

## 科学研究費助成事業 研究成果報告書

平成 29 年 6 月 8 日現在

機関番号：32661

研究種目：研究活動スタート支援

研究期間：2015～2016

課題番号：15H06637

研究課題名（和文）言語統計解析に基づく日本語と中国語の帰納的推論の比較研究

研究課題名（英文）The comparative study of inductive reasoning between Japanese and Chinese based on a statistical analysis of Japanese and Chinese corpora

研究代表者

張 寓杰 (ZHANG, Yujie)

東邦大学・理学部・博士研究員

研究者番号：70759894

交付決定額（研究期間全体）：（直接経費） 2,200,000円

研究成果の概要（和文）：本研究の目的は、言語統計解析に基づく計算モデルを用いて、日本語と中国語の帰納的推論の比較研究を行うことである。まず中国語と日本語における名詞と動詞の関係に名詞と形容詞の関係を加えた大規模言語データの統計解析に基づき、より精度の高い両言語の帰納的推論の計算モデルを構築し、心理学実験によりモデルの妥当性を検証した。さらに同じ入力に対する両言語の計算モデルの出力結果を比較し、両言語の背景にある文化や社会システムの共通性と差異を検討し、文化比較の新しい方法論を提案した。

研究成果の概要（英文）：The purpose of this research is to compare the inductive reasoning between Japanese and Chinese using computational models based on a statistical analysis of Japanese and Chinese corpora. Firstly, based on the statistical analysis of large-scale language data including not only the relationship between nouns and verbs but also the relationship between nouns and adjectives in Chinese and Japanese, we constructed computational models of inductive reasoning of both languages with higher accuracy. Secondly, using psychological experiment, we verified the validity of the models. Furthermore, through the comparison of computational models' simulation for the same input, we examined the commonality and differences of culture and social system in the background of both languages, and proposed a new methodology of cultural comparison.

研究分野：推論、言語処理、認知科学

キーワード：帰納的推論 計算モデル 比較

### 1. 研究開始当初の背景

(1)帰納的推論とは、いくつかの個別知識から、一般法則を導き出す推論を意味する。帰納的推論は単に科学的推論に限らず、広く日常生活でも用いられることが多い、極めて基本的な人間の思考過程の一つである。本研究では心理学や認知科学の分野で、広く一般的に用いられている以下の帰納的推論の形式を取り扱う(例えば Rips, 1975; Osherson, Smith, Wikie, Lopez, and Shafir, 1990; Sakamoto & Nakagawa, 2007, 2008, 2010 など)。

前提:

Aさんはステーキが好きである。(正事例)

Aさんはうどんが好きではない。(負事例)

結論: Aさんはハンバーグが好きである。

この形式では、線分の上部が前提命題で、線分の下部が結論命題である。この形式の帰納的推論に関しては、従来さまざまな理論が提唱されている。たとえば Osherson(1990)は、この種の形式における推論を「カテゴリに基づく帰納的推論(category-based induction)」という仮説に基づき考察し、Sloman(1993)は類似性に基づく理論を提案している。これらの帰納的推論の各仮説では前提事例からカテゴリへの一般化や共有属性からの単語間の類似性の推定といった何らかの内的な心理学的メカニズムを想定している。しかし、そのような内的なメカニズムを直接心理学実験だけで実証することは困難である。帰納的推論の内的メカニズムを説明するために、今までに様々な計算モデルが提案されてきた(Rips, 1975; Osherson, 1990; Sloman, 1993; Sanjana, 2002)。しかし、これらのモデルは、全て共通して、非常に限られた知識領域のみを対象とした帰納的推論以外は検証していないという問題点を含んでいる。

(2)Sakamoto & Nakagawa (2007, 2008, 2010)は以上の既存のモデルの問題点に対して、心理学実験による評価に依存せず、大規模言語データの統計解析を用いて数万語を含む確率的言語知識構造を構築し、より広範な概念についての予測が可能な帰納的推論の計算モデルを構築した。ただし、坂本の研究はすべて日本語に限られており、日本語以外での計算モデルの可能性については、全く考慮されていない。

(3)この問題点を考慮し、張ほか(2013)は中国語の帰納的推論の計算モデルを構成して、このモデルの日本語以外での有用性を明らかにし、日本語と中国語における帰納的推論の計算モデルを比較し、「人間の帰納的推論は必ずしも個々の言語表現に直接依存しておらず、両言語に共通する内的メカニズムに基づいている」という仮説を実証し、両言語の背景にある文化や社会システムの共通性や差異を考察している。

しかし、これら両言語のモデルは共に、名詞と動詞の関係しか用いていないという問題点を含んでいる。

### 2. 研究の目的

(1)先行研究の問題点を解消し、モデルの精度を向上させるため、中国語の言語コーパスを拡張し、名詞と形容詞(中国語では実質的に形容詞の役割を担っている「名詞修飾語」)の関係を加え、その統計解析に基づき、確率的言語知識構造を構築し、帰納的推論の計算モデルを構成する。

(2)日本語についても中国語と同様に、名詞と形容詞の関係を加えた大規模言語データの統計解析に基づき、確率的言語知識構造を構築し、帰納的推論の計算モデルを構成する。

(3)中国語、日本語の各々のモデルのシミュレーション結果と、中国人、および日本人各々の実験参加者による心理学実験の結果を比較し、各々のモデルの妥当性を検証する。

(4)上記の過程で構築され、その心理学的妥当性が検証された日本語と中国語の計算モデルのシミュレーション結果と先行研究における、両言語のシミュレーション結果を各々比較する。その比較結果を踏まえて形容詞と名詞の関係を加えることで両言語のモデルとも精度が向上していることを検証する。

(5)本研究で新しく構築した日本語と中国語における帰納的推論の計算モデルに同じ入力を行いその出力結果を比較する。両言語の背景にある文化や社会システムについて先行研究より、より精細に幅広く比較検討し、日本と中国の社会はどのように変化したのかを踏まえて、その共通性や差異を明らかにする。

### 3. 研究の方法

(1)中国語の言語コーパスを拡張する。現在使われているコーパスは ChineseTreebank4.0 (2010 取得)、人民日報タグ付きコーパス (1998)、新京報電子版(2010 取得)、文学作品の電子テキスト(2010 取得)、合計 651.44MB であるが、近年はインターネットでブログ、SNS、ツイッターなどの利用者が急激に増加し、これらのデータを収集する必要がある。さらに、近年の新聞や文学作品のデータを増やし、より現実的に社会全体を反映するようにコーパスを拡張する。

(2)中国語における以下の手順に従って確率的言語知識構造を構成する。

形態素解析

係り受け解析

単語間共起頻度の抽出

Kameya and Sato(2005)のアルゴリズムに基づくクラスタリング

本研究ではまず上記の と の係り受け解析の結果得られた、「形容詞と名詞」、「名詞(目的語)と動詞」、「名詞(主語)と動詞(述語)」の各対について、 で全言語データ中の共起頻度を計算する。次に各対の共起頻度に基づき、 の方法を用いて各対の共起確率と各条件付き確率、潜在クラスの確率の最尤値を推定する。ここで「形容詞と名詞」、「名詞(目的語)と動詞」、「名詞(主語)と動詞(述語)」の各対について推定された条件付き確率と潜在クラスの確率の総体を確率的言語知識構造と呼ぶ。

(3)中国語の確率的言語知識構造に基づき、中国語の帰納的推論の計算モデルを構成する。計算モデルは以下のようなカーネル関数に基づき構成される。

$$v(N_i^c) = aSIM_+(N_i^c) + bSIM_-(N_i^c) - h \quad (1)$$

$$SIM_+(N_i^c) = \sum_j^{n^+} e^{-\beta d_{ij}^+} \quad (2)$$

$$SIM_-(N_i^c) = \sum_j^{n^-} e^{-\beta d_{ij}^-} \quad (3)$$

$$d_{ij}^+ = \sqrt{\sum_k^m (P(c_k|N_i^c) - P(c_k|N_j^+))^2} \quad (4)$$

$$d_{ij}^- = \sqrt{\sum_k^m (P(c_k|N_i^c) - P(c_k|N_j^-))^2} \quad (5)$$

ここで  $v(N_i^c)$  は結論  $(N_i^c)$  の尤もらしさの値、 $P(c_k|N_i^c)$ 、 $P(c_k|N_j^+)$ 、 $P(c_k|N_j^-)$  は各々結論  $(N_i^c)$ 、正事例  $(N_j^+)$ 、負事例  $(N_j^-)$  が与えられたときの潜在クラス  $(c_k)$  の条件付確率を示している。本研究ではこのモデルの形式を用いて、中国語の帰納的推論の計算モデルを構築し、シミュレーションを行う。

(4)中国人の実験参加者に対して心理学実験を実施し、シミュレーション結果と実験結果を定量的に比較し、計算モデルの妥当性を実証する。

(5)本研究で構築した新しい中国語のモデルの妥当性と先行研究、張ほか(2013)の中国語のモデルの妥当性を比較し、新しいモデルの優越性を検証する。

(6)日本語に対しても中国語と同じように名

詞と形容詞の関係を加え、確率的言語知識構造を構成し、帰納的推論の計算モデルを構築し、シミュレーションを行う。

(7)日本語のモデルについても、シミュレーション結果と実験結果を比較し、モデルの妥当性を実証する。さらに先行研究と比較し新しいモデルの優越性を検証する。

(8)先行研究の課題の事例を拡張し、より広い範囲で様々な分野の課題を選び、日本語と中国語各々の計算モデルに入力してシミュレーションを行い、その都度出力される単語の意味内容を比較することで、単なる心理学実験や調査研究では計り知れない、日本と中国、両国の文化や社会の特徴を、より幅広く比較考察する。すなわちこのようにして構成された、より精度の高い計算モデルとシミュレーションに基づき、全く新しい文化比較の客観的方法を提案する。

#### 4. 研究成果

(1)まず、中国語における構築したモデルと先行研究のモデルのシミュレーション結果の比較として1つの例を挙げる。

表1の正事例と負事例を入力し、出力結果の尤もらしさの値の上位10個の単語を抽出した比較結果である。

本研究の結果は全体的に「芸術」潜在クラスに属する単語が出力されている一方、先行研究の結果では、「防風林」、「感嘆符」、「台湾にある町」など、「芸術」潜在クラスに全く関係ない単語が含まれている。

このように本研究と先行研究のシミュレーション結果を比較し、本研究のモデルは先行研究よりモデルの精度が良くなったと考えられる。

さらに、構築した計算モデルの妥当性を検証するために、帰納的推論の心理学実験を実施した。実験材料は計算モデルのシミュレーションに基づき、張ほか(2013)の研究で使用した8組の課題の正事例と負事例を用い、結論には課題ごとにシミュレーション結果から上位10個、中位10個、下位10個の単語を抽出し、合計30個の単語を用いた。

この8課題の実験材料を使い、中国人の被験者38名に対して、インターネットでのアンケート調査を実施した。調査での質問内容は課題ごとに異なっており、現実の様々な場面に対して、帰納的推論を行うように設定した。たとえば「趣味」課題の質問は、「ある人はバスケットボールとサッカーが好き。社会学と政治学が好きではない。この人が以下のものが好きである可能性に対して、どう思いますか？」である。結論の評定には「かなりあり得る～まったくありえない」の5段階評定を用いた。

心理学実験結果のデータから、「かなりあり得る」を5、「あり得る」を4、「わからない」を3、「ありえない」を2、「全くありえない」

を1として各結論ごとに被験者の平均を算出した。

表1. 中国語における本研究と先行研究のシミュレーション結果の比較

正事例	日本語訳
芭蕾	バレエ
絵画	絵画
負事例	日本語訳
物理学	物理学
科学	科学

上位10個 (本研究)    上位10個 (日本語訳)    上位10個 (先行研究)    上位10個 (日本語訳)

芭蕾	バレエ	絵画	絵画
絵画	絵画	芭蕾	バレエ
诗词	詩	防风林	防風林
交谊舞	社交	感叹号	感嘆符
音乐	音楽	连江	台湾にある町
作曲家	作曲家	本体主义	存在論
舞蹈	ダンス	古体诗	古体詩
油画	油絵	国医	中国医学
歌坛	音楽業界	练习曲	練習曲
乐曲	楽曲	古兰经	クルアーン

表2. 中国語の計算モデルのシミュレーション結果と心理学実験における評定平均値との相関係数 (\*\*:p<.01)

課題	相関係数
課題1(身分)	0.8075**
課題2(衣料)	0.4871**
課題3(交通)	0.7511**
課題4(趣味)	0.6947**
課題5(会議)	0.7728**
課題6(業界)	0.8420**
課題7(商品)	0.8601**
課題8(場所)	0.8095**
平均	0.6661**

さらに各課題ごとに評定の平均値とシミュレーション結果(a=1, b=-1, h=0, =1の

場合)の相関係数を算出した。

表2に示すように、当該モデルのシミュレーション結果と心理学実験における評定平均値との相関係数は実験で用いた8課題ともに高い値を示しており、すべて検定結果も1パーセント水準で有意である。この結果から、中国語の帰納的推論の計算モデルの心理学的妥当性が実証されたと言える。

つまり、中国語における名詞と動詞の関係に名詞と形容詞の関係を加え、計算モデルを拡張し、心理学実験によりモデルの妥当性を検証した。

(2)日本語の名詞と動詞の関係に名詞と形容詞の関係を加え、計算モデルを拡張し、心理学実験によりモデルの妥当性を検証する。

表3. 日本語における趣味課題のシミュレーション結果

	日本語	もってもらしさ
	バスケットボール	0.8972
	サッカー	0.8262
	野球	0.5245
最上位十個	卓球	0.5195
	テニス	0.4969
	ラグビー	0.4852
	ゲーム	0.4844
	読書	0.4671
	スポーツ	0.4508
	ソフトボール	0.4365
	水泳	0.2635
	稽古	0.2576
	洗濯	0.2560
	整理	0.2330
中位十個	俳句	0.2323
	恋愛	0.2304
	雑談	0.2243
	歌舞伎	0.2048
	園芸	0.2038
	相撲	0.1423
	経営学	-0.3840
	経済学	-0.3872
	工学	-0.3876
最下位十個	病理学	-0.3998
	教育学	-0.4050
	地質学	-0.4091
	建築学	-0.4178
	心理学	-0.4299
	社会学	-0.7576
	政治学	-0.8704

表3は日本語の帰納的推論の計算モデルのシミュレーション結果の一例である。「趣味」課題には、正事例の単語「バスケットボール」と「サッカー」を入力し、負事例の単語「社会学」と「政治学」を入力し、出力結果から「最上位十個」、「中位十個」、「最下位十個」合計30個の単語を抽出し、表3に示した結果になる。上位の結果は「スポーツ」カテゴリに属する単語が出力され、下位の結果には「学問」カテゴリに属する単語が出力され、中位の単語の中に、「俳句」、「歌舞伎」、「相撲」のような、日本の独特な活動を表す単語が出力された。一方、「洗濯」、「雑談」、「恋愛」のような日常生活感がある単語も入っている。張ほか(2013)の研究より、本研究のシミュレーションの結果はさらにきめ細かく日本の文化を反映し、名詞と動詞の関係に名詞と形容詞の関係を加えたことにより、日本語の帰納的推論の計算モデルの精度が良くなったと考えられる。

さらに、構築した計算モデルの妥当性を検証するために、帰納的推論の心理学実験を実施した。実験材料は計算モデルのシミュレーションに基づき、張ほか(2013)の研究で使用した8組の課題の正事例と負事例を用い、結論には課題ごとにシミュレーション結果から上位10個、中位10個、下位10個の単語を抽出し、合計30個の単語を用いた。

この8課題の実験材料を使い、日本語を母語とする大学生と大学院生17名に対して、インターネットでのアンケート調査を実施した。

表4. 日本語における評定の平均値とシミュレーション結果の相関係数( $r$ : $p<.01$ )

課題	相関係数
課題1(身分)	0.7356**
課題2(衣料)	0.8462**
課題3(商品)	0.8798**
課題4(交通)	0.7263**
課題5(会合)	0.6952**
課題6(趣味)	0.8926**
課題7(業界)	0.8613**
課題8(場所)	0.9108**

表4に示すように、日本語における構築した計算モデルのシミュレーション結果と心理学実験における評定平均値との相関係数は、実験で用いた8課題ともに高い値を示しており、すべて検定結果も1%水準で有意である。この結果から、本研究のモデルの心理学的妥当性が実証されたと言える。

(3)日本語と中国語の帰納的推論の検索システムを構成するため、ソフトウェアXAMPPを使い、PHPのローカル環境で検索システムを構成した。図1は日本語の検索システムの一例である。

例である。

ローカルサーバで検索システムの画面を開くと、以下の図1に示した画面で正事例、負事例、表示する単語数の部分が空欄で入力待ち状態となる。次に任意の正事例の単語と負事例の単語を入力する。たとえば、本研究で使った「趣味」課題に、課題の内容「ある人はバスケットボールとサッカーが好き。社会学と政治学が好きではない。この人はどんなものが好きか?」という事態を想定し、図1に示したように、正事例の単語「バスケットボール」と「サッカー」を入力し、負事例の単語「社会学」と「政治学」を入力する。さらに出力する単語数、例として30を入力し、「実行」ボタンをクリックすると、出力結果が画面に表示される。

図1. 検索システムの画面例 (単語入力と結果出力)

検索システムに用意された二万以上の単語から任意の単語を選んで、正事例、負事例として入力すると、それに対応して人間の帰納的推論を模擬する結果が出力される。どんな組み合わせで、どんな結果が出力されるのかは簡単に予測できず、応用上、大変興味深いシステムと考えられる。この検索システムはまだ試作段階であり、現在はローカル環境に限定されている。

今後はシミュレーションでの課題の種類や実験参加者数を増やし、システムの妥当性と検索システムや教育システムへの応用を視

野に入れて、システムの有用性も検証していきたい。さらに、検索システムの実用化を目指して、ローカル環境だけではなく、一般のインターネット上に接続できるようにする予定である。

また、日本語、中国語、英語の言語ビッグデータの統計解析に基づき、3カ国語の帰納的推論の計算モデルを構築し、そのコンピュータシミュレーションを通じて、帰納的推論の国際比較を行う予定である。

#### <引用文献>

Kameya, Y., & Sato, T. (2005) "Computation of probabilistic attributes using a statistical analysis of Japanese corpora", Proceedings of Symposium on Large-scale Knowledge Resources, 65-68.

Kayo Sakamoto, Asuka Terai, Masanori Nakagawa, (2007) "Computational models of inductive reasoning using a statistical analysis of a Japanese corpus", Cognitive Systems Research, 8, 282-299.

Kayo Sakamoto, Masanori Nakagawa, (2008) "A Computational Model of Risk-Context-Dependent Inductive Reasoning Based on a Support Vector Machine", T. Tokunaga and A. Ortega (Eds.): LKR2008, LNAI 4938, Springer-Verlag Berlin Heidelberg, pp.295-309.

Kayo Sakamoto, Fang Xie, Masanori Nakagawa, (2010) "Syntactic Dependency Analysis Reveals Semantic Concept Structure Underlying Inductive Reasoning: Towards a Domain-Inclusive Structure that Enables Context-Dependent Knowledge Selection", Cognitive Studies, Vol.17, No.1, 143-168.

Osherson, D. N., Smith, E. E., Wilkie, O., López, A., and Shafir, E., (1990) "Category based induction", Psychological Review, 97, 185-200.

Rips, L. J., (1975) "Inductive judgment about natural categories", Journal of Verbal Learning and Verbal Behavior, 14, 665-681.

Sloman, S., A., (1993) "Feature based Induction", Cognition, 49, 67-96.

張寓杰, 寺井あすか, 董媛, 王月, 中川正宣, (2013) "日本語と中国語における帰納的推論の比較研究 言語統計解析に基づく計算モデルを用いて", 認知科学, No. 20 Vol. 4, 439-469.

#### 5. 主な発表論文等

(研究代表者, 研究分担者及び連携研究者には下線)

[学会発表](計 4 件)

Masanori Nakagawa, Yujie Zhang, Yali

Zhang, Asuka Terai, Wanying Wang, Kenichi Kikuchi The Construction of Computational Model for Inductive Reasoning in Chinese with a focus on predicate verb. The 12<sup>th</sup> East Asian Science Technology and Society (EASTS) Network Conference. 2016年11月19日. 清華大学(中国・北京市)

張寓杰, 張亜麗, 寺井あすか, 王婉瑩, 菊地賢一, 中川正宣. 中国語における述語動詞を中心とした帰納的推論の計算モデルの構築. 日本認知科学会第33回大会, 2016年9月18日. 北海道大学フロンティア応用科学研究棟1階(北海道・札幌市).

Yujie Zhang, Kenichi Kikuchi, Asuka Terai, Luning Ruan, Masanori Nakagawa. A Computational Model of Inductive Reasoning Based on a Statistical Analysis of Japanese Corpora-An Examination of Similarity Functions. The 5<sup>th</sup> Annual International Conference on Cognitive and Behavioral Psychology (CBP2016), 2016年2月22日. Singapore (Singapore).

張寓杰, 孫星越, 菊地賢一, 中川正宣. 中国語における帰納的推論の計算モデルの構成 名詞と形容詞及び名詞と動詞の関係を以て. 日本認知科学会第32回大会, 2015年9月20日. 千葉大学総合校舎D号館(千葉県千葉市).

#### 6. 研究組織

##### (1) 研究代表者

張 寓杰 (ZHANG, Yujie)

東邦大学・理学部・博士研究員

研究者番号: 70759894

##### (4) 研究協力者

上西 秀和 (KAMINISHI, Hidekazu)