

共起語分布の言語間差異を手がかりとする日英対照ブログ分析支援*

中崎 寛之[†] 川場 真理子[†] 山崎 小有里[‡] 宇津呂 武仁[†] 福原 知宏[§]

筑波大学大学院 システム情報工学研究科[†] ,

筑波大学 第三学群工学システム学類[‡] , 東京大学 人工物工学研究センター[§]

1 はじめに

本論文では、ある同一のトピックについてまとめた規模の記述が書かれたブログサイトを、日英各言語について検索し、その記述内容を二言語間で対照分析する方式を提案する(図1)[中崎09]。あるトピックの日英二言語表現を得る際には、Wikipediaの日英二言語エントリを用いる。ブログサイトの検索においては、特定トピックを表すキーワードを用いて商用検索エンジンAPIにより上位のブログサイトを収集し、これを、特定トピックを表すキーワード、および Wikipediaから収集した関連語の出現数順にランキングする方法を用いる[川場08]。この方法により、そのトピックについての記述が多く含まれる有用なブログサイト、および、それらのブログサイト中における有用な記事を上位にランキングすることが可能となる。さらに、これまでに行った評価実験では、それらのブログサイトの内容を日英二言語間で対照分析することにより、ブログ特有の個人レベルの情報や意見における国間差異が多数観測されている。

2 評価用トピック

評価用トピック候補として、Wikipediaにおいて、日英 Wikipediaエントリが存在し、日英ブログ空間におけるエントリ名のヒット数が一定数以上となるものを50個程度選定した。このトピック候補の中から、評価用トピックとして、「捕鯨」、「臓器移植」、「喫煙」、「サブプライムローン」の社会系トピック4種類を選定した。これらの評価用トピックの要約と日英ブログにおける評価用トピックに対する主な意見を表1に示す。

3 二言語対照ブログ分析

3.1 ブログサイト検索

Wikipediaエントリをトピックとするブログサイトの検索においては、日本語ブログの検索には、Yahoo!Japan

* Assisting Comparative Analysis of Japanese/English Blogs based on Cross-Lingual Gaps between Keyword Distributions

[†]Hiroyuki Nakasaki, Mariko Kawaba, Takehito Utsuro, Graduate School of Systems and Information Engineering, University of Tsukuba

[‡]Sayuri Yamazaki, College of Engineering Systems, Third Cluster of Colleges, University of Tsukuba

[§]Tomohiro Fukuhara, Research into Artifacts, Center for Engineering, University of Tokyo

検索APIを、英語ブログの検索には米Yahoo!検索APIを利用し、日本語ブログでは大手11社¹、英語ブログでは大手12社²のブログ会社のドメインに限って検索を行った。検索の際には、Wikipediaエントリのエントリ名を検索クエリとして、複数のブログホストを一度に指定して検索し、1000件の記事を取得する。しかしAPIの検索ではブログ記事単位の検索になるので、同一著者のブログ記事は一つのブログサイトにまとめるという作業を行った。その結果、一トピックあたり約200前後のブログサイトを取得することができた。その後、各ブログサイトにおいて、Wikipediaエントリのエントリ名のヒット数を求め、ヒット数が下限未満(本論文では、10)のブログサイトを削除した。

3.2 ブログ記事検索

次に、検索した日英ブログサイト集合の中から、トピックについて詳しく書かれたブログ記事を検索する。手法としては、トピック名がタイトルである各言語のWikipediaエントリのリダイレクト、さらにWikipediaエントリの本文から太字、他エントリリンクをブログ記事検索のための関連語として抽出する。そして、抽出した関連語のいずれかが出現する各言語のブログ記事をブログサイト集合内からそれぞれ検索する。各トピックのWikipedia関連語数、各トピックで検索したブログサイト数、検索したブログサイト中でWikipedia関連語のいずれかが出現したブログ記事数、検索したブログ記事本文に含まれる総形態素数および総単語数を表2に示す。

3.3 ブログ記事からの共起語抽出

本研究では、対照分析の方法として、各言語のブログに出現する共起語を用いる。まず、検索した日本語ブログ記事からは名詞句を抽出し、検索した英語ブログ記事からは一単語、二単語連語、三単語連語を抽出し、それぞれの頻度統計と出現確率を求める。日本語名詞句 X_J の日本語ブログにおける出現確率 $P_J(X_J)$ と、英語一単語・二単語連語・三単語連語 Y_E の英語ブログにおける出現

¹FC2.com, yahoo.co.jp, rakuten.ne.jp, ameblo.jp, goo.ne.jp, livedoor.jp, Seesaa.net, jugem.jp, yaplog.jp, webry.info.jp, hatena.ne.jp

²blogspot.com, msnblogs.net, spaces.live.com, livejournal.com, vox.com, multiply.com, typepad.com, aol.com, blogsome.com, wordpress.com, blog-king.net, blogster.com

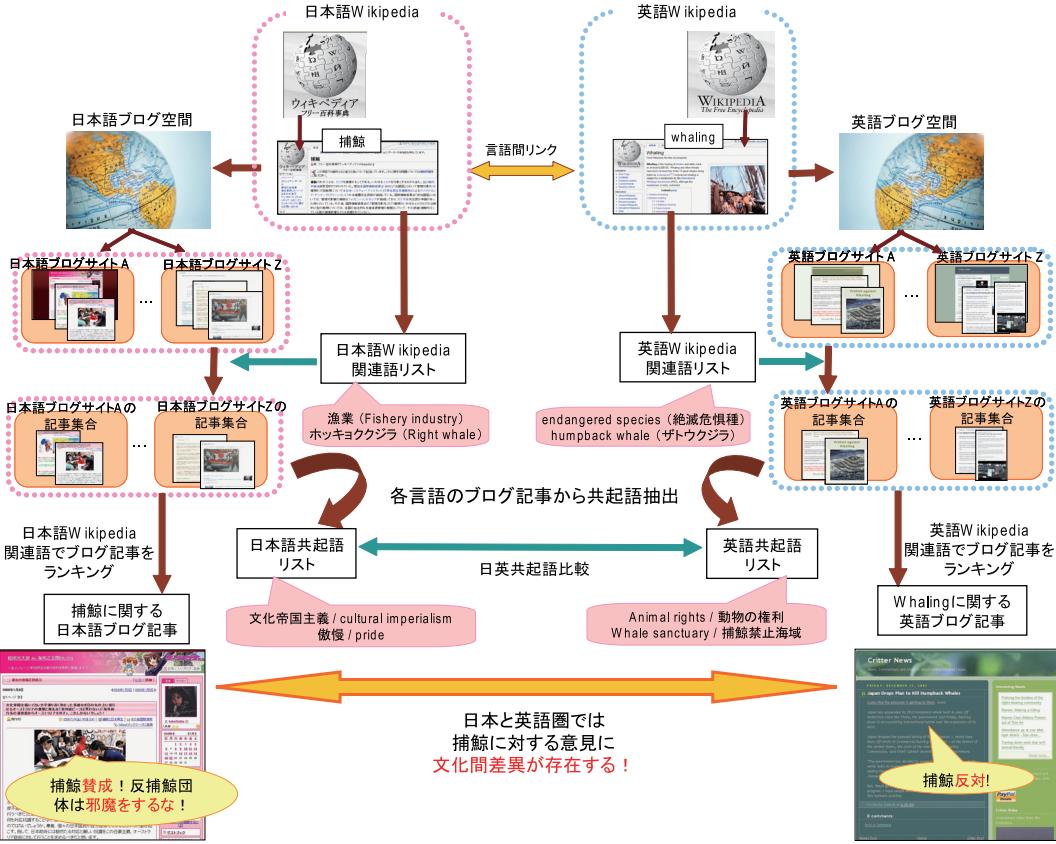


図 1: 二言語対照ブログ分析の全体的枠組み

確率 $P_E(Y_E)$ を以下のようにそれぞれ定義する。

$$P_J(X_J) = \frac{X_J \text{ の出現頻度}}{\text{対象日本語ブログサイト集合内の総形態素数}}$$

$$P_E(Y_E) = \frac{Y_E \text{ の出現頻度}}{\text{対象英語ブログサイト集合内の総単語数}}$$

また、抽出した語句の訳語が相手言語ブログに出現するか調べるために、Wikipedia の言語間リンクを使用して語句の訳語を求める。Wikipedia で語句の対訳を取得できない場合は、英辞郎³で語句の対訳を取得する。さらに、抽出した語句の出現率と対訳語句の出現率から、相手言語ブログと比較した出現確率比を求める。本研究では、抽出した日本語名詞句 X_J と X_J の英訳 X_E の出現確率比 $R_J(X_J, X_E)$ と、英語単語・二単語連語・三単語連語 Y_E と Y_E の和訳 Y_J の出現確率比 $R_E(Y_E, Y_J)$ を以下のように定義した。

$$R_J(X_J, X_E) = \frac{P_J(X_J)}{P_E(X_E)}, R_E(Y_E, Y_J) = \frac{P_E(Y_E)}{P_J(Y_J)}$$

そして、定義した出現確率比で各言語の共起語をランクインし、それぞれの言語で高い出現確率比の共起語

を比較することで、共起語単位でブログ空間におけるトピックの文化間ギャップ発見を支援することができる。

トピック「臓器移植」において、日英ブログ記事から抽出した共起語例を表 3 に示す。表 3 から、片言語で特徴的な共起語もあれば、両言語で多く出現する共起語も存在することがわかる。ここで定義した各共起語の出現確率と出現確率比を 3.5 節の共起語マップの座標として用いる。

3.4 ブログサイト・ブログ記事の順位付け

各言語のブログサイト群およびブログ記事群の順位付けにおいては、3.2 節で抽出した Wikipedia 関連語を用いる。ブログ記事は、以下のスコアの降順に順位付けする。

$$\text{PostScore}(p) = \sum_t (\text{weight}(\text{type}(t)) \times \text{freq}(t))$$

$\text{weight}(\text{type}(t))$ は、Wikipedia 関連語 t の種類 $\text{type}(t)$ に付与する重みで、 $\text{freq}(t)$ は、ブログ記事 p 内における Wikipedia 関連語 t の出現頻度である。また、Wikipedia 関連語 t の種類 $\text{type}(t)$ がリダイレクトの場合は重みを 3、太字の場合は重みを 2、他エントリリンクの場合は重みを 0.5 とする。また、ブログサイトは、各ブログサイトに含まれるブログ記事のスコアの総和の降順に順位

³<http://www.eijiro.jp/>

表 1: 各トピックに関連する日英ブログにおける意見の要約

トピック — 概要		
日英ブログ間の意見の差異		
(日本語ブログ)	(英語ブログ)	
捕鯨 (Whaling) — 捕鯨問題において、捕鯨賛成派と捕鯨反対派が対立している。多くのブログが捕鯨賛成派、反捕鯨団体を激しく非難している。また、捕鯨について書いているプロガーには、右寄りの考えを持つ人が多くみられた。	多くのブログが捕鯨反対派。特に日本の捕鯨を激しく非難している。また、いくつかのプロガーはホエールウォッッチングについて書いている。	
臓器移植 (Organ transplant) — 治療のために、提供されたドナーの臓器を患者に移植する医療法		
多くのブログは日本の臓器移植法改正の必要性を訴えている。また、いくつかのブログでは、日本の医者によって行われた病気腎移植問題のことに注目している。	多くのブログで、臓器不足という現状から、臓器移植のドナー登録を強く推奨している。また、いくつかのブログでは中国の違法臓器摘出を非難している。	
喫煙 (Tobacco smoking) — 喫煙することで、人の健康を損なうということで知られている。		
多くのブログで、健康や喫煙マナーの悪さを理由に喫煙に反対しているが、一部のブログは喫煙賛成派である喫煙者のプロガーであった。	多くのブログで、肺がんの原因である喫煙に反対している。	
サブプライムローン (Subprime lending) — 近年発生した世界金融危機の大きな原因の一つ		
多くのブログで、米国のサブプライム問題による影響で日本経済が悪化したと指摘。	多くのプロガーが経済学者で、サブプライムローンによって発生した住宅バブルや、現在の金融危機や経営危機の発生原因など考察している。	

表 2: Wikipedia から抽出した関連語数、ブログサイト・記事数、ブログ記事中の形態素・単語数

トピック	Wikipedia 関連語数	ブログサイト数	ブログ記事数	総形態素数/ 総単語数
捕鯨	162 / 174	121 / 239	2232 / 6532	5024966 / 2611942
臓器移植	100 / 231	89 / 206	696 / 1301	995927 / 781476
喫煙	399 / 276	86 / 252	1481 / 400	1323767 / 492727
サブプライムローン	39 / 68	134 / 205	1088 / 1216	980552 / 883450

表 4: 日英ブログ記事上位 10 件の概要（「臓器移植」、記事順位 / サイト順位 / 概要）

日本語	英語
(記事) 2 位, 3 位, 4 位, 6 位, 9 位 / (サイト) 1 位 / 病気腎移植のニュースを取り上げている。病気腎移植に反対する日本移植学会を批判。	(記事) 2 位, 6 位 / (サイト) 8 位 / 中国の違法臓器摘出を批判しているニュースを紹介
(記事) 7 位 / (サイト) 14 位 / 脳死移植に反対。臓器移植法の改正は慎重に行うべきと主張。	(記事) 3 位 / (サイト) 2 位 / 違法臓器摘出を批判しているニュースを紹介
(記事) 8 位, 10 位 / (サイト) 7 位 / 病気腎移植に反対。患者が完治するとは思えないと主張。	(記事) 8 位 / (サイト) 7 位 / 臓器提供に関するニュース記事を紹介。ドナー登録することを強く推奨。

付けする。

トピック「臓器移植」の場合のブログ記事上位 10 件の概要を表 4 に示す。多くのブログ記事はトピックについて詳しく書いてあるブログ記事であり、いくつかのブログ記事ではそのトピックに対してブログ著者の賛成意見や反対意見が述べられていることがわかる。

3.5 文化間差異の発見支援のための共起語マップ

本論文では、同一トピックの日英ブログにおける文化間差異の発見支援ツールとして、日英ブログの共起語マップを利用する。トピック「臓器移植」について、日英ブログから抽出した共起語例を共起語マップに配置した結果

を図 2 に示す。共起語マップの横軸は、各共起語の出現確率比を表し、縦軸は各共起語の単言語における出現確率を表す。日本語ブログから抽出した日本語共起語 X_J は、座標 $(-R_J(X_J, X_E), P_J(X_J))$ に表示される。このとき、日本語共起語 X_J の英訳 X_E が英語ブログに出現しない場合は、日本語ブログで特徴的な共起語として左端軸上に表示される。一方、英語ブログから抽出した英語共起語 X_E は、座標 $(R_E(X_E, X_J), P_E(X_E))$ に表示される。そして、英語共起語 X_E の和訳 X_J が日本語ブログに出現しない場合は、英語ブログで特徴的な共起語として右端軸上に表示される。このマップ上では、片言語のみで特徴的である話題から抽出された共起語群は、縦軸から大きく離れて表示される傾向にある。逆に、両言語で共通している話題から抽出された共起語群は、日英対訳となる共起語の組が縦軸から近い位置に表示される。

トピック「臓器移植」においては、英語ブログ特有で出現した共起語には、中国の違法臓器移植を批判している意見をあらわすものが多かった。反対に、日本語ブログで特徴的な共起語は、病気腎移植問題と関わりの強いものが多くあらわれた。このことから、日英ブログから抽出した共起語が日英ブログの文化間差異の発見支援となることがわかった。

表 3: 日英ブログから抽出した共起語例（「臓器移植」）

日本語 共起語	出現頻度 (日本語)	英語 共起語	出現頻度 (英語)	$R_J(X_J, X_E)$	日本語 共起語	出現頻度 (日本語)	英語 共起語	出現頻度 (英語)	$R_E(X_E, X_J)$
日本移植 学会	75	対訳なし (The Japan Society for Transplantation)	0	∞	対訳なし (臓器摘出)	0	organ harvesting	270	∞
病気 腎移植	442	対訳なし (transplant using diseased kidney)	0	∞	人権	9	human rights	508	71.93
脳死移植	366	brain-dead transplant	0	∞	ドナーカード	27	donor card	124	5.85
臓器 移植法	123	organ transplant law	3	32.17	臓器提供	200	organ donation	673	4.29

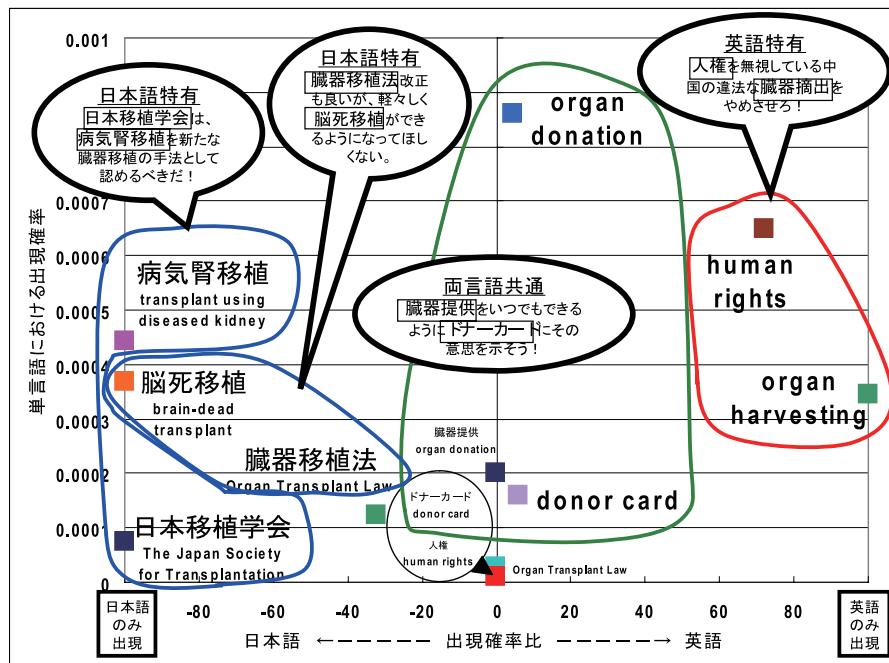


図 2: 日英ブログから抽出した共起語例を用いた共起語マップ（「臓器移植 / Organ transplant」）

4 おわりに

本稿では、Wikipedia エントリを用いてトピックに関する日英ブログサイトを検索し、その記述内容を二言語間で対照分析する方式を提案した。今後は、日英ブログから主観情報・経験情報を多く含む箇所を抽出 [Wiebe05, 乾08] することにより、文化間差異測定尺度の高度化に取り組む。また、多言語ブログバースト分析 [福原07]、複数情報源からのニュースの差異分析 [吉岡07]、Wikipedia百科事典、ニュース、ブログといった異種情報源の相補的利用 [佐藤09] との連携を行う。

参考文献

[福原07] 福原知宏, 宇津呂武仁, 中川裕志: 複数言語間の語彙出現傾向比較による言語横断型ウェブログ関心解析システムの開発、言語処理学会第 13 回年次大会「大規模 Web 研究基盤上での自然言語処理・情報検索研究」ワークショップ論文集, pp. 40–43 (2007).

[乾08] 乾健太郎, 原一夫: 経験マイニング: Web テキストからの個人の経験の抽出と分類、言語処理学会第 14 回年次大会論文集, pp. 1077–1080, 言語処理学会 (2008).

[川場08] 川場真理子, 中崎寛之, 宇津呂武仁, 福原知宏: Wikipedia エントリとブログサイトの対応付けによる日本語ブログ空間のトピック分布推定、情報処理学会研究報告, Vol. 2008, No. (2008-NL-187), pp. 83–90 (2008).

[中崎09] 中崎寛之, 川場真理子, 山崎小有里, 宇津呂武仁, 福原知宏: 同一トピックの日英ブログにおける文化間差異の発見支援、DEIM フォーラム論文集 (2009).

[佐藤09] 佐藤由紀, 中崎寛之, 川場真理子, 宇津呂武仁, 福原知宏: Wikipedia を知識源とするニュース・ブログ間の相補的ナビゲーション、DEIM フォーラム論文集 (2009).

[Wiebe05] Wiebe, J., Wilson, T. and Cardie, C.: Annotating Expressions of Opinions and Emotions in Language, *Language Resources and Evaluation*, Vol. 39, No. 2-3, pp. 165–210 (2005).

[吉岡07] 吉岡真治: 複数のニュース源の差異を考慮したニュース分析の研究、言語処理学会第 13 回年次大会「大規模 Web 研究基盤上での自然言語処理・情報検索研究」ワークショップ論文集, pp. 27–20 (2007).