

社会学におけるテキスト分析 とNDL全文データの可能性

#NDL全文使ってみた ～「次世代デジタルライブラリー」 & 「NDL Ngram Viewer」

東京大学

瀧川裕貴

全文テキストに関わる過去の自分自身の研究

1. 「国会会議録」全文テキストを用いた道徳社会学の研究
2. 『社会学評論』全文テキストを用いた戦後社会学史の研究

「国会会議録」全文テキストを用いた道徳社会学の研究

国会会議録検索システム

第1回国会（昭和22年5月）からの本会議・委員会の会議録を、テキスト又は画像で閲覧できます。



- データ：1947-2016年の衆参両院，本会議と各委員会会議録のほぼすべて．発言数8,328,057
- 目的：日本の政治家による言論の感情および道徳原理の使い方の立場（与党・野党）による違いや変化を探る

「国会会議録」全文テキストを用いた道徳社会学の研究

国会会議録検索システム

第1回国会（昭和22年5月）からの本会議・委員会の会議録を、テキスト又は画像で閲覧できます。

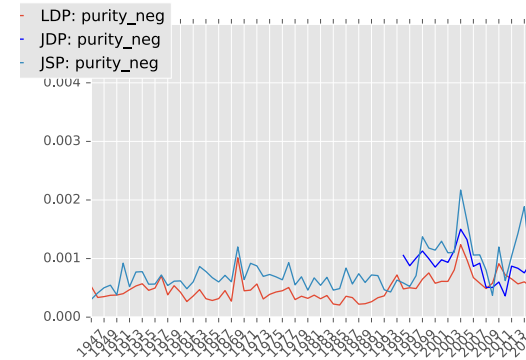
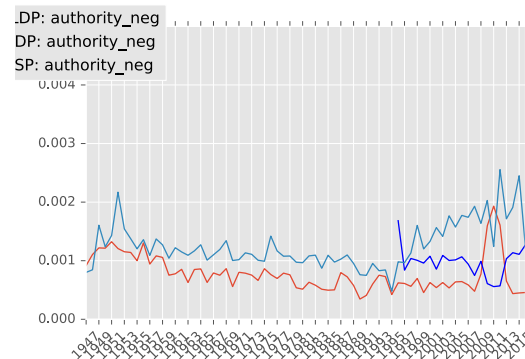
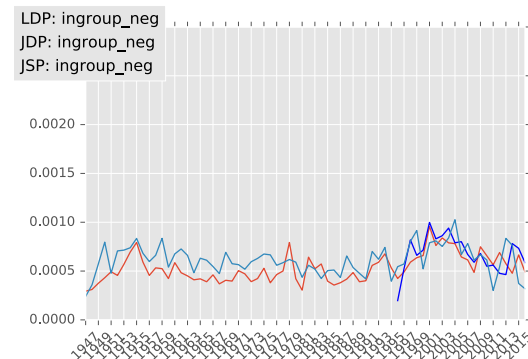
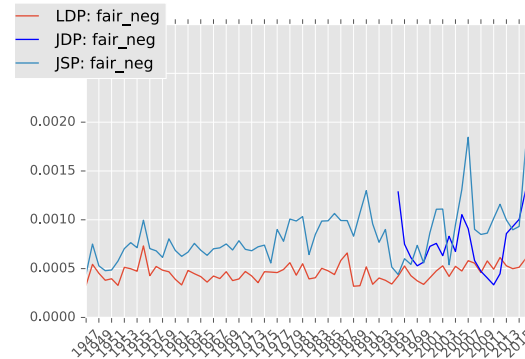
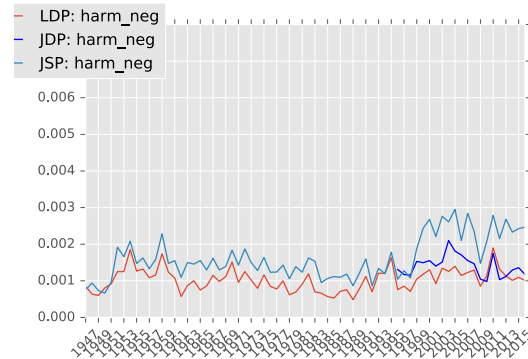


- The moral dictionary and its translated version

- 方法：辞書を用いた頻度分析
 - 日本語評価極性辞書（東山ほか2008）
 - J.Haidtの道徳基盤辞書の独自翻訳版

安心* 01	safe*
平和* 01	peace*
思いやり* 01	compassion*
共感* 01	empathy*
同情* 01	sympathy*
ケア01	care
思いやり01	caring
保護* 01	protect*
保護01	shield
避難所01	shelter
友好01	amity
安全* 01	secure*
利益* 01	benefit*
防衛* 01	defen*
守る* 01	guard*
保存する01 07 0	preserve

「国会会議録」全文テキストを用いた道徳社会学の研究



- 日本の主要政党による道徳原理の使い方とその変化

『社会学評論』全文テキストを用いた戦後社会学史の研究

- データ：『社会学評論』に掲載された学術論文. 1号（1951年）から65号（2015年）までの全論文（書評等を除く）1646本
- 目的：戦後日本社会学が扱った研究テーマや研究領域・研究関心の変遷を明らかにする



『社会学評論』全文テキストを用いた戦
後社会学史の研究

- 方法：トピックモデルという機械学習の方法を用いて，全文データから自動的にトピックを抽出





社会学における全文テキスト研究

- データ：The Google Ngram corpus（全文ではなく，5-gram[5つの連続する単語の集まり]を使用）
- 目的：社会階級の文化的意味の変遷を探る

Check for updates



The Geometry of Culture: Analyzing the Meanings of Class through Word Embeddings

American Sociological Review
2019, Vol. 84(5) 905–949
© American Sociological
Association 2019
DOI: 10.1177/0003122419877135
journals.sagepub.com/home/asr



Austin C. Kozlowski,^a  Matt Taddy,^b
and James A. Evans^{a,c} 

Abstract

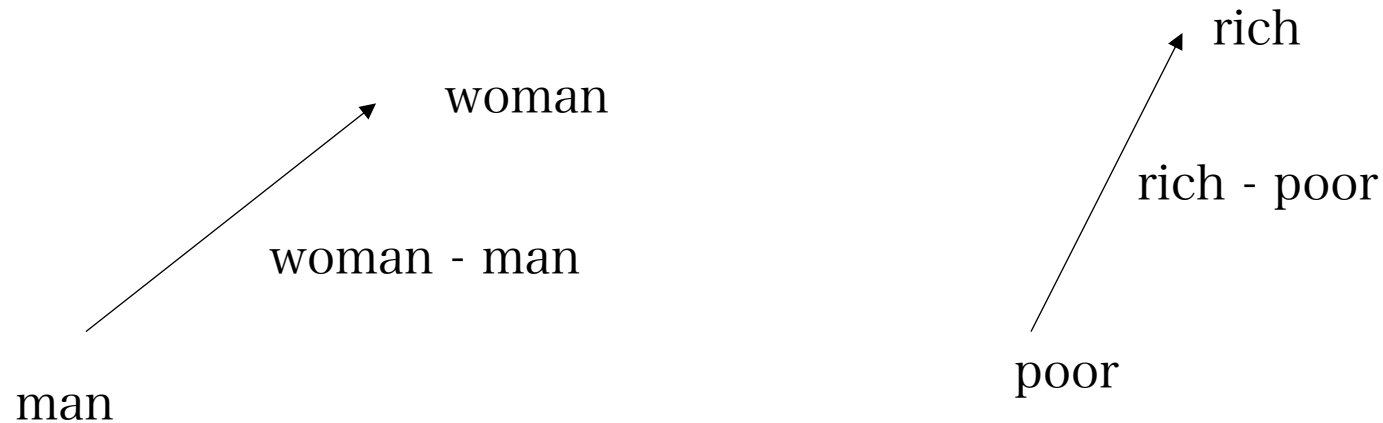
We argue word embedding models are a useful tool for the study of culture using a historical analysis of shared understandings of social class as an empirical case. Word embeddings represent semantic relations between words as relationships between vectors in a high-dimensional space, specifying a relational model of meaning consistent with contemporary theories of culture. Dimensions induced by word differences (*rich* – *poor*) in these spaces correspond to dimensions of cultural meaning, and the projection of words onto these dimensions reflects widely shared associations, which we validate with surveys. Analyzing text from millions of books published over 100 years, we show that the markers of class continuously shifted amidst the economic transformations of the twentieth century, yet the basic cultural dimensions of class remained remarkably stable. The notable exception is education, which became tightly linked to affluence independent of its association with cultivated taste.

Keywords

word embeddings, *word2vec*, culture, computational sociology, methodology, text analysis, content analysis

社会学における全文テキスト研究

- 方法：単語埋め込みモデルの応用による文化的意味次元の抽出
 1. 単語の持つ超高次元情報を低次元ベクトルに効率的に縮約（分散意味表現）
 2. さらに対義語の単語ベクトルを組み合わせることで特定の文化的意味次元を表す解釈可能なベクトルを構成

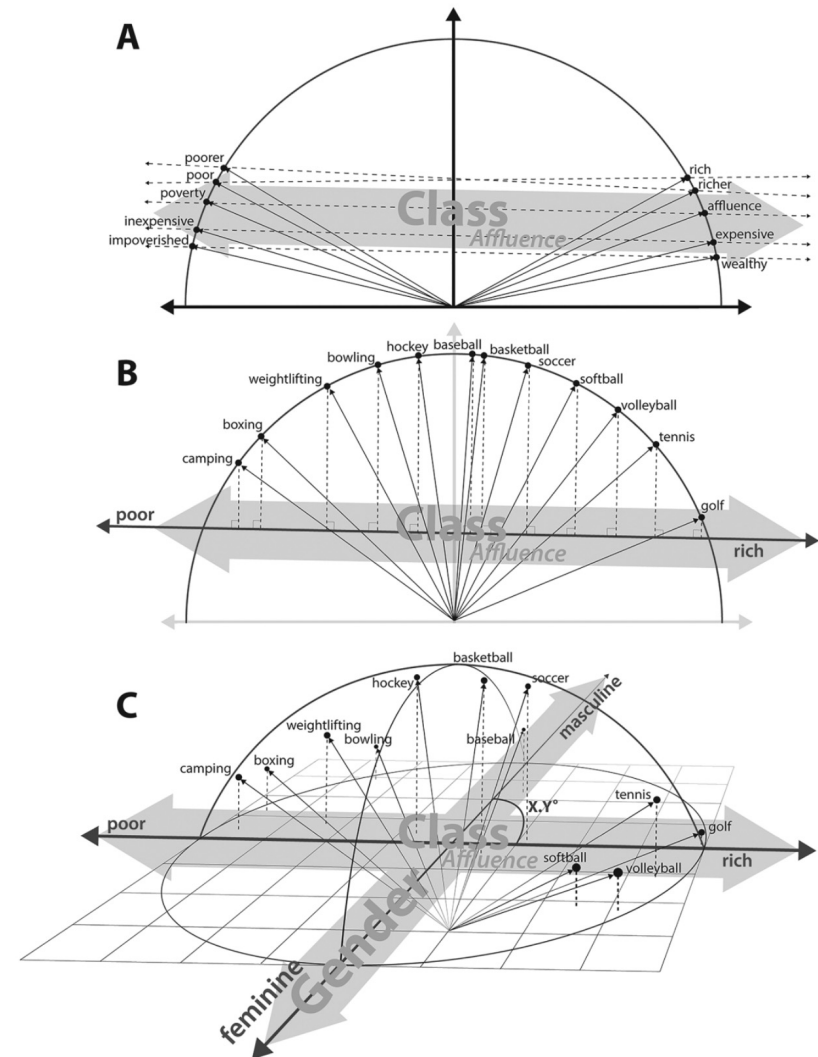


社会学における全文テキスト研究

A: 経済的豊かさの軸の構成

B: 経済的豊かさという意味次元
上でのスポーツの位置

C: 経済的豊かさおよびジェン
ダーという意味次元上でのスポー
ツの位置



NDL全文テキストを用いた研究計画

- 単語埋め込みモデルを用いた社会的実践・思想・概念の文化的意味づけの構造と歴史的変遷の分析