# English Lexical Semantic Evolution across Time

Sangpil Youm, He Zhou

# Introduction

Human language is undergoing a constant evolution driven by the ongoing change in the real world: new concept, language contact, etc.

Language change: phonetic&phonological, morphological, syntactic, semantic

Lexical semantic change: diachronic evolution of lexicon usage, which is easier to change compared with other components of language.

"Every word has its own history."

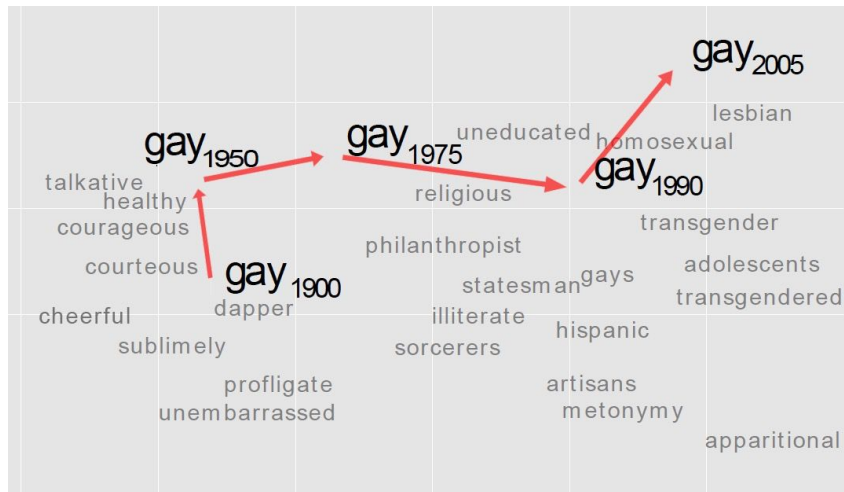Theoretically, lexical semantic change mainly includes 3 ways:

1. Broadening   2. Narrowing    3. Shift

THE GAY GENIUS

*The Life and Times of Su Tungpo*

BY

LIN YUTANG



Questions:

Can we figure out what words have changed in meaning through time?

If a word has changed its meaning, when and how?

Can we analyze the diachronic change from network science perspectives?

# Related Work

Adam Jatwot et al. 2014.

*A Framework for Analyzing Semantic Change of Words across Time*

A visual analytics framework for discovering and visualizing lexical change at 3 different levels -- individual words, word pairs, sentiment orientation

Data: Google Book 5-gram, COHA

Word Representations:

1) Normal Word Representation
2) Positional Word Representation
3) Latent Semantic Analysis based Representation

Vivek Kulkarni et al, 2015

*Statistically Significant Detection of Linguistic Change*

Word Evolution Modeling: frequency, syntactic, distributional

Statistical Soundness: use change point detection in time series to assign significance of change scores to each word

Cross-Domain Analysis: books, tweets and online reviews

Data: Google Book N-gram, Tweets, Amazon reviews

William Hamilton et al. 2016

*Diachronic Word Embeddings Reveal Statistical Law of Semantic Change*

Developed a robust methodology for quantifying semantic change by evaluating word embeddings (PPMI, SVD, word2vec) against known historical changes.

Data: 6 historical corpora spanning 4 languages and 2 centuries.

Two quantitative laws of semantic change:

1) the law of conformity- frequent words change more slowly

2) the law of innovation - polysemous words change more quickly

# Data

❖ English Literature from Gutenberg Project

➢ 17th Century 200 books (5,741,155)

➢ 19th Century 200 books (141,422,786)

➢ 20th Century 200 books (till around 1970s) (86,618,416)

❖ 21st: Google News Embeddings

Preprocessing:

❖ Convert all capital letters to lower-case

❖ Remove all punctuations

❖ Remove stopwords with NLTK

# Method

- ❖ Word-embeddings (word2vec) for words in each time period

  - Default setting, bigram detection

  - Analogy test

  - ex) "King - Man + Woman = Queen"

- ❖ Measure distance between words with cosine similarity

  - Find out similar meaning words

  - Figure out semantic change among different time periods

# Method

❖ Build a word list in which all words show up in both periods

- 17th vs 20th, 19th vs 20th

- Make vocab containing words appeared in both time perionds

❖ Randomly sample 1,000 words from the word list

# Method

Network Analysis

❖ Generate graph with selected words and similar words based on cosine similarity

❖ Calculate Jaccard Coefficient for each word in different time periods

  - 17th vs 20th, 19th vs 20th

❖ Plot the correlation between Jaccard coefficient and frequency

$$J(A, B) = \frac{|A \cap B|}{|A \cup B|}$$

# Method

Network Analysis

- Cosine similarity threshold 0.8

# Result     pretty (ADJ → ADV)



‘pretty’ - 17th century                         ‘pretty’ - 19th century

'pretty' - 20th century

'pretty' - 21st century

# check (VERB → NOUN)



'check' - 17th century

'check' - 19th century

'check' - 20th century

'check' - 21st century

# nice (negative → positive)



'nice' - 17th century

'nice' - 19th century

'nice' - 20th century

'nice' - 21st century

# gay



'gay' - 17th century
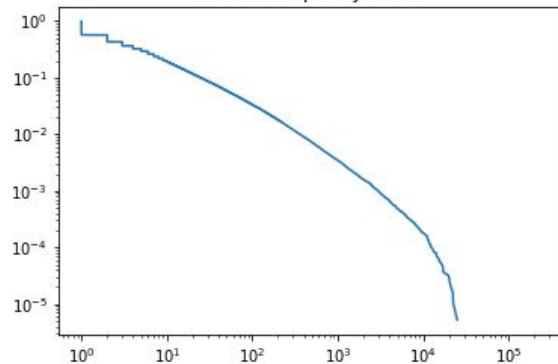
'gay' - 19th century

'gay' - 20th century

'gay' - 21st century

# Frequency

17and20Century_all
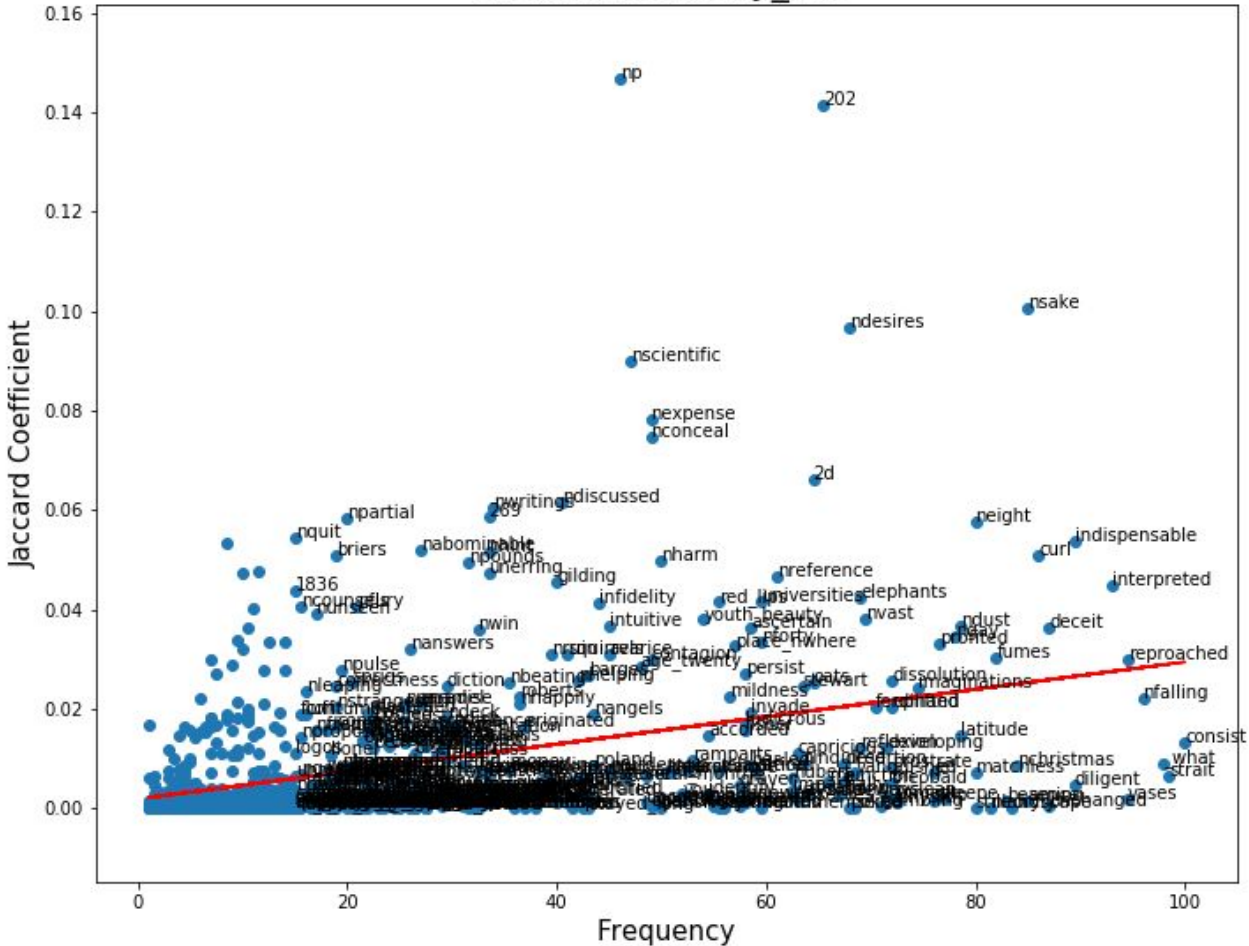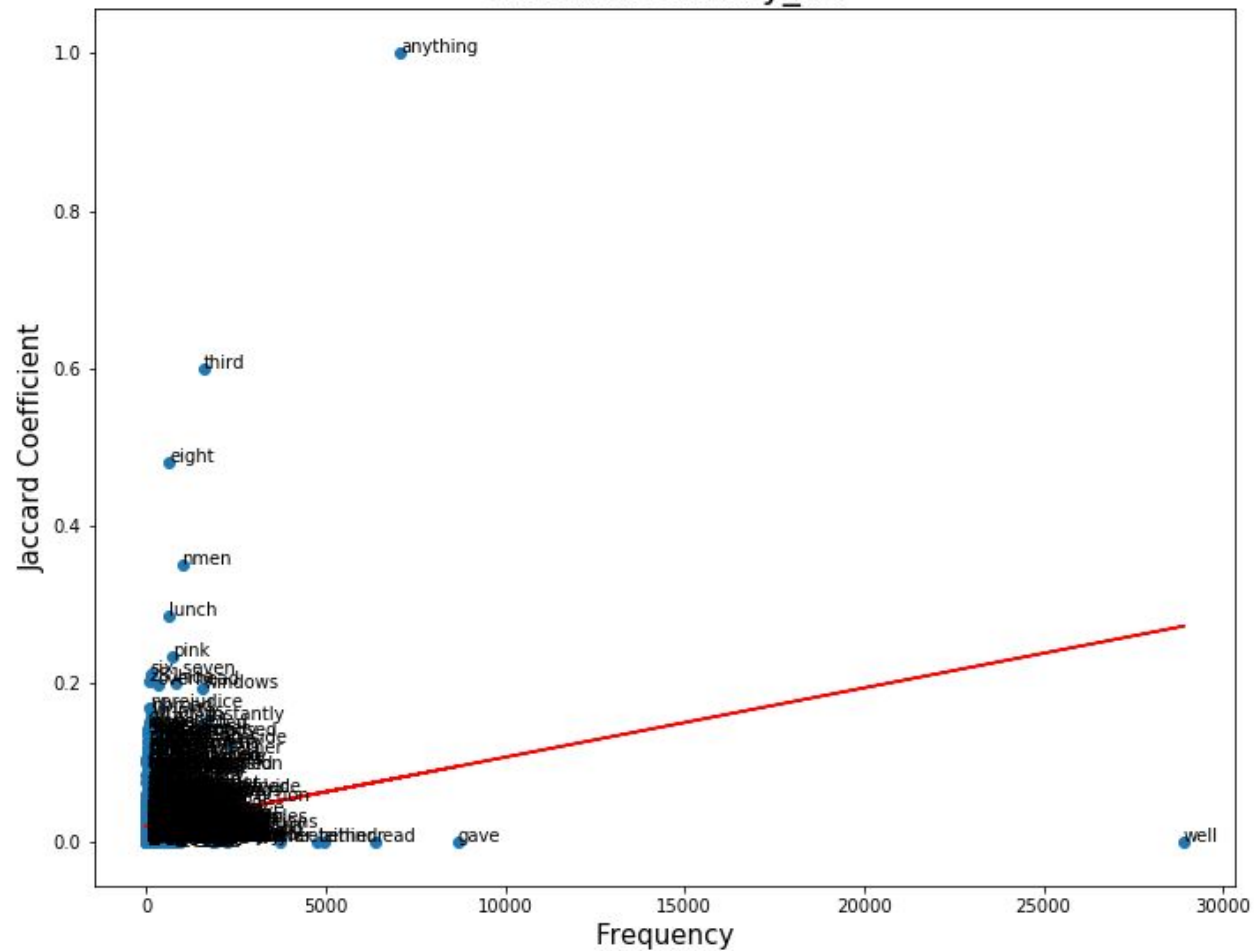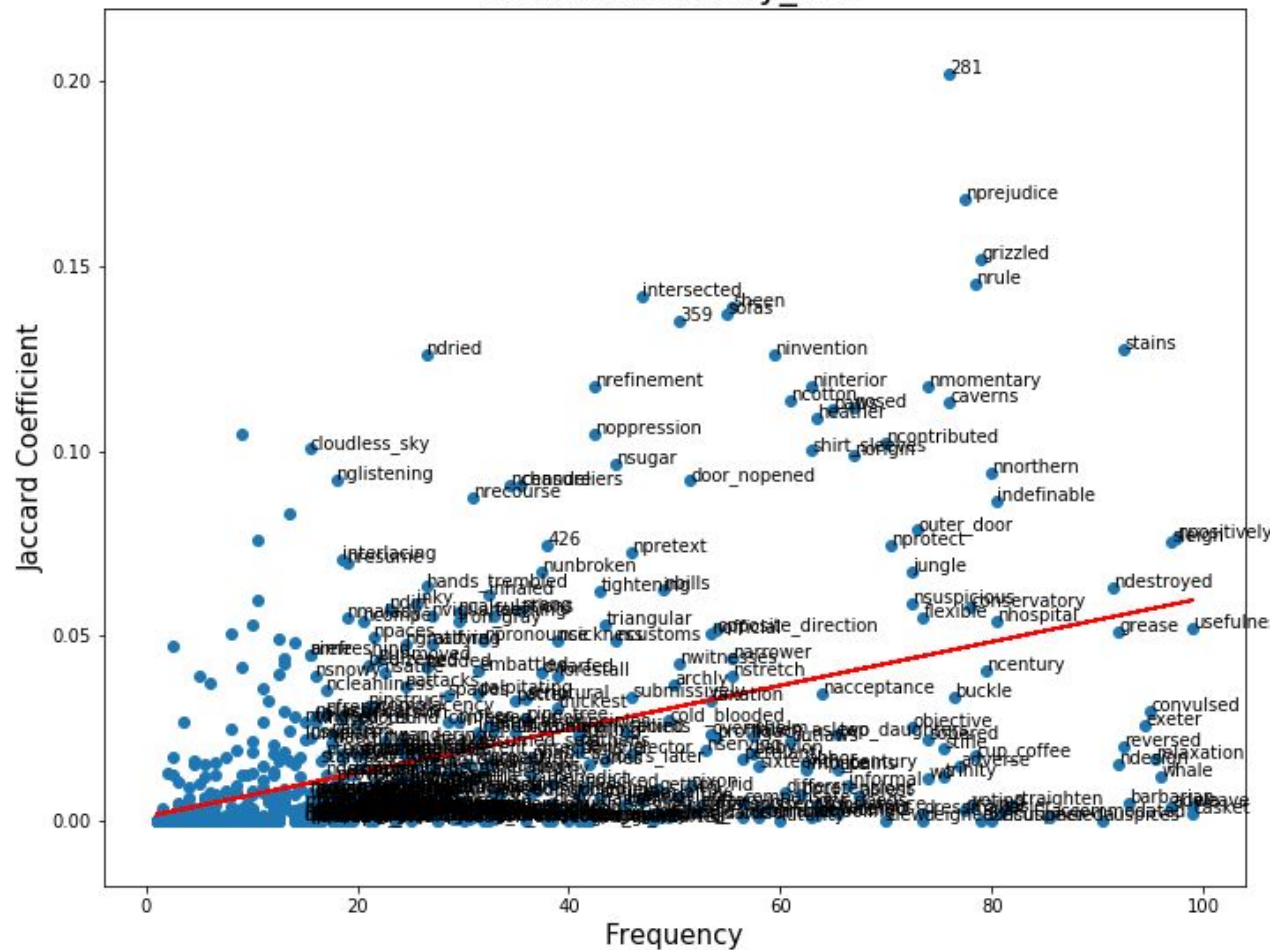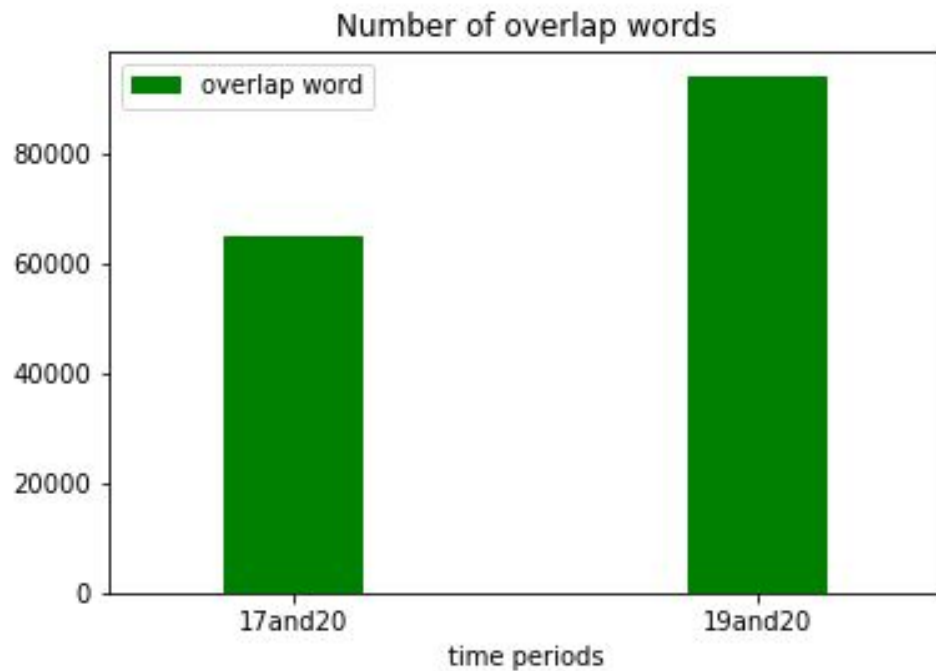
17and20Century_100
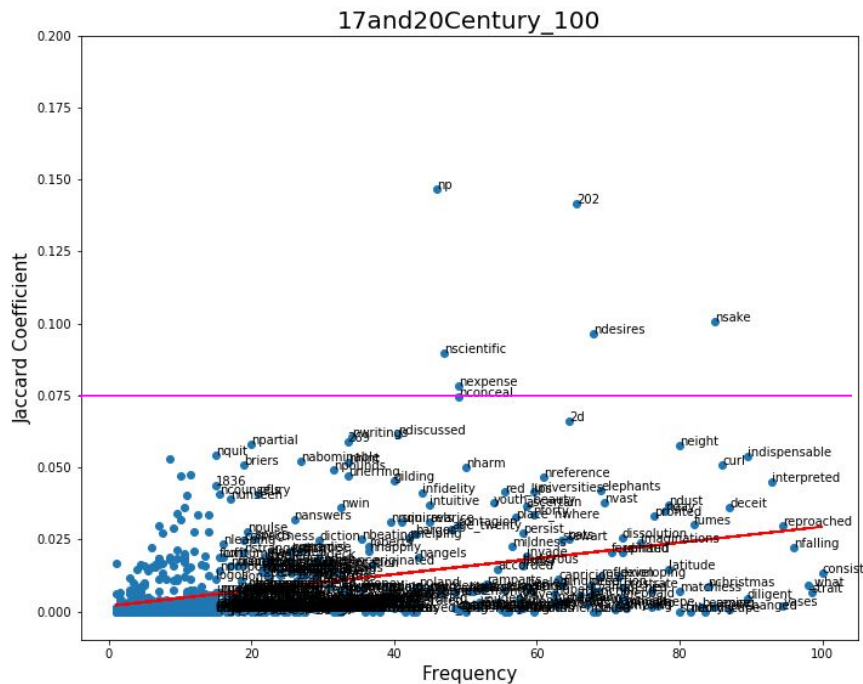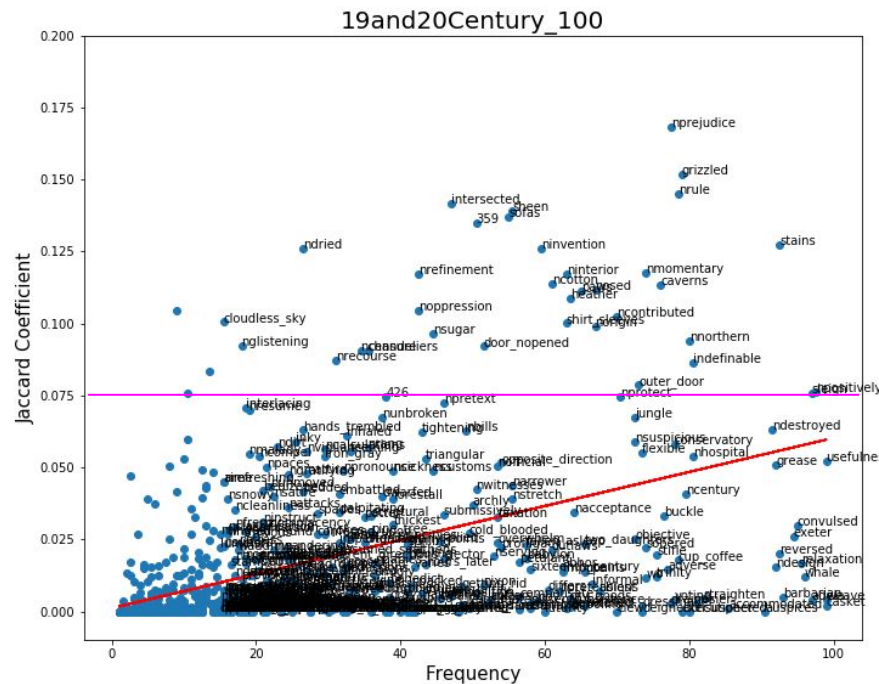
19and20Century_all

19and20Century_100

# Number of overlap words

# Words have changed overtime



863 words

868 words

# Conclusion & Future Work

Conclusion:

1. Semantic change can be detected within the limited amount of data.

   - Semantic change over time

2. Relation between semantic change and frequency can be shown with network analysis.

   - Higher frequency words have less semantic change
   - Words have been changed over time

# Conclusion & Future Work

Future Work:

1. Data is never enough

2. Better way to represent word meaning, e.g.: transfer learning on Google Book N-gram

3. More measures

# Thank you

# &

# Any feedback is welcomed!