

社会学理論と言語処理技術の接続： ブルデュー理論に基づく言語処理「界」の分析を事例に

高橋祐貴¹ 塚越柚季²

¹フリー ²東京大学大学院人文社会系研究科

biz.yuki.takahashi@gmail.com yuzuki@l.u-tokyo.ac.jp

概要

本研究は、人文学の系譜に位置付く社会学理論の実証の文脈において、いかに自然言語処理が応用可能かを示すべく、言語処理の研究者集団を対象にブルデュー流の「界」の分析を行い、その分化の構造を明らかにする。ACL Anthology と researchmap を用いて構築した有力な言語処理研究者 578 名のデータセットを用いて、多重対応分析により「界」の空間を構築したところ、資本総量の差は見られたものの、解釈可能な資本構成の交差配列構造は認められなかった。BERTopic により析出した研究分野との関連もそれほど明確とは言えず、理工系の分野においてはブルデュー理論が有効とは言えないことが示唆された。

1 はじめに

哲学の流れを汲む社会理論に端を発しながらも、実証科学としての側面を強調する学問として発展してきた社会学では、言葉による定式化されていない理論を、定量・定性それぞれの分析により、経験的に検証する試みが積み重ねられてきた。その中で、特に定量的な実証分析において近年関心が高まっているのが、自然言語処理の技術である。その範囲は、LDA を始めとするトピックモデルの応用[1,2]から、単語埋め込みベクトルによる単語の意味変遷の推定[3]、LLM を用いた政治的反応の予測実験[4]など、非常に幅広い。

社会学が自然言語処理を応用するに当たって、特に強く意識されるのは、分析により得られた何らかの数値的結果が社会学理論上ではどのような意味を持つのか、その対応を明確化することである。例えば、LDA により得られるトピックの理論的な位置付けを、Alfred Schutz のレリヴァンス理論を参照し、「解釈的レリヴァンスの計量化」として位置付ける、

といった具合である[1]。本研究は、このような社会学における自然言語処理の応用の一例として、2 節にてフランスの社会学者 Pierre Bourdieu の「界」理論における「態度決定」という概念とトピックモデルを対応させた分析枠組みを提示する。そして3 節で実際の分析例として、日本の言語処理「界」の分析を行うためのデータと方法を提示し、4 節にて人文社会系の分野の分析手法として提案されたブルデューの学問の「界」の分析視角が、理工系分野である言語処理の分析でも有効性を持つかを検証する。5 節では、得られた知見を基に、今後の研究の展開を議論する。

2 背景と関連研究

2.1 ブルデューの「界」理論と学問分析

「文化資本」の概念などで知られるブルデューの「界」理論は、元々は作品の創作者の社会関係と関連づけて作品の創作と受容を考察するべく、芸術分析の流れの中で提示された。まず「資本」の変数を用いて多重対応分析を行い、2~3 次元の「界」の空間を析出する。「資本」は、例えば文芸界なら文学賞の受賞歴や部数、科学界なら論文数や学会賞歴など、それぞれの「界」の中で有利な立場に立つために必要な有形・無形の所有物が分析ごとに定められる。続いてこの空間における座標を行為者の「位置」とし、座標軸の意味を解釈することで「界」の構造を読み解く。多くの場合、最も分散を説明する1 軸目に資本の多寡が反映され、2 軸目に資本の構成割合の分化が見られるという解釈になる。そして、「界」における「位置」と、行為者が表明する創作物の内容や政治的態度といった「態度（位置）決定」は関連する、という図式を通じて、文化的生産物と行為者の社会的位置の関係性を読み解くのが、ブルデューの「界」理論の枠組みである[5]。

ブルデュー自身は『ホモ・アカデミクス』[6]において、フランスの人文社会科学界の分析を行っている。ここで、ブルデューは多様な学問的所有物を基準に選別した、界の中で高い位置を占めると想定されるパリの「文学系」正教授のサンプルから、対応分析により文学・人間科学部の権力空間を構築した。ブルデューが見出したのは、「純粋に大学的な権力」と、「多少とも制度化された知的名声」（学問的威信と知識人界での名声という象徴資本の双方を含む）の2種類の権力の対立である。加えてもう一つ主要に見出された区分として、単に大学的・学問的権威の多い年長の教授と、大学内権力としては下位に位置付く若年の教授たちという区分があった。この年齢的な区分は、資本総量の分化に対応するものとして捉えられる。こうしてブルデューは、文学・人間科学部界の中にも、相対的に文化資本を多く持つグループと、経済資本を多く持つグループがあるという対立構造を見出した。

さらにブルデューが踏み込んで検討したのは、文学・人間科学部界における知的生産物の内容や形式と、こうした権力空間との間に見られる、構造的な関連性についてである。文学テキストの正統的注釈を主張した保守派のレイモン・ピカールと現代主義的注釈を主導した革新派のロラン・バルトが、対応分析の図では対照的な位置にいたることが、その一例であった。

2.2 経験的な先行研究

後続の学問の界の研究では、フランスの経済学界[7]、アメリカの社会学界[8]などの分析が行われ、比較的明瞭な資本構成の対立構造が見られること、また界における「位置」と研究分野や方法といった「態度決定」とが関連することが示唆されてきた。特に「態度決定」の分析は、ブルデュー自身は定性的に読解した著作物の内容と、対応分析の結果とを照らし合わせて行っていたのに対し、後続の研究は軸の生成に寄与しないサプリメンタリ変数として内容に対応する変数（例えば研究者個人の研究分野）を投影する方法へと変化していった。ここで、「態度決定」が「創作物の空間」であるとしたブルデューの主張に立ち返り、トピックモデルを用いて論文や書籍の本文から「態度決定」変数を導き出す手法を取ったのが、日本の社会学界を分析した高橋[9]である。ここでは論文誌の本文データに構造トピックモデルを適用し、析出したトピック割合を基

に個々の研究者の主要な研究分野と方法の変数を作成している。ブルデューが定性的に行っていた分析をより大規模に拡張する上で、トピックモデルにより析出したトピックを「態度決定」として見なす手法には応用可能性がある。

一方で後続の先行研究が欠いてきたのが、いわゆる「理系」の学問分野の「界」の分析である。元々芸術創作との類似性から人文社会系の学問に多く適用されてきた「界」理論が、理工系や医薬系といった別の学問分野においても有効性を持つかは検証されてこなかった。そこで本研究は、理工系の一分野と見做せる自然言語処理を対象に、「界」の構造の分析を行い、人文社会系分野と同様の資本構成の対立構造や「位置」と「態度決定」との対応関係が見られるかを検証する。

3 データと方法

3.1 データ

本研究で用いるのは、ACL Anthology と researchmap を基に構築された、有力な言語処理研究者 578 名の「資本」と「態度決定」に関するデータセットである。まず、ACL Anthology より日本人と判別できる氏名に紐づく発表実績を抽出し、その氏名や発表物のタイトルから researchmap の個人ページを特定し、合計 1443 名の研究者の研究実績や役員歴、受賞歴などのデータを収集した。そこからさらに researchmap で「論文」の項目に登録のある研究者にデータを絞り込んだ結果、サンプルサイズは 578 となった。構築されたサンプル集団は、ACL Anthology に何らかの論文を掲載しておりかつ researchmap に業績を管理しているという意味でバイアスがかかっている点には注意されたい。変数構築の詳細は付録に記載している。

3.2 方法

まず、「資本」の変数を用いて界の空間を構築し、各変数の平均的な配置から軸の意味を読み解く。続いて「態度決定」変数として研究分野の変数をサプリメンタリ変数として投影し、研究分野と「界」の構造に何らかの関連が見られるかを分析する。

研究分野変数は、ACL Anthology のアブストラクトのデータに対し BERTopic[10]を適用し、析出したトピックの割合が多い順に 2 分野を、その研究者の主要研究分野と見做して作成した。析出したトピッ

クとその代表語は表1の通りである。

表1 BERTopicによるトピックの解釈

解釈	Representation
言語処理	['translation', 'language', 'models', 'model', 'based', 'method', 'paper', 'data', 'task', 'results']
機械翻訳	['translation', 'machine', 'english', 'nmt', 'parallel', 'japanese', 'sentence', 'corpus', 'sentences', 'paper']
埋め込みモデル	['models', 'word', 'embeddings', 'tasks', 'language', 'model', 'bias', 'method', 'performance', 'embedding']
対話	['dialogue', 'responses', 'user', 'utterances', 'systems', 'human', 'response', 'users', 'task', 'dialog']
日本語処理	['japanese', 'corpus', 'patent', 'simplification', 'word', 'words', 'method', 'english', 'functional', 'expressions']
構文解析	['parsing', 'argument', 'dependency', 'structure', 'event', 'resolution', 'coreference', 'structures', 'discourse', 'parser']
固有表現認識、関係認識	['entity', 'entities', 'relation', 'ner', 'extraction', 'knowledge', 'model', 'relations', 'embeddings', 'task']
要約生成	['summarization', 'summaries', 'generation', 'headline', 'summary', 'document', 'model', 'text', 'information', 'data']
音声認識	['speech', 'recognition', 'asr', 'acoustic', 'language', 'translation', 'sign', 'censrec', 'spoken', 'japanese']
SNS分析	['twitter', 'polarity', '19']
マルチモーダル(画像、動画)	['image', 'video', 'images', 'visual', 'captions', 'multimodal', 'caption', 'videos', 'captioning', 'dataset']
機械翻訳データセット(辞書)	['mt', 'dictionaries', 'translation', 'lrs', 'resources', 'web', 'www', 'machine', 'user', 'language']
質問応答	['question', 'answer', 'questions', 'answering', 'qa', 'rc', 'model', 'reading', 'mrc', 'scoring']
文法修正	['gec', 'correction', 'grammatical', 'error', 'evaluation', 'learners', 'models', 'examples', 'learner', 'manual']
AI人狼	['game', 'agents', 'agent', 'werewolf', 'games', 'llms', 'logics', 'contest', 'circumstances', 'holding']
議事録処理	['minutes', 'assembly', 'local', 'retrieval', 'spoken', 'meeting', 'corpus', 'document', 'task', 'queries']
音声感情分析	['emotion', 'emotional', 'speech', 'dialogue', 'labeling', 'fusion', 'labels', 'modalities', 'corpus', 'learning']
フィードバック生成	['ranking', 'recommendations', 'generation', 'recommendation', 'online', 'task']
金融	['stock', 'prices', 'market', 'price', 'comments', 'numerical', 'weather', 'changes', 'data', 'generating']
実況生成	['commentaries', 'commentary', 'game', 'live', 'sports', 'races', 'domain', 'real', 'world', 'spectators']
LLM	['reasoning', 'logical', 'syllogistic', 'llms', 'cot', 'biases', 'models', 'abilities', 'symbolic', 'step']
トピックモデル	['topic', 'model', 'topics', 'adverbs', 'tree', 'dynamic', 'hierarchical', 'words', 'structured', 'interactive']
物語生成	['story', 'event', 'stories', 'events', 'module', 'sequence', 'generation', 'narrative', 'salience', 'schema']
広告	['appeals', 'concerns']
料理レシピ	['procedural', 'food', 'dish']
数値処理、Numerical Commonsense	['numerals', 'numeral', 'numbers', 'digit', 'ncs', 'quantitative', 'nn', 'arithmetic', 'magnitudes', 'models']
自動運転	['driving', 'channel', 'dialog', 'utterances', 'driver', 'rate', 'utterance', 'human', 'car', 'context']
フランス語(その他言語)	['des', 'les', 'la', 'le', 'une', 'et', 'du', 'dans', 'pour', 'en']

4 分析結果

「資本」の変数を用いて多重対応分析を行なった結果、1軸の説明率が高く、修正分散説明率は2軸で8割を超えた(表2)。そのため、2軸までの結果を主に解釈する。

表2 多重対応分析の分散説明率

次元	固有値	各軸の分散説明率		
		分散説明率	ベンゼクリ補正比率	累積ベンゼクリ補正比率
dim 1	0.19	13.13	77.83	77.83
dim 2	0.07	5.22	8.89	86.72
dim 3	0.06	4.00	4.33	91.04
dim 4	0.05	3.20	2.23	93.27
dim 5	0.04	3.00	1.81	95.08

4.1 軸の解釈

1軸への寄与率が平均以上のモデルを2-1軸上

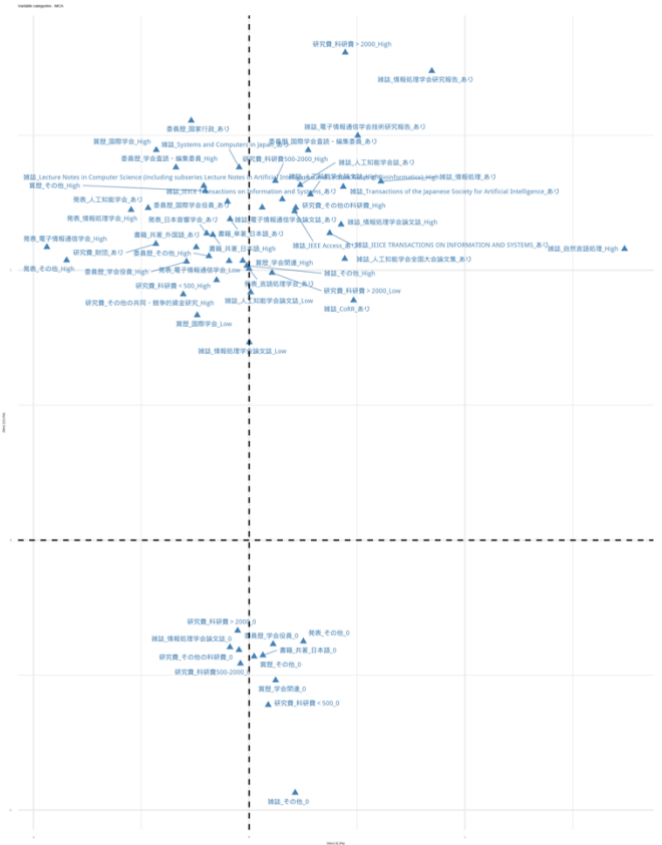


図1 1軸への寄与率が高いモデル(2-1軸)の空間に投影したものが図1である。

これを見ると、1軸方向には資本総量と解釈できそうな分化が反映されていることが分かる。1軸方向に正の座標には雑誌論文や学会発表、研究費の獲得が高いモデルが並び、負の方向には0のモデルが並ぶ。

他方、2軸の構成に平均以上に寄与しているモデルをプロットしたのが図2である。これを見ると、まず寄与率が高い変数が少ない上に、査読のない一部の学会発表が多い群と、査読付きの一部の雑誌論文や海外学会での発表が多い群の分化を示している、程度にしか解釈できず、資本構成の明確な分化が認められるとは言い難いことが分かる。

4.2 分野変数の投影と解釈

2-1軸の「界」の空間上に、「態度決定」としての研究分野変数をサプリメンタリ変数として投影したのが図3である。多くの分野が原点近くに配置され、

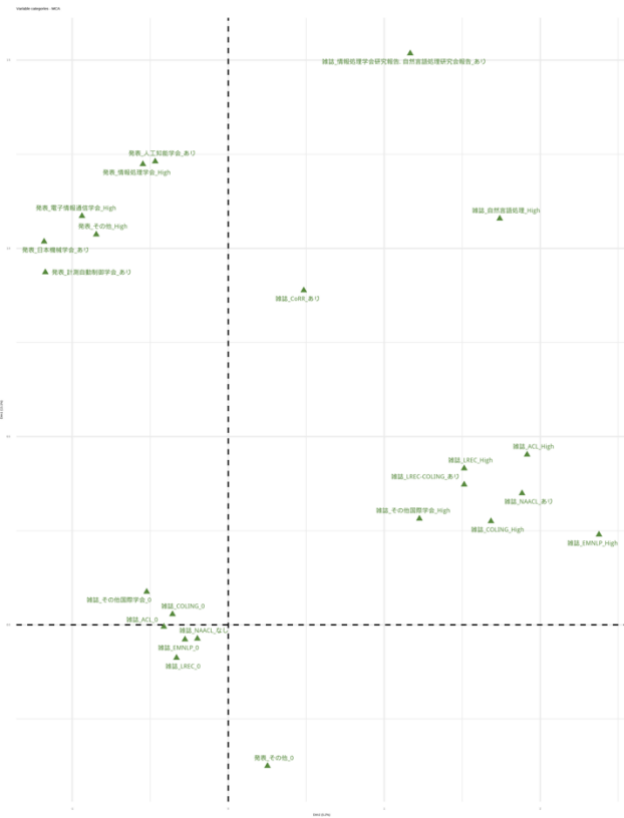


図 2 2軸への寄与率が高いモデル(2-1軸)

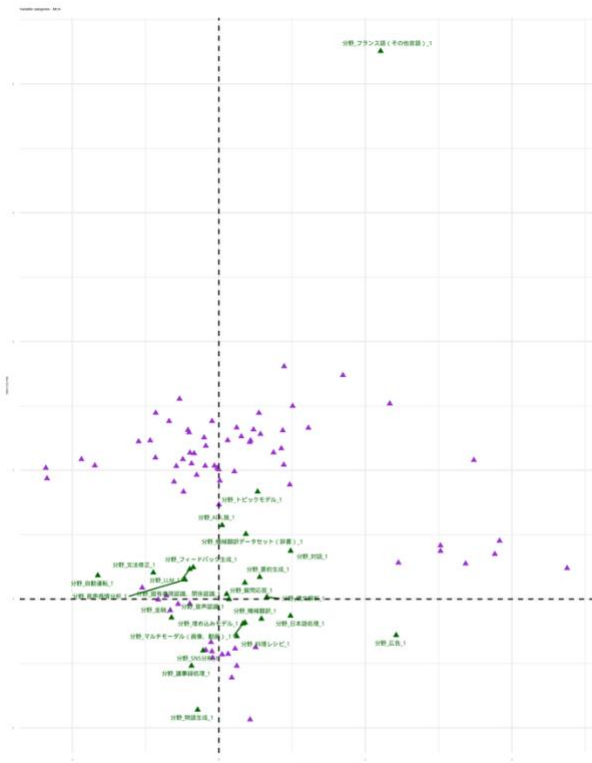


図 3 分野変数の投影(2-1軸)

界における位置と態度決定の対応があると言い難い。

5 おわりに

本研究では、社会学における自然言語処理の応用例を示すべく、ブルデュー理論に基づいた言語処理「界」の分析を試みた。結果として、言語処理分野においては、ブルデュー理論と整合的な界の分化は確認できないことが明らかになった。同様に researchmap のデータを基に社会学者に関する分析を行なった高橋[9]では解釈可能性の高い資本構成の分化が見られたのとは対照的である。これは、人文社会系の分野に見られた内部対立構造が言語処理分野には見られない可能性を示唆する。一方でこの結果自体は、人文社会系分野で想定されるような「資本」変数では分化の構造を捉えられないからとも考えられるし、researchmap をデータソースとして用いることでそもそも分野の全体を捕捉できなかったという問題も考えられる。科研費申請に用いられる researchmap は企業所属の研究者を捕捉できない問題が指摘されている[11]が、多くを大学所属の研究者が占める人文社会系の分野と異なり、企業所属の研究者も中心的な役割を果たしている言語処理においてこれは大きな欠点と言えよう。また、個人が主体の基礎単位となる人文社会系とは異なり、理工系分野では研究室単位で研究主体を捉えた方が適切だとも考えられるだろう。今後ブルデュー派の分析を理工系分野に展開していくに当たっては、人文社会系とのこうした性質の違いを考慮し、より適切なデータソースを検討していく必要がある。

また今後の展開としては、その他の自然言語処理技術とブルデューの理論との接合を模索する道筋も考えられる。Mochihashi [12]が提案した Researcher2Vec や、研究タイトルから研究者の埋め込みを得た Nagao & Katsurai[13]のように「書かれたものの内容から研究者の相対的な位置を把握する」技術はすでに複数存在している。例えば「態度決定」の空間を Researcher2Vec と通じて構築し、その上で関心のあるその他の変数の平均的な配置が分化するかを見る手もあるだろう。自然言語処理の社会学への応用はまだ始まったばかりであり、今後も発展が期待される。

参考文献

1. 小田中悠・中井豊, 2019, 「意味世界の計算社会科学的分析に向けて——社会学におけるトピックモデルの意義の検討」『理論と方法』34(2): 280-95.
2. 瀧川裕貴, 2019, 「戦後日本社会学のトピックダイナミクス」『理論と方法』34(2): 238-61.
3. **Kozlowski, Austin C., Matt Taddy & James A. Evans**, 2019, "The Geometry of Culture: Analyzing the Meanings of Class through Word Embeddings," *American Sociological Review*, 84(5): 905-49.
4. **Kozlowski, Austin, Hyunku Kwon & James Evans**, 2024, "In Silico Sociology: Forecasting COVID-19 Polarization with Large Language Models."
5. 磯直樹, 2020, 『認識と反省性——ピエール・ブルデューの社会学的思考』法政大学出版局.
6. **Bourdieu, Pierre**, 1984, *Homo academicus*, MINUIT. (石崎晴己・東松秀雄訳, 1997, 『ホモ・アカデミクス』藤原書店.)
7. **Lebaron, Frédéric**, 2001, "ECONOMISTS AND THE ECONOMIC ORDER The Field of Economists and the Field of Power in France," *European Societies*, 3(1): 91-110.
8. **Warczak, Tomasz & Stephanie Beyer**, 2021, "The Logic of Knowledge Production: Power Structures and Symbolic Divisions in the Elite Field of American Sociology," *Poetics*, 87: 101531.
9. 高橋祐貴, 2023, 「現代日本の社会学界——ブルデュー派アプローチによる『界』の構造と学知の関連の分析」東京大学大学院教育学研究科 2023 年度修士論文.
10. **Grootendorst, Maarten**, 2022, "BERTopic: Neural Topic Modeling with a Class-Based TF-IDF Procedure," *arXiv preprint arXiv:2203.05794*.
11. 渡辺健太郎, 2023, 「日本の研究者を対象とした無作為抽出調査は可能か」科学技術社会論学会第 22 回年次研究大会発表資料.
12. **Mochihashi, Daichi**, 2023, "Researcher2Vec: Neural Linear Model of Scholar Recommendation for Funding Agencies," *International Society for Scientometrics and Informatics (ISSI 2023)*, 2: 329-335.
13. **Nagao, Hiroyoshi & Marie Katsurai**, 2024, "Researcher Representations Based on Aggregating Embeddings of Publication Titles: A Case Study in a

Japanese Academic Database," *Proceedings of the Fourth Workshop on Scholarly Document Processing (SDP 2024)*: 277-82.

