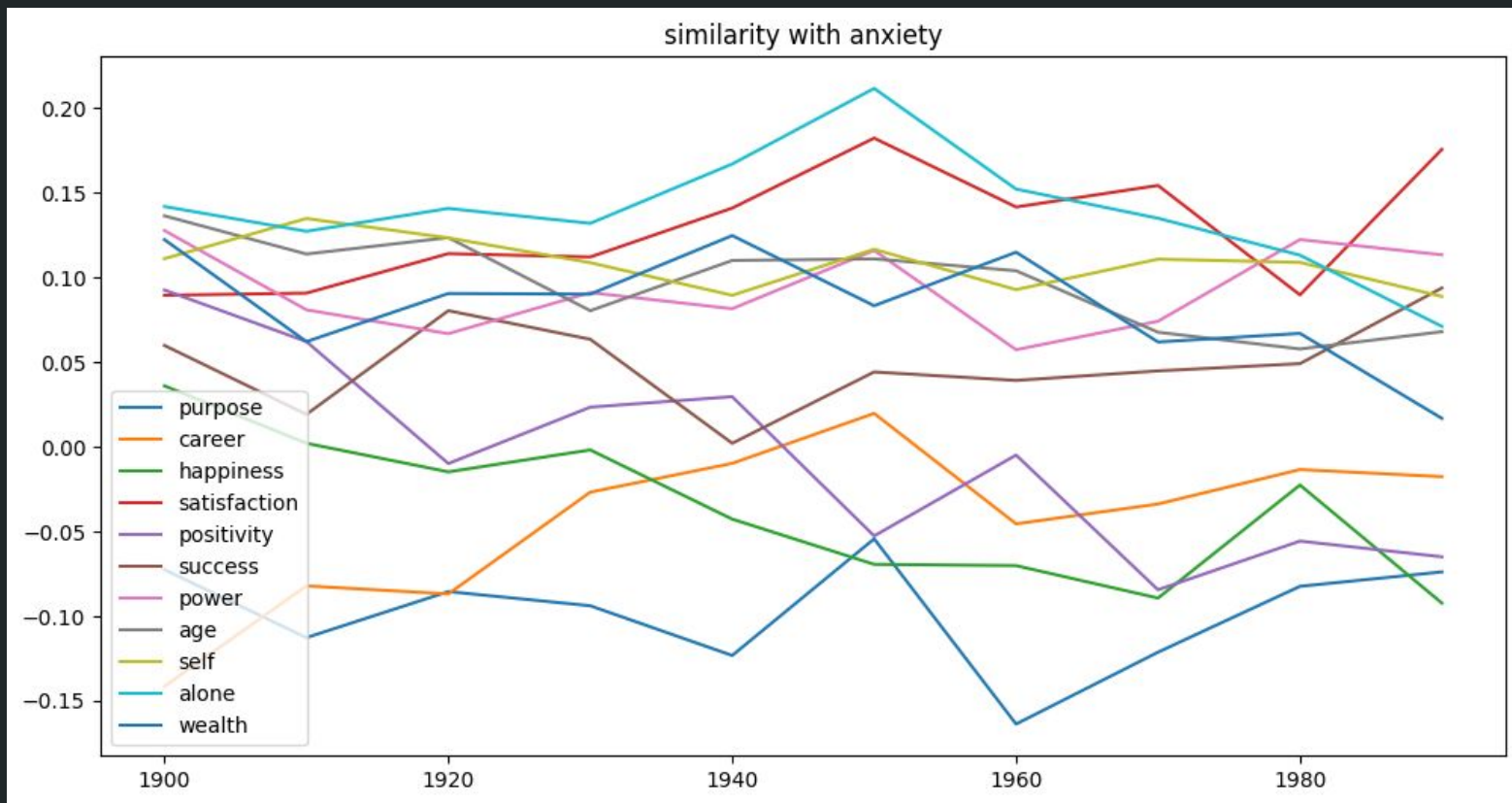
Axis similarity: youth



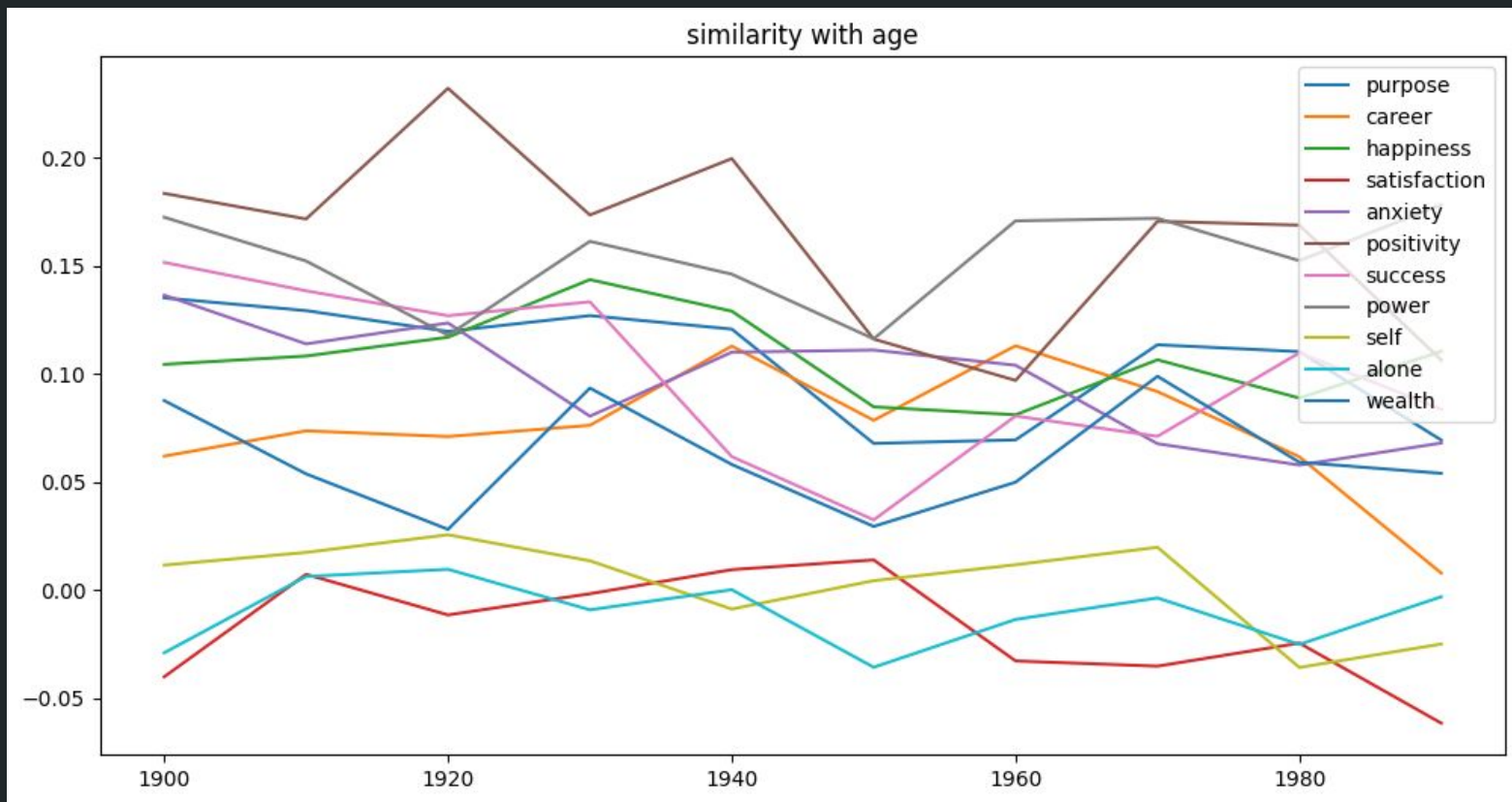
Axis similarity: aloneness



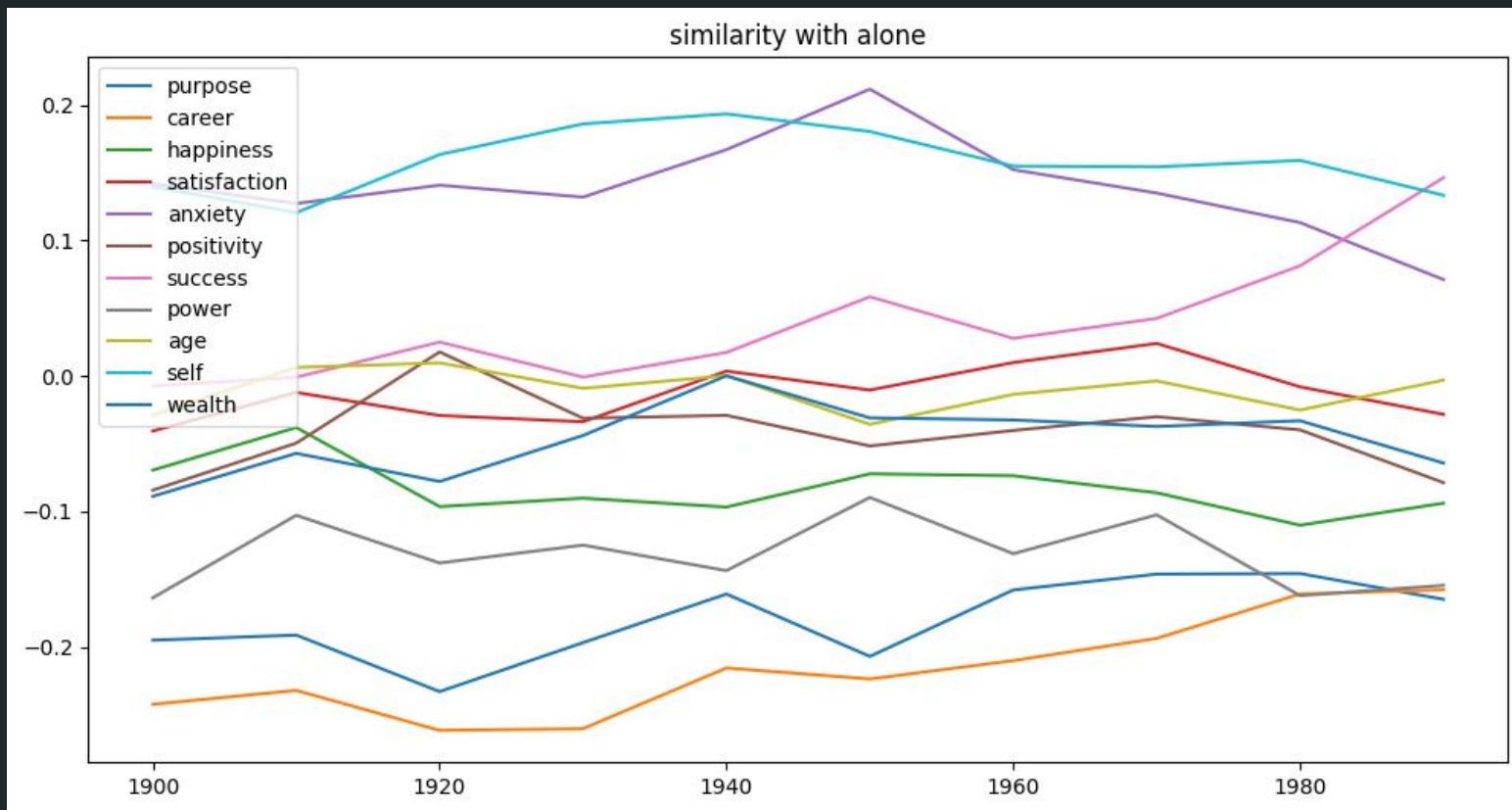
Axis similarity: anxiety



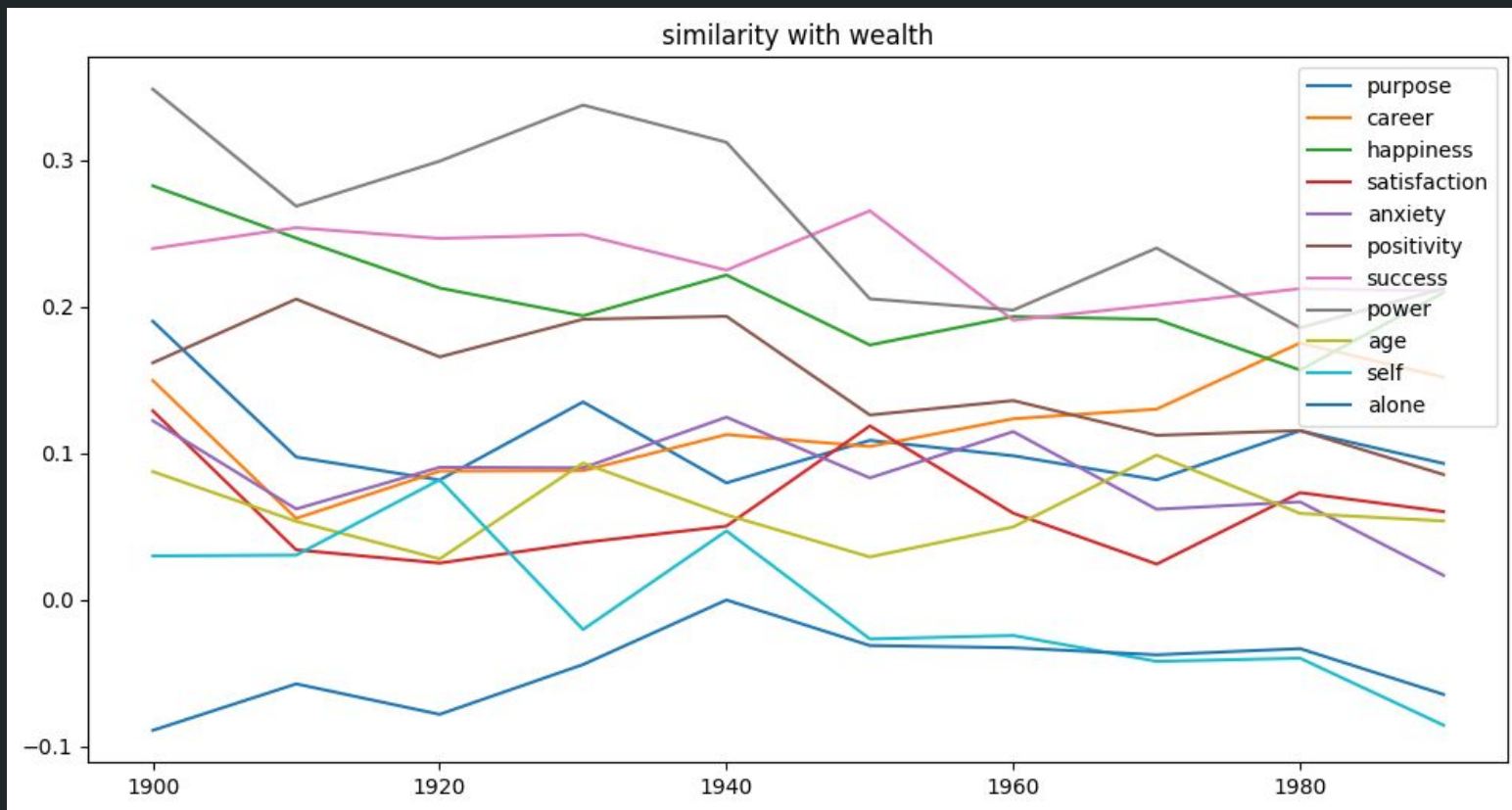
Axis similarity: age



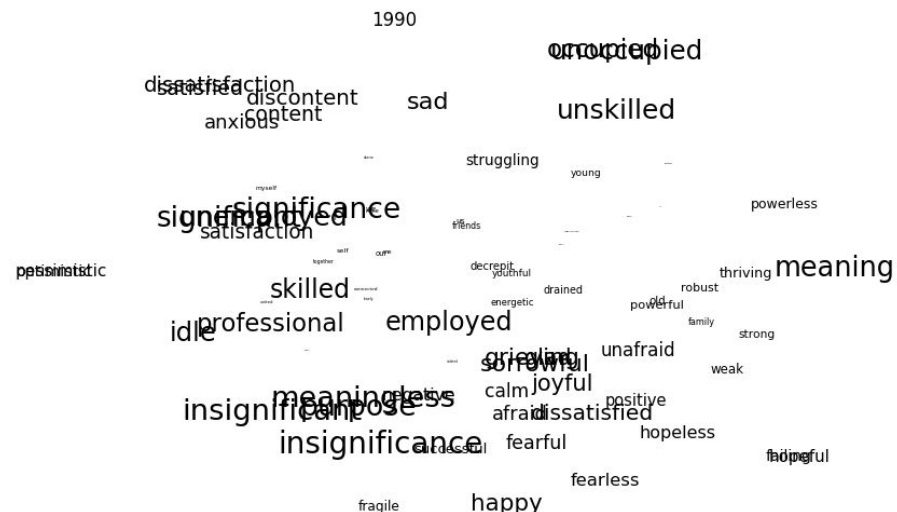
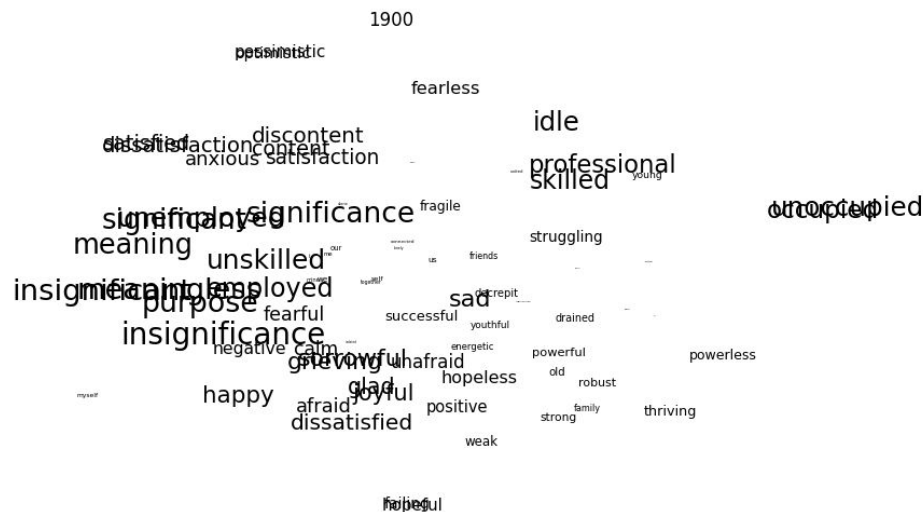
Axis similarity: aloneness



Axis similarity: wealth



Linguistic change



Next steps

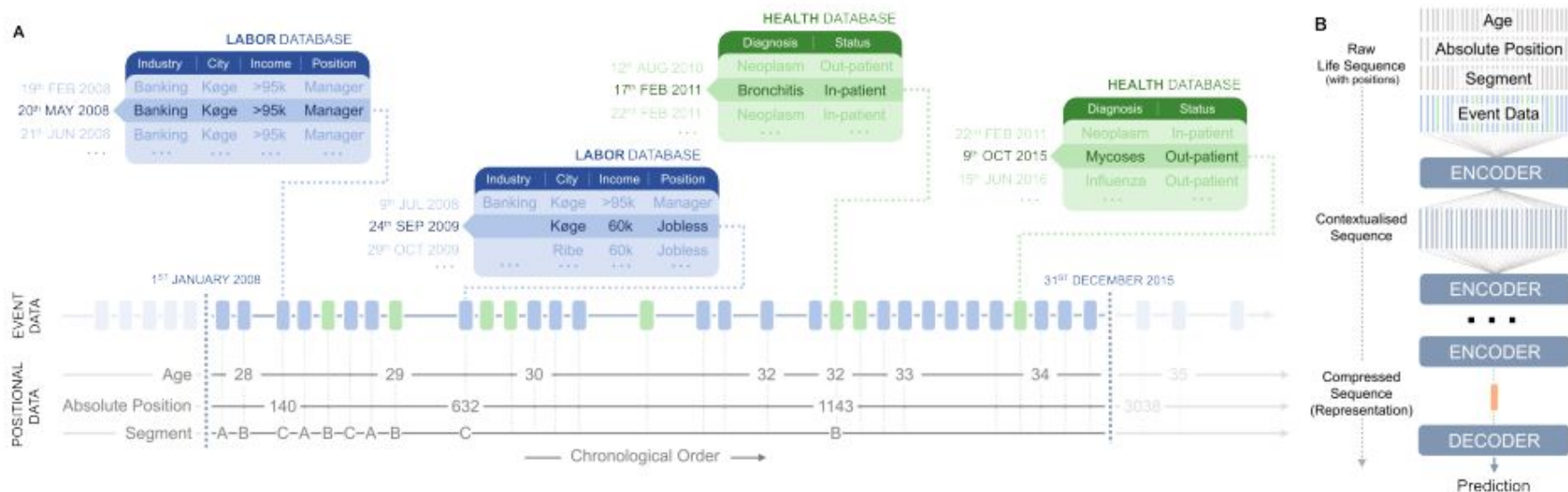
-Historical thesaurus

-Word projections

-Doc2Vec

-Discourse atoms

-Life2Vec



Next steps

-Historical thesaurus

-Word projections

-Doc2Vec

-Discourse atoms

-Life2Vec

