

# 大規模ウェブテキストを用いた評価表現辞書の自動構築

鍛冶伸裕<sup>†</sup> 喜連川優<sup>†</sup>

近年の情報爆発時代において、Sentiment Analysis や評価情報分析と呼ばれる処理が注目を集めている。こうした処理を行うためには、評価表現とその極性（好評/不評）の組を登録した辞書（評価表現辞書）が必要不可欠である。そのため、大規模な評価表現辞書の構築が重要な研究課題となっている。これまでに、シソーラスや国語辞典などの語彙資源を利用して評価表現辞書を自動構築する手法が提案されている。しかし、そうした語彙資源を利用した手法には、網羅性に欠けるという問題や、句を扱えないという問題がある。一方、検索エンジンを使って種単語との共起頻度を求めるという方法も提案されているが、こちらは計算コストが大きく、大規模な評価表現辞書を構築するには不向きである。また、種単語との共起頻度という考え方はシンプルで分かりやすいが、その精度には疑問が残る。そこで、本論文では、HTML 文書から自動構築した評価文コーパスを用いて評価表現辞書を自動構築する方法を提案する。

## Building Lexicon for Sentiment Analysis from Massive Collection of HTML documents

NOBUHIRO KAJI<sup>†</sup> and MASARU KITSUREGAWA<sup>†</sup>

Recognizing polarity requires a list of polar words and phrases. For the purpose of building such lexicon automatically, a lot of studies have investigated (semi-) unsupervised method of learning polarity of words and phrases. In this paper, we explore to use structural clues that can extract polar sentences from Japanese HTML documents, and build lexicon from the extracted polar sentences. The key idea is to develop the structural clues so that it achieves extremely high precision at the cost of recall. In order to compensate for the low recall, we used massive collection of HTML documents. Our experiment demonstrated the effectiveness of the proposed method.

### 1. はじめに

近年、評価や感情が記述されたテキストを解析する技術が注目を集めている。そのようなテキストを解析するためには、評価表現とその極性（好評/不評）の組を登録した辞書（評価表現辞書）が必要不可欠となる。そのため、大規模な評価表現辞書の構築が重要な研究課題となっている。

これまでに、シソーラスや国語辞典などの語彙資源を利用して評価表現辞書を自動構築する手法が提案されている<sup>13)19)4)5)6)</sup>。しかし、そうした語彙資源を利用した場合、そのエントリに登録されている単語しか扱うことができない。そのため、既存の語彙資源に登録されていない新語や口語（「しょばい」など）に対応できないなど、網羅性に欠けるという問題がある。さらに、句を扱えないということも問題である。例えば「質が高い」は好評極性を持つが「コストが高い」

は不評極性を持つというような事例を正しく把握するためには、句とその評価極性が登録された評価表現辞書が必要となる。

この問題に対する解決策の1つとして、Turney の提案する方法を挙げるができる<sup>21)22)</sup>。Turney は検索エンジンを用いて種単語（「excellent」「poor」など）との共起頻度を求めることで、語句の評価極性の強さを求める手法を提案している。しかし、共起頻度を取得するためには検索エンジンを使わなくてはならないため計算コストが大きく、大規模な評価表現辞書を構築するには不向きである。また、種単語との共起頻度という考え方はシンプルで分かりやすいが、その精度には疑問が残る。

そこで、以下のような辞書構築手法を考案した（図1）。まず、Kajiら<sup>11)</sup>の考案した手法を用いて、評価極性を持つ文を大規模なHTML文書集合から自動抽出する（step 1）。以下ではこのような文を評価文と呼び、抽出された評価文は、好評極性を持つ文（好評文）と不評極性を持つ文（不評文）に分けて保持される。こうして作られたデータセットを評価文コーパスと呼ぶ。

<sup>†</sup> 東京大学 生産技術研究所

Institute of Industrial Science, University of Tokyo

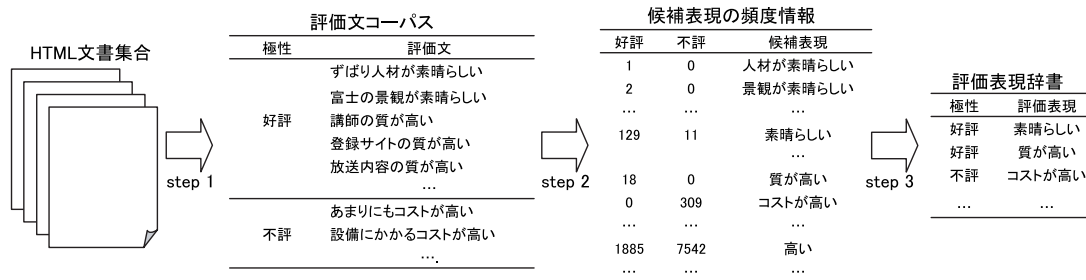


図 1 評価表現辞書構築の流れ

次に、評価文コーパスから評価表現の候補を抽出し、その頻度情報を集計する (step 2)。最後に、得られた頻度情報を利用して、候補表現の中から評価表現だけを選別して辞書に登録する (step 3)。

## 2. 評価文コーパスの自動構築

はじめに、HTML 文書集合から評価文を自動抽出する手法について説明する (step 1)。自動抽出には、HTML 文書中のレイアウト構造やテキスト構造にもとづく手がかりを利用する<sup>11)</sup>。

### 2.1 レイアウト構造の利用

レイアウト構造は箇条書き形式と表形式の 2 種類を利用した。例えば、図 2 のような箇条書きは「良い点」「悪い点」という見出しを持っているため、箇条書きに評価文が記述されていることを判定できる。本論文では「良い点」「悪い点」のような、評価文の存在を示唆する表現を手がかり表現と呼ぶ。手がかり表現リストを手で作成して、それと HTML タグを利用して評価文を自動抽出した。表 1 に手がかり表現リストを示す。

これらは予備実験を通して人手で選定した。表には「良い」のような用言も含まれているが、これらは「所」「点」「面」という 3 つの名詞と組み合わせる。すなわち「良い所」「良い点」「良い面」の 3 つを手がかり句として使うことを意味する。「長所」や「メリット」のような名詞は、単語そのものを手がかり句として使う。なお、詳細は省略しているが「駄目な所」と「ダメな所」または「良い所」と「良いところ」のような表記揺れも網羅的に人手で記述している。

表形式も箇条書き形式の場合とほぼ同様である (図 3)。表の 1 列目に手がかり表現 (気に入った点、イヤな点) が存在していて、これが見出しの働きをしている。そして 2 列目には評価文が記述されている。

### 2.2 テキスト構造の利用

次に、定型的なテキスト構造に着目した。

良い点
<ul style="list-style-type: none"> <li>● 変に加工しない素直な音を出す。</li> <li>● 曲の検索が簡単にできる。</li> <li>● お気に入りのプレイリストを作って楽しめる。</li> </ul>
悪い点
<ul style="list-style-type: none"> <li>● リモコンに液晶表示がない。</li> <li>● ボディに傷や指紋が付きやすい。</li> <li>● ライトを点灯し続けると直ぐに電池がなくなる。</li> </ul>

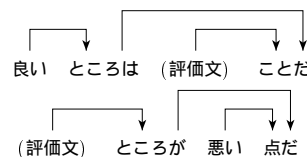
図 2 箇条書き形式で記述された評価文

燃費 (市街地)	7.0km/litter
燃費 (高速)	9.0km/litter
満足度	95%
気に入った点	4ドアなのにカッコよすぎる。
イヤな点	シートがほろくライトが暗い、色がはげてきてる。

図 3 表形式で記述された評価文

- (1) この 良いところは 計算が速いことだ。
- (2) 慣れるまで時間がかかる ところが、悪い点だ。例文 (1) には「計算が速い」、例文 (2) には「慣れるまで時間がかかる」という評価文が含まれている。いずれも「良いところは～こと」「～ところが悪い点」といった定型的なテキスト構造を使って記述されている。

このような評価文は、以下のような語彙統語パターンを用いて自動抽出する。



矢印は文節間の依存関係を表していて、(評価文) はマッチした部分木が評価文として抽出されることを表す。実際のコーパス構築では上記の語彙統語パターンをそのまま用いるのではなく、手がかり表現の部分 (良いところ、悪い点) を、前述の手がかり表現リストを用いて汎化したものを使った。

厳密には文ではなく節と呼ぶべきだが、レイアウト構造を用いて抽出される評価文との整合性を考えて文と呼ぶ。

表 1 実験で使用した手がかり句

好評手がかり句	不評手がかり句
良い, いい, 素晴らしい, 楽しい, 面白い, 最高だ, 満足だ, 便利だ, グッドだ, プラスだ, 良好だ, 大好きだ, 好きだ, 優れる, 長所, 利点, メリット	悪い, 最悪だ, 残念だ, 不満足だ, 不満だ, 嫌だ, 大嫌いだ, 嫌いだ, 駄目だ, いまいちだ, いまひとつだ, バッドだ, マイナスだ, 不便だ, 困る, 短所, デメリット, 問題点, 不満点, 欠点, 難点, 弱点

### 2.3 評価文コーパス

約 10 億件の HTML 文書を用いて評価文コーパスの構築を行った。その結果、約 50 万文からなる評価文コーパスを構築することができた。その内訳は好評文が 220,716 文、不評文が 288,755 文である。表 2 に実際に抽出された評価文の例を示す。構文解析には KNP を用いた。以下の実験でも同様である。

極性	評価文の例
好評	順応性が素晴らしくある。 使い方がわかりやすい。 何と言っても、料金が良心的だ。 費用が高い。
不評	いい加減な意見、ふざけた意見などが出てくる。 エンジンが非力で少々うるさい。

自動構築されたコーパスの質を確認するため、コーパス中の 500 文を 2 人の被験者 (被験者 A, B と呼ぶ) が個別に調べた<sup>11)</sup>。その結果、被験者 A は 91.8%(459/500) の文を適切である判断した。同様に被験者 B は 92%(460/500) の文を適切であると判断した。被験者間での判断の一致率は 93.4%(467/500) であり、このことから、高い精度で評価文が獲得できたと結論づけることができる。

不適切であると判断された評価文を観察した結果、そのほとんどは、評価極性が文脈に依存する文であった。例えば、コーパスには「何しろ情報量が多い」が好評文として登録されていたが、被験者は 2 人ともこれを不適切と判断していた。

## 3. 評価表現の獲得

### 3.1 候補表現の頻度集計

自動構築した評価文コーパスから評価表現の候補 (候補表現と呼ぶ) を抽出する。そして、各候補表現の頻度情報を集計する (step 2)。

しばしば指摘されるように、形容詞は評価極性を持ちやすい。また、1 単語で評価極性が決まる形容詞もあれば、そうでない語 (「高い」など) も存在する。こ

れらを踏まえて、全ての形容詞と形容詞句 (名詞 + 格助詞 + 形容詞) を候補表現とした。ただし、一部の機能語表現を特別処理をする。まず、形容詞に否定を表す機能語 (「ない」と「ぬ」) が付属している場合は、否定をあらわすタグを形容詞に付与する。また、動詞に接尾辞の「やすい」「にくい」が付属している場合、その「動詞 + 接尾辞」を 1 つの形容詞として扱った。例えば「使いにくい」は 1 つの形容詞と考える。

各候補表現について、好評文と不評文における出現頻度を集計した。単純には、好評表現 (不評表現) は好評文 (不評文) に多く出現すると考えられるが、以下のような例外的な場合が存在する。

- (1) 面倒な 準備やテクニックは不要で、非常に簡単です。

この例文は文全体としては肯定的な評価を示しているため、好評文である。しかし、その中に「面倒だ」といった不評表現が出現している。そこで「好評文 (不評文) の主節には、好評表現 (不評表現) が出現しやすい」と仮定して、主節における頻度のみを集計した。

### 3.2 評価表現の選別

候補表現の評価極性の強さを数値化して (この数値を評価極性値と呼ぶ)、それにもとづき候補表現の中から評価表現だけを選別する (step 3)。

各候補表現  $c$  に対して、次のような分割表を作成することができる。

	表 3 分割表	
	<i>pos</i>	<i>neg</i>
$c$	$f(c, pos)$	$f(c, neg)$
$\neg c$	$f(\neg c, pos)$	$f(\neg c, neg)$

$f(c, pos)$  は候補表現  $c$  の好評文における頻度、 $f(\neg c, pos)$  は  $c$  以外の候補表現の頻度の和である。 $f(c, neg)$  と  $f(\neg c, neg)$  も同様である。この分割表から  $c$  の評価極性値を規定する。実験では比較のため、次の 2 種類の手法を試した。

#### $\chi^2$ 値にもとづく評価極性値

候補表現  $c$  の出現の偏りを見積るために  $\chi^2$  値を利用した。表 3 から求めた  $\chi^2$  値は次のようになる。

<http://www.tkl.iis.u-tokyo.ac.jp/~kaji/acp/>  
<http://nlp.kuee.kyoto-u.ac.jp/nl-resource/knp.html>

$$\chi^2(c) = \sum_{x \in (c, -c)} \sum_{y \in (pos, neg)} \frac{\{f(x, y) - \hat{f}(x, y)\}^2}{\hat{f}(x, y)}$$

$\hat{f}(x, y)$  は、候補表現  $c$  の出現確率が好評文と不評文で独立であると仮定したときの  $f(x, y)$  の期待値である。

$\chi^2(c)$  には、 $c$  が好評文と不評文のどちらに多く出現しているのかという情報は反映されていない。そこで  $\chi^2(c)$  を用いて、以下のように評価極性値を設定した。

$$PV_{\chi^2}(c) = \begin{cases} \chi^2(c) & \text{if } P(c|neg) < P(c|pos) \\ -\chi^2(c) & \text{otherwise} \end{cases}$$

$P(c|pos)$  は  $c$  の好評文における出現確率であり、 $P(c|neg)$  は不評文における出現確率である。

$$P(c|pos) = \frac{f(c, pos)}{f(c, pos) + f(-c, pos)}$$

$$P(c|neg) = \frac{f(c, neg)}{f(c, neg) + f(-c, neg)}$$

#### PMI にもとづく評価極性値

PMI (Pointwise Mutual Information) を用いると、候補表現  $c$  と好評文  $pos$  (不評文  $neg$ ) の関連の強さは次のように定義できる。

$$PMI(c, pos) = \log_2 \frac{P(c, pos)}{P(c)P(pos)}$$

$$PMI(c, neg) = \log_2 \frac{P(c, neg)}{P(c)P(neg)}$$

この2つの数値の差を評価極性値とした。これは Turney と同様の考え方である<sup>21)</sup>。

$$\begin{aligned} PV_{PMI}(c) &= PMI(c, pos) - PMI(c, neg) \\ &= \log_2 \frac{P(c, pos)/P(pos)}{P(c, neg)/P(neg)} \\ &= \log_2 \frac{P(c|pos)}{P(c|neg)} \end{aligned}$$

#### 評価表現の選別

上記のように定義した評価極性値と閾値  $\theta (> 0)$  を用いて、ある候補表現  $c$  が評価表現であるかどうかを判定する。まず  $\theta < PV(c)$  であれば、その候補表現は好評表現と考える。同様に、 $PV(c) < -\theta$  であれば不評表現とする。それ以外は評価表現ではないと考える。

#### 4. 実験結果

自動構築した評価表現辞書を用いて、テストデータから評価表現を抽出する実験を行った。

テストデータは、ウェブテキストから無作為に抽出した405の形容詞句に評価極性(好評/不評/中立)をタグ付けして作成した。タグ付けの結果、好評/不評/

中立の数はそれぞれ158/150/97であった。同一データを二人の被験者がタグ付けしたところ Kappa 値は0.73であった。

このテストデータから評価表現を抽出する。基本的には、テストデータに含まれる形容詞句を辞書引きしていくことになる。ただし、辞書に登録されている形容詞(「素晴らしい」など)は、その形容詞を含む全ての形容詞句(「景色が素晴らしい」など)とマッチさせる。

比較のために、Turney<sup>21)</sup>の提案する評価極性値を用いて評価表現辞書を構築し、同様の実験を行った。Turneyの手法は「excellent」「poor」のような種単語が必要となるがここでは「最高」「最低」を用いた。検索エンジンには我々の研究室で開発したローカル検索エンジンとGoogleの2つを試した。前者は約1億5,000万件のHTML文書をインデックスしている。

評価表現抽出の結果を適合率と再現率で評価した(図4)。上のグラフは好評表現抽出の適合率と再現率を、閾値  $\theta$  を変化させながら観察したものである。下のグラフは同様のことを不評表現に対して行った結果である。実験の結果、提案手法はTurneyの手法よりもうまく働くことが確認できた。また、PMIにもとづく評価極性値は、 $\chi^2$  値にもとづくものよりも優れていることが分かった。表6に評価表現の具体例を示す。また、獲得された評価表現数を表4と表5に示す。1行目が閾値  $\theta$ 、2行目が獲得された評価表現数である。

コーパスの大きさによる精度の比較など、さらに詳細な実験結果に興味のある読者は文献<sup>12)</sup>を参考にされたい。

表6 評価表現の具体例

評価表現	$PV_{\chi^2}(c)$	$PV_{PMI}(c)$
謙虚だ	38.3	11.9
支障が無い	34.5	11.8
エキサイティングだ	13.5	10.4
漏れが少ない	9.2	9.8
能力が高い	113.0	6.9
ダサイ	-2.9	-3.1
厄介だ	-11.9	-3.8
消耗が早い	-17.7	-4.3
魅力が無い	-19.4	-4.4
しょぼい	-55.3	-9.1

#### 5. おわりに

本論文ではHTML文書集合から評価表現辞書を自動構築する手法を提案した。そして実験を行い手法の有効性を検証した。今後は、この辞書を用いて、実際のテキストから評価情報抽出を行う予定である。

$\theta$	0	10	20	30	40	50	60
評価表現数	9,670	2,056	1,047	698	533	423	335

$\theta$	0	0.5	1.0	1.5	2.0	2.5	3.0
評価表現数	9,670	9,320	9,039	8,804	8,570	8,398	8,166

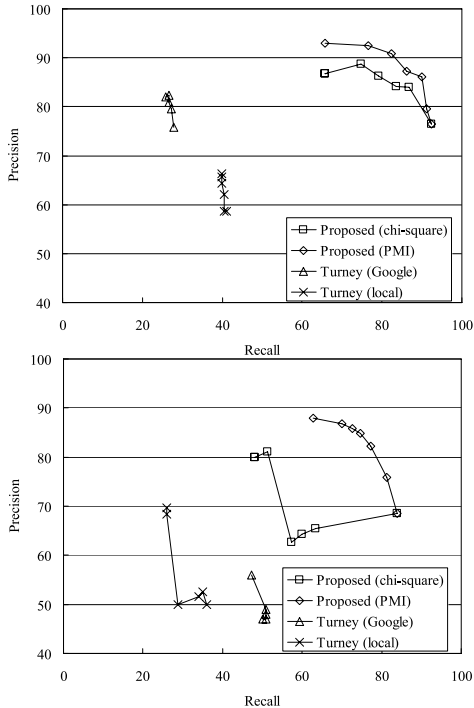


図 4 再現率-適合率曲線 (上: 好評表現, 下: 不評表現)

## 参考文献

- 1) Kenneth Ward Church and Patric Hanks, “Word Association Norms, Mutual Information, and Lexicography”, In Proceedings of ACL, pp. 76-83, 1989.
- 2) Yejin Choi and Claire Cardie and Ellen Riloff and Siddharth Patwardhan, “Identifying Sources of Opinions with Conditional Random Fields and Extraction Patterns”, In Proceedings of HLT/EMNLP, 2005.
- 3) Yejin Choi and Eric Breck and Claire Cardie, “Joint Extraction of Entities and Relations for Opinion Recognition”, In Proceedings of EMNLP, pp. 431-439, 2006.
- 4) Andrea Esuli and Fabrizio Sebastiani, “Determining the Semantic Orientation of Terms through Gloss Classification”, In Proceedings of CIKM, 2005.
- 5) Andrea Esuli and Fabrizio Sebastiani, “Determining Term Subjectivity and Term Orientation for Opinion Mining”, In Proceedings of EACL, pp.193-200, 2006.
- 6) Andrea Esuli and Fabrizio Sebastiani, “SENTIWORDNET: A Publicly Available Lexical Resource for Opinion Mining”, In Proceedings of LREC, 2006.
- 7) Christiane Fellbaum, “WordNet: An Electronic Lexical database”, MIT Press, Cambridge, 1998.
- 8) Vasileios Hatzivassiloglou and Kathleen R. McKeown, “Predicting the Semantic Orientation of Adjectives”, In Proceedings of ACL, pp.174-181, 1997.
- 9) Minqing Hu and Bing Liu, “Mining and Summarizing Customer Reviews”, In Proceedings of KDD, pp.168-177, 2004.
- 10) 乾孝司, 奥村学, “テキストを対象とした評価情報の分析に関する研究動向”, 自然言語処理 13(3), pp.201-242, 2006.
- 11) Nobuhiro Kaji and Masaru Kitsuregawa, “Automatic Construction of Polarity-tagged Corpus from HTML Documents”, In Proceedings of COLING/ACL (Poster Sessions), pp.452-459, 2006.
- 12) Nobuhiro Kaji and Masaru Kitsuregawa, “Building Lexicon for Sentiment Analysis from Massive Collection of HTML Documents”, In Proceedings of EMNLP-CoNLL, pp.1075-1083, 2007.
- 13) Jaap Kamps and Maarten Marx and Robert J. Mokken and Maarten de Rijke, “Using WordNet to Measure Semantic Orientations of Adjectives”, In Proceedings of LREC, 2004.
- 14) Hiroshi Kanayama and Tetsuya Nasukawa, “Fully Automatic Lexicon Expansion for Domain-oriented Sentiment Analysis”, In Proceedings of EMNLP, pp.355-363, 2006.
- 15) Soo-Min Kim and Eduard Hovy, “Identifying and Analyzing Judgement Opinions”, In Proceedings of NAACL-HLT, 2006.
- 16) Soo-Min Kim and Eduard Hovy, “Extracting Opinions, Opinion Holders, and Topics Expressed in Online News Media Text”, In Proceedings of Workshop on Sentiment and Subjectivity in Text, pp.1-8, 2006.
- 17) Moshe Koppel and Jonathan Schler, “The Importance of Neutral Examples for Learning Sentiment”, In Proceedings of Workshop on the Analysis of Informal and Formal Informa-

- tion Exchange during Negotiations, 2005.
- 18) 小林のぞみ, 乾健太郎, 松本裕治, 立石健二, 福島俊一, “意見抽出のための評価表現の収集”, 自然言語処理 12(3), pp.203-222, 2005.
  - 19) Hiroya Takamura and Takashi Inui and Manabu Okumura, “Extracting Semantic Orientation of Words using Spin Model”, In Proceedings of ACL, pp.133-140, 2005.
  - 20) Hiroya Takamura and Takashi Inui and Manabu Okumura, “Latent Variable Models for Semantic Orientation of Phrases”, In Proceedings of EACL, pp.201-108, 2006.
  - 21) Peter D. Turney, “Thumbs up or Thumbs Down ? Semantic Orientation Applied to Unsupervised Classification of Reviews”, In Proceedings of ACL, pp.417-424, 2002.
  - 22) Peter D. Turney and Michael L. Littman, “Unsupervised Learning of Semantic Orientation from a Hundred-Billion-Word Corpus”, National Research Council, Institute for Information Technology, Technical Report ERB-1094(NRC #44929), 2002.
  - 23) Theresa Wilson and Janyce Wiebe and Paul Hoffmann, “Recognizing Contextual Polarity in Phrase-level Sentiment Analysis”, In Proceedings of HLT/EMNLP, 2005.
  - 24) Hong Yu and Yasileios Hatzivassiloglou, “Towards Answering Opinion Questions: Separating Facts from Opinions and Identifying the Polarity of Opinion Sentences”, In Proceedings of the EMNLP, 2003.

謝辞 本研究は文部科学省リーディングプロジェクト e-society:先進的なウェブ解析技術によって支援されている。本研究にあたり, 生産技術研究所協力研究員の田村孝之氏に大変お世話になりました。感謝致します。