

# もう一つの意味的極性「活性／不活性」と知識獲得への応用

橋本 力 鳥澤 健太郎 De Saeger, Stijn 呉 鍾勲 風間 淳一

(独) 情報通信研究機構 ユニバーサルコミュニケーション研究所 情報分析研究室

## 1 はじめに

我々は「活性／不活性」という「助詞＋用言」に関する意味的極性とその獲得法を提案する。また、情報分析、質問応答等にとって重要な、矛盾句対（「ガンを破壊する ⇔ ガンを進行させる」等）と因果関係（「需要が減る ⇒ 失業が増える」等）、コーパスに明記されていない因果関係仮説（「需要が拡大する ⇒ 失業を減少させる」等）の獲得への活性／不活性知識の応用法を提案する。矛盾は 100 万対を適合率 70% 以上で獲得できた。因果関係は 50 万対を適合率約 70%、因果関係仮説は 100 万対を適合率約 57% で獲得した。

## 2 「活性／不活性」とは？

活性／不活性は、「助詞＋用言」（以下、template と呼ぶ）を「活性」、「不活性」、「中立」に分類する。

**活性** 項の指す対象の主たる機能、効果、目的、役割、影響が準備あるいは活性化されることを含意する。（例：「を引き起こす」、「を使う」、「を買う」、「を進行させる」、「を輸入する」、「が増える」）

**不活性** 項の指す対象の主たる機能、効果、目的、役割、影響が抑制あるいは不活性化されることを含意する。（例：「を防ぐ」、「を捨てる」、「を治療する」、「が減る」、「を破壊する」、「が不可能になる」）

**中立** 活性でも不活性でもないもの。（例：「を考える」、「を探す」、「に比例する」）

例えば、「地震を引き起こす」は「地震」の影響が活性化されることを、「津波を防ぐ」は「津波」の影響が不活性化されることを含意する。

template はその活性／不活性の度合いを示す活性／不活性値を持つ。活性 template は +1 までの正の値を、つまり正の極性を持つ。不活性 template は -1 以上の負の値を、つまり負の極性を持つ。活性／不活性値の絶対値の小さいものは中立 template と見なす。なお本研究では、述語が否定されるとその template の極性は反転すると仮定する。例えば「を使う」の活性／不活性値が 0.9 なら、「を使わない」の値は -0.9 となる。

活性／不活性はいわゆる評価極性 (good/bad) [2] とは独立である。例えば「が上達する」も「を発症する」も活性だが後者のみが bad で、「を治療する」も「が頓挫する」も不活性だが後者のみが bad である。

## 3 活性／不活性 template 獲得

我々は活性／不活性 template を、自動構築した template ネットワークと人手で用意した少数の seed template により獲得する。ネットワークのノードは活性か不活性の template である。（本獲得法では全ての template を活性か不活性と見なす。獲得後、活性／不活性値の絶対値の小さいものを中立と見なす。）リンクはその両端の template の極性が同じか反対かを示す。本手法全体は template ネットワーク構築と活性／不活性値の計算を交互に行う、ブートストラップ的手法である。

### 3.1 活性／不活性 template の性質

- (1) 動脈硬化を引き起こす 〜ので 脳梗塞を発症する  
（「を引き起こす」「を発症する」とともに活性。）
- (2) 動脈硬化を防ぐ 〜ので 脳梗塞を免れる  
（「を防ぐ」「を免れる」とともに不活性。）
- (3) 動脈硬化を引き起こす 〜が 脳梗塞を免れる  
（「を引き起こす」は活性。「を免れる」は不活性。）
- (4) 動脈硬化を防ぐ 〜が 脳梗塞を発症する  
（「を防ぐ」は不活性。「を発症する」は活性。）
- (5) 抗癌剤を使用する 〜ので 癌が治る  
（「を使用する」は活性。「が治る」は不活性。）
- (6) 抗癌剤を使用する 〜が 癌が治らない  
（「を使用する」「が治らない」とともに活性。）
- (7) (非文) 動脈硬化を防ぐ 〜が 脳梗塞を免れる  
（「を防ぐ」「を免れる」とともに不活性。）

図 1: ネットワーク構築時の制約の例：(動脈硬化, 脳梗塞) は正の関係の名詞対、(抗癌剤, 癌) は負の関係の名詞対

	極性同一	極性反対	正の関係	負の関係	他
順接	正の関係	負の関係	極性同一	極性反対	—
逆接	負の関係	正の関係	極性反対	極性同一	—

表 1: 名詞対の制約 (左) と template 対の制約 (右)

本手法の出発点は、一文中で共起する 2 つの template とそれらを繋ぐ接続詞、2 つの template と当該文中で係り受け関係を結ぶ名詞対の間に存在する、我々が発見した意味的な制約である。表 1 にその制約を示し、図 1 で具体例を示す。我々はまず、あるタイプの名詞対が template の極性と接続詞に密接な関係を持つことを確認した。そうしたタイプの名詞対を以下では、正の関係もしくは負の関係と呼ぶことにする。正の関係の名詞対は、順接接続詞（「〜ので」等）で繋がる活性／不活

性の極性が同一の template 対か、逆接接続詞（「～が」等）で繋がる極性が反対の template 対のみと一文中で係り受け関係を結ぶ（図 1 の (1)-(4) を参照）。正の関係の名詞対の典型例は、(動脈硬化, 脳梗塞) 等の因果関係を表すものや (牛肉, 牛丼) 等の材料関係を表すもの等である。上記制約に違反し、なおかつ自然であるような文は考えにくい。実際、因果関係や材料関係などを表す正の関係の名詞対が、上記制約に違反する template 対と接続詞の組と共起した場合、図 1(7) のような不自然な文になることが殆どであり、正の関係の名詞対、接続詞、template の極性は一種の意味的制約をお互いに課すことが分かる。一方、負の関係の名詞対は、順接接続詞で繋がる活性／不活性の極性が反対の template 対か、逆接接続詞で繋がる極性が同一の template 対のみと一文中で係り受け関係を結ぶ（図 1 の (5)-(6)）。負の関係の名詞対には、(抗癌剤, 癌) 等の「負の因果関係」と言える関係を表すものが含まれる。このような「負の因果関係」を示す負の関係の名詞対が上記の制約を破るような自然な例文の一部になると想像するのは難しい。

以上の制約を表 1 左に要約する。名詞対から見たこの制約を template 対から見ると表 1 右のようになる。例えば、一文中で、正の関係の名詞対と係り受け関係を結び、順接で繋がれている 2 つの template の極性は同一でなくてはならない。表 1 の制約は、名詞対の正負がわかればそれと共起する template 対の極性の異同がわかり、template 対の極性の異同がわかればそれと共起する名詞対の正負がわかる、ということを意味する。本研究では、この template の極性の異同を用いて、具体的な極性を計算することになる。なお、接続詞は予め人手で順接と逆接に分類されているものとする。

## 3.2 template 獲得の全体像

3.1 節の制約に基づき、template 間に「極性同一」（二つの template が同じ極性を持つことを示す）か「極性反対」（二つの template が反対の極性を持つことを示す）のリンクを張ることで template ネットワークを構築する。しかし実際の計算では、名詞対の正負か template 対の極性の異同がわからないとネットワーク構築は始まらない。そこで、人手で活性と不活性それぞれの seed template を少数用意し、それを手がかりとして正負の名詞対を、さらにそれを手がかりとして seed 以外の template 対を獲得する。その結果から template ネットワークを構築し、そのネットワークに 3.3 節で述べる手法（高村法 [2]）を適用して各 template の活性／不活性値を計算する。計算の結果得られた活性／不活性値の絶対値の高い template を最初に人手で用意した seed template に追加した上で、再度、template ネットワークの構築を行い、再度、活性／不活性値の計算を行う。このブートストラップによる template 獲得の全体像を図 2 に示す。提案手法はこのブートストラップを M 回繰り返して終了する。M は予備実験の結果に基づき 7 とした。

1. 人手で seed template を用意し、活性／不活性値として  $+1/-1$  を付与。活性、不活性それぞれ 36 個、10 個を用意した。
2. 2 つの seed template の全ての可能な組み合わせを生成。
3. seed template の対と順接または逆接の接続詞と共に共起する名詞対をコーパスから抽出し、表 1 左に従い正の関係か負の関係に分類。正の関係としてのみ出現する名詞対、負の関係としてのみ出現する名詞対だけを各々正の関係名詞対、負の関係名詞対として残す。但し出現頻度が F 回以下なら除外。F は 5 とした。
4. 正の関係または負の関係名詞対と順接または逆接の接続詞と共に共起する (seed 以外の) template 対を抽出。表 1 右に従い、各 template 対を極性同一と極性反対に分類。template 対が極性同一、極性反対両方で出現している場合は多数決でどちらか決定。
5. 上記 template 対からネットワークを構築。リンクには極性同一、極性反対のラベルを付与。但し、リンクで繋がれる template の数が D 未満の template は除外。D は 5 とした。
6. 高村法をネットワークに適用し活性／不活性値を計算。
7. 活性／不活性値上位、下位それぞれの  $N \times (i-1)$  個の template を上記の結果から取得 ( $N$  は 30 とした。 $i$  は iteration 回数)。これらは活性／不活性値  $+1/-1$  が付与された上で、次の iteration の seed として (元の seed と共に) 使われる。Step 2 へ。

図 2: ブートストラップによる template 獲得

## 3.3 活性／不活性値計算の詳細

template の活性／不活性値は高村法により計算される。高村法は物理学におけるスピンモデルに依拠する。スピンモデルでは、電子は正 (up) か負 (down) のスピンを取り、電子のネットワーク上で次のエネルギー関数が定義され、そのエネルギーを最小化することで具体的なスピンの値が求まる。

$$E(\mathbf{x}, W) = -1/2 \times \sum_{ij} w_{ij} x_i x_j$$

$x_i$  と  $x_j$  ( $\in \mathbf{x}$ ) は電子  $i, j$  のスピんで、行列  $W = \{w_{ij}\}$  は各電子間のリンクの重みを表す。我々は template を電子、活性／不活性値をスピン（活性を up、不活性を down）と見なす。リンクの重みは template 間のリンクのタイプに基づいて以下のように決定されるが、これは、エネルギーを最小化すると、 $w_{ij}$  の値が正なら  $x_i$  と  $x_j$  は同じ極性に、負なら  $x_i$  と  $x_j$  は反対の極性になる傾向があることに基づいている。

$$w_{ij} = \begin{cases} 1/\sqrt{d(i)d(j)} & \text{if } (i, j) \text{ は極性同一} \\ -1/\sqrt{d(i)d(j)} & \text{if } (i, j) \text{ は極性反対} \end{cases}$$

$d(i)$  は template  $i$  とリンクする template の数である。高村法でエネルギー関数を最小化することにより、template の活性／不活性値が決まる。活性／不活性の初期値として、seed template は  $+1/-1$  を、それ以外は 0 を与えられる。我々は高村法の実装としては SUPPIN (<http://www.lr.pi.titech.ac.jp/~takamura/pubs/SUPPIN-0.01.tar.gz>) を用いた。

## 4 矛盾獲得

活性／不活性を用いて矛盾句対を獲得する。獲得する矛盾句対は「癌を破壊する  $\Leftrightarrow$  癌を進行させる」のように、1 つの名詞と極性が反対の template 対（「を破壊する」は不活性、「を進行させる」は活性）から成る。

活性／不活性について著しい対立があり（つまり、極性が反対の template 対の活性／不活性値の絶対値の大きくて）、分布類似度の高い template 対から成る句対ほど矛盾の可能性が高いと仮定し、次のスコアで矛盾句対を順位付ける。

$$Ct(p_1, p_2) = |s_1| \times |s_2| \times sim(t_1, t_2)$$

$p_1$  と  $p_2$  は矛盾句対で、 $t_1$  と  $t_2$  は  $p_1$  と  $p_2$  の template、 $|s_1|$  と  $|s_2|$  は  $t_1$  と  $t_2$  の活性／不活性値の絶対値である。 $sim(t_1, t_2)$  は  $t_1$  と  $t_2$  の分布類似度 [1] である。

## 5 因果関係獲得

活性／不活性を用いて「需要が減る ⇒ 失業が増える」等の因果関係句対を獲得する。獲得する因果関係句対は一文中で順接で繋がれた template 対と、それと同一文中で共起する名詞対から成る。我々は、ある名詞の指す対象（「需要」等）の主たる機能、効果、目的、役割、影響が強く活性化／不活性化される程同一文中で順接接続詞を介して共起する別の名詞の指す対象（「失業」等）が強く影響される傾向があり、この傾向が強いほど（名詞対と共起する template 対の活性／不活性値の絶対値が大きい程）、当該句対が因果関係である可能性が高いと仮定し、次のスコアで因果関係句対を順位付ける。

$$Cs(p_1, p_2) = |s_1| \times |s_2|$$

$p_1$  と  $p_2$  は因果関係をなす句対、 $|s_1|$  と  $|s_2|$  はその template 対の活性／不活性値の絶対値である。

なお、活性／不活性の有効性が不明確になる恐れがあるため、評価実験の際、順接接続詞のうち「～ので」等因果関係を明示するものと共起する句対は除外した。

## 6 因果関係仮説生成

獲得した因果関係句対と矛盾句対から、コーパス中に記載のない因果関係を仮説として自動生成する。我々は、獲得した因果関係「 $p_1 \Rightarrow p_2$ 」（「需要が減る ⇒ 失業が増える」等）が妥当なら、 $p_1$ 、 $p_2$  をそれと矛盾する句  $q_1$ 、 $q_2$  で置換したもの（「需要が増える ⇒ 失業が減る」等）は妥当な因果関係仮説である可能性が高いと仮定し、次のスコアで因果関係仮説を順位付ける。

$$Hp(q_1, q_2) = Ct(p_1, q_1) \times Ct(p_2, q_2) \times Cs'(p_1, p_2)$$

$q_1$  と  $q_2$  は因果関係仮説をなす句対、 $p_1$  と  $p_2$  は元の因果関係句対、つまり「 $p_1 \Leftrightarrow q_1$ 」と「 $p_2 \Leftrightarrow q_2$ 」は矛盾句対で、 $Ct(p_1, q_1)$  と  $Ct(p_2, q_2)$  はその矛盾スコアである。 $Cs'(p_1, p_2)$  は元の因果関係のスコアである。これは 5 節の  $Cs$  でもよいが、本稿では次式に基づく。

$$Cs'(p_1, p_2) = |s_1| \times |s_2| \times npfreq(n_1, n_2)$$

$npfreq(n_1, n_2)$  は、因果関係仮説中の名詞対  $(n_1, n_2)$  が正の関係なら極性が同じ template 対と共起する頻度、負なら極性が異なる template 対と共起する頻度である。

なお、文内共起している（コーパス中で記載されている）因果関係仮説の句対は出力から除外する。

## 7 評価実験

以下の実験では、公平を期するため、全手法の出力を shuffle した上で著者以外の評価者 3 名が評価し、最終的な判定は多数決で決めた。また評価サンプルには seed template から成るものは含めないようにした。コーパスは KNP (<http://nlp.ist.i.kyoto-u.ac.jp/index.php?KNP>) で解析した Web 6 億文書を用いた。

### 7.1 活性／不活性 template 獲得

提案手法 (Proposed) により活性 template (活性／不活性値  $> 0$ ) を 8,685 個、不活性 template (活性／不活性値  $< 0$ ) を 2,140 個獲得した。評価には活性、不活性それぞれ 100 個のランダムサンプルを用いた（比較対象の他の手法も同様）。評価には活性、不活性、中立の 3 ラベルを用い、中立と判定されたものは活性、不活性いずれの評価でも不正解としてカウントされる。

図 3 は適合率のグラフである。Proposed は、活性の

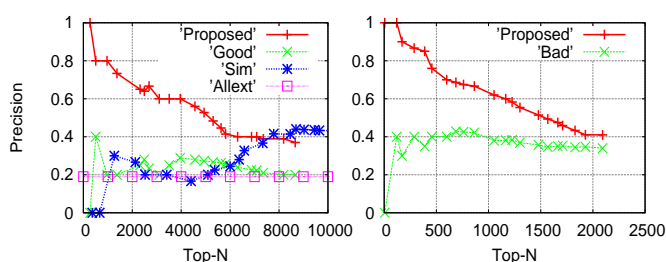


図 3: template 獲得の適合率: 活性 (左) と不活性 (右)

場合上位 2,000 で、不活性の場合上位 500 で 70% 以上の適合率を達成した。ベースラインとして、活性 template が不活性 template を数で上回るという観察に基づき、全 template を活性と見なす手法を評価した (Allert)。結果、適合率は約 20% に留まった。また、紙面の都合上詳細は割愛するが、分布類似度 [1] の高い template を極性同一としてリンクすることでネットワークを構築する手法も評価した (Sim)。結果、不活性 template が全く得られず、活性の適合率も上位 2,000 で約 30% に留まった。これは、類似度は高いが極性が反対の template 対を極性同一としてリンクしてしまうのが主な原因である。以上から、提案手法が高精度であると結論づけられる。さらに、活性を good、不活性を bad と見なして評価した。結果、いずれの場合も適合率が低いため、活性／不活性と評価極性の評価における評価者間の Fleiss' kappa はそれぞれ 0.48 と 0.6 だった (moderate な一致)。

### 7.2 矛盾獲得

提案手法で獲得した活性 template 上位 2,000 と不活性 template 上位 500 を用いて、4 節の提案手法で矛盾句対を獲得した (Proposed)。そのスコア上位 100 万から 200 のランダムサンプルを取得して評価した。比較手法の Proposed-na と Wordnet も同様。Random は

100 サンプルで評価。) 図 4 は適合率のグラフである。Proposed は上位 100 万で 70%以上を達成したため、提

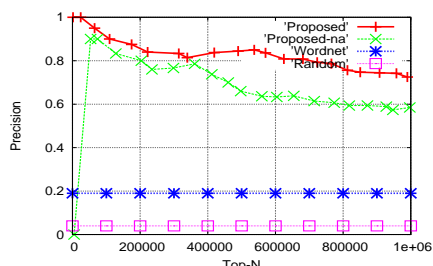


図 4: 矛盾獲得の適合率

案手法は高精度であると結論づけられる。活性／不活性値が矛盾獲得に有効であることを示すため、順位付けに活性／不活性値を使わない以外 Proposed と同じ手法 Proposed-na を評価した (つまり、template の極性を反対にするという条件は課すが、順位付けは  $sim(t_1, t_2)$  のみで行う)。上位 100 万で 10%程適合率が下がったため、活性／不活性値は矛盾獲得に有効であると言える。句対をランダムに生成する手法は適合率が約 4%だった (Random)。他、WordNet の反義関係を使った手法を評価した (Wordnet)。詳細は割愛するが、日本語 WordNet には反義関係がまだ付与されていないため、英語版の反義関係を日本語版に移植して評価したが、適合率は約 20%だった。Proposed の 200 サンプル中の template 対の異なり数は 194 だった。つまり多様な矛盾が獲得できた。表 2 は提案手法が獲得した矛盾の例である。評価者

103,581	「運転を助ける ⇔ 運転を妨げる」
317,028	「騒音がひどくなる ⇔ 騒音は減少する」
487,496	「痛みが発症する ⇔ 痛みを減らす」
529,173	「アクセスが生ずる ⇔ アクセスを抑制する」
848,331	「ガンを破壊する ⇔ ガンを進行させる」
982,980	「ウイルスが死滅する ⇔ ウイルスが活性化する」

表 2: 提案手法で獲得した矛盾の例とその順位

間の Fleiss' kappa は 0.78 だった (substantial な一致)。

### 7.3 因果関係獲得

5 節の提案手法の結果上位 100 万から 100 個のランダムサンプルを取得して評価した (Proposed)。(比較手法も同様。) 評価者には因果関係抽出元の文も提示した。図 5 は適合率を示す。Proposed は上位 50 万で約

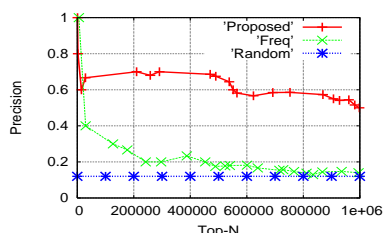


図 5: 因果関係獲得の適合率

70%を達成した。ランダムな句対の適合率 (Random)

は約 12%だった。Freq は句対をその出現頻度で順位付ける以外 Proposed と同じだが、適合率は上位 50 万で約 20%だった。以上から提案手法が高精度であると結論づけられる。Proposed の 100 サンプル中の template 対異なり数は 91 であり、多様な因果関係が獲得できた。表 3 は提案手法で得た因果関係の例である。評価者間の Fleiss' kappa は 0.68 だった (substantial な一致)。

20,819	「体脂肪を分解する ⇒ エネルギーを生む」
154,762	「免疫力を高める ⇒ 疲労回復を促進する」
174,325	「採光性が高まる ⇒ 開放感がアップする」
206,609	「揚力が下がる ⇒ 機首が下がる」
644,966	「血行を促進する ⇒ 新陳代謝を助ける」
829,988	「体力が落ちる ⇒ 抵抗力が減少する」

表 3: 提案手法で獲得した因果関係の例とその順位

### 7.4 因果関係仮説生成

5 節の提案手法で得た因果関係上位 10 万を元に生成した因果関係仮説の上位 100 万からランダムに 100 個サンプリングし評価した。評価者には元の因果関係とその抽出元文も提示した。結果、上位 100 万で約 57%の適合率を得た。100 サンプル中の template 対異なり数は 99 だった。表 4 に提案手法が生成した仮説の例を挙げる。Fleiss' kappa は 0.51 だった (moderate な一致)。

18,886	「ストレスが減少する ⇒ 不眠が改善される」 (「ストレスが増加する ⇒ 不眠が続く」)
205,486	「犯罪を減らす ⇒ 不安が無くなる」 (「犯罪が増加する ⇒ 不安が高まる」)
450,353	「需要が拡大する ⇒ 失業を減少させる」 (「需要が減る ⇒ 失業が増える」)
512,592	「店は減る ⇒ 活気が減る」 (「店が増える ⇒ 活気がある」)
874,036	「消費をもたらす ⇒ 売上を増やす」 (「消費が減少する ⇒ 売上を減少させる」)

表 4: 提案手法で生成した因果関係仮説とその順位 (括弧内は仮説生成元の因果関係)

## 8 おわりに

NICT で開発中の WISDOM (wisdom-nict.jp) 等の情報分析システムにとって、情報の間の矛盾や、ある出来事の原因や結果を自動検出する機能は重要である。我々が獲得した矛盾と因果関係、因果関係仮説は今後 WISDOM に組み込む予定である。その結果、「TPP 締結の結果製造が難しくなると予想される製品は何か、その予想の根拠とされる情報と矛盾する情報はないか」等の自動検出の性能が飛躍的に向上すると期待できる。

## 参考文献

- [1] Dekang Lin. Automatic retrieval and clustering of similar words. In *COLING-ACL*, 1998.
- [2] Hiroya Takamura, Takashi Inui, and Manabu Okumura. Extracting semantic orientation of words using spin model. In *ACL*, 2005.