

# 事象の意味と時間関係の偏りを考慮した時間的順序関係推定

高木宏伸<sup>†</sup>嶋田和孝<sup>‡</sup><sup>†</sup>九州工業大学大学院 情報工学府<sup>‡</sup>九州工業大学大学院 情報工学研究院

{h\_takagi, shimada}@pluto.ai.kyutech.ac.jp

## 1 はじめに

自然言語処理のタスクの一つに時間情報解析がある。文書中のある2つの事象についてそのどちらが先に起こったか、または同時に起こったかを推定する時間的順序関係推定も時間情報解析に含まれる。このような事象の発生した時間やそれらの時間関係を同定することは、因果関係知識の獲得や含意関係認識への応用など、深い言語理解の実現に必要である。

海外における時間情報解析については、時間情報が付与されたコーパスである TimeBank[1] が公開され、TempEval[2] と呼ばれるワークショップを中心に、英語、スペイン語、中国語といった主要な言語についてさまざまな研究が行われている [3, 4]。しかしながら、日本語を対象とした研究はまだ少ない。

そのような背景の中、近年日本語での時間情報が付与された BCCWJ-TimeBank[5, 6] が構築され、機械学習による統計的解析や、定量的な性能評価が可能になった。しかしながら、現時点においてコーパスのデータ量は十分とはいえず、1事例のみの表現も多く存在する。したがって、表現を特徴量に落とし込む際には抽象化やグループ化などの処理が必要になる。そこで本稿では、事象の表層表現に頼らない意味と外部の知識を利用した時間関係の偏りに関する素性を提案し、その有効性について報告する。

## 2 関連研究

英語での時間情報が付与された TimeBank は、TimeML 準拠の時間表現と事象表現、及びそれらの時間的順序関係について13種のラベルが付与されている。この TimeBank を用いて、Bethard[7] は時間表現、事象表現の抽出、またそれらの時間的順序関係を推定する機械学習モデルを提案している。このモデルでは、事象の品詞や前後の単語、構文木での位置などといった語彙的・統語的素性を用いることで高精度の

結果を報告している。

BCCWJ-TimeBank を用いた時間的順序関係推定手法としては、吉川ら [8] が英語での時間的順序関係推定に一般的に使われる素性を日本語で使用できるよう整備し、これらの素性が日本語での時間的順序関係推定においても有用であることを示した。また、稲田ら [9] は、同一文内の事象間に焦点を当てたモデルを提案し、類義語関係や大規模データの頻度情報を用いた素性なども使用することで、高精度の結果を報告している。

## 3 使用データと研究対象

時間的順序関係分類モデルを作成するにあたり、事象間の時間関係が付与された BCCWJ-TimeBank を使用する。BCCWJ-TimeBank は、「現代日本語書き言葉均衡コーパス<sup>1</sup>」のコアデータ54文書に対して、時間表現を示す「TIMEX3」、事象表現を示す「EVENT」、それらの間の時間関係を示す「TLINK」が付与されている。

時間関係を示す「TLINK」には、隣接する事象間 (E2E)、隣接文の主述部間 (MAT)、時間表現と事象間 (T2E)、文書作成日時と事象間 (DCT) の4種類の関係について「before, after, overlaps」など17種類のラベルが3人の作業員によって付与されている。

本稿では、事象間の時間的順序関係推定 (E2E) を行うが、一文内の範囲を超えた事象対では、文書全体の文脈情報が必要になるなど取り扱う問題が異なってくるため、今回は同一文内の事象対のみを対象とする。また、時間関係を示すラベルの一部はわずかしき出現せず、疎データ問題を引き起こす。そのため、17種類のラベルを表1に示す4種類に簡略化する。また、実験では、作業員3人の付与したラベルが一致したもののみを用いる。現時点で同一文内の事象対への時間関係が付与されているのは1241事例である。表2に事

<sup>1</sup>[http://www.ninjal.ac.jp/corpus\\_center/bccwj/](http://www.ninjal.ac.jp/corpus_center/bccwj/)

表 1: 時間的順序関係ラベルの簡略化

元ラベル	変換後
after	AFTER
met-by	AFTER
before	BEFORE
meets	BEFORE
contains	OVERLAP
during	OVERLAP
equal	OVERLAP
finished-by	OVERLAP
finishes	OVERLAP
identity	OVERLAP
includes	OVERLAP
is_included	OVERLAP
overlaped-by	OVERLAP
overlaps	OVERLAP
started-by	OVERLAP
starts	OVERLAP
vague	VAGUE

表 2: 時間的順序関係が付与された事象対数

ラベル	事象対数
AFTER	212
BEFORE	601
OVERLAP	387
VAGUE	41
合計	1241

例の内訳を示す。

## 4 提案手法

同一文内の事象対に対して、時間的順序関係を出力する多値分類問題として捉え、教師あり学習による分類モデルを作成する。学習器としてはSVMを利用する。SVMは多くの時間情報解析研究でも利用されており、再現が容易なことに加え、性能が高いことが知られている。表3に用いる素性一覧を示す。4.1節で分類モデルの基本となる素性について述べ、4.2節及び4.3節で新たに提案する素性について説明する。

### 4.1 表層的な素性

表3のうち、1~7は英語における時間関係認識の代表的なシステムが使用している基本的な素性や、日本語における研究において有効性が報告されている素性を既存の解析器を用いて作成したものである。

表 3: 使用した素性

id	素性
1	事象と前後2形態素の基本形
2	事象と前後2形態素の品詞
3	事象と前後2形態素の活用形
4	BCCWJ-TimeBankのclass属性
5	事象間の係り受け関係の有無
6	事象間の係り受け関係の距離
7	事象間の時間表現の有無
8	事象が含まれる節の種類
9	事象の語尾
SEM	事象の意味役割(大分類)
	事象の意味役割(中分類)
	事象の意味役割(小分類)
LTV	事象の時間的前後関係の頻度
	事象の時間的同時関係の頻度

1~3は事象表現に関する形態素の基本的な素性であり、4はBCCWJ-TimeBankの付加属性を用いた素性である。また、5~7は英語・日本語問わず有効性が示されている事象間の構文的関係についての素性である。これらの素性については、Mecab<sup>2</sup>及びCabocha<sup>3</sup>の解析結果を使用する。

8は日本語独特な表層の特徴として、我々がエラー分析を行った際[10]に事象が含まれる従属節の種類(補足節・連体節・副詞節・並列節)と時間的順序関係と相関が見られたため採用する。

9について、日本語における時間的順序関係推定の研究のひとつである吉川らのモデルでは、事象の時制・相を捉えるために事象表現の語尾が「る、た、ている、ていた」であるかという素性を用いている。我々のモデルでも、この素性を採用する。

### 4.2 意味に関する素性

時間的順序関係推定において表面的な構造のみからでは識別が困難な場合がある。

- (a) ホンジュラス戦を 控え<sub>E1</sub>、練習に 臨む<sub>E2</sub>。  
overlap(E1,E2)
- (b) 死球を 受け<sub>E1</sub>、ベンチに 退く<sub>E2</sub>。  
before(E1,E2)

この例は、表層構造は全く同じであるが、(a)では“控える”という状態のまま練習に“臨む”ためこの2

<sup>2</sup><http://taku910.github.io/mecab/>

<sup>3</sup><http://taku910.github.io/cabocha/>

表 4: 時間的順序関係抽出の手がかり表現

時間関係	手がかり表現
AFTER	前, まえ
BEFORE	後, あと, 結果
OVERLAP	ながら, 同時

対の事象間の時間関係は重なっている。しかしながら (b) では、死球を“受け”た結果、ベンチに“退く”という (継起/因果) の意味を持ち、この 2 対の事象間の時間関係は前後の関係である。これらは表層構造からでは区別ができず、また表層表現を素性とするには学習データ量が非常に乏しいため、従来の手法ではこれらの識別を行うことが難しかった。

そこで、表層からでは読み取れない動作の意味を述語項構造シソーラス [11] を用いて付与し、異なった表現でも同じ意味の事象をグループ化して学習させることを試みる。具体的には、BCCWJ-TimeBank の事象に対し、述語項構造シソーラスに記述されている大分類 1・2, 中分類, 小分類 1・2 をそれぞれ素性とする。

#### 4.3 時間関係の偏りに関する素性

表 2 に示した通り、時間関係認識に使用できるデータ数は十分ではない。そこで、外部の知識から事象ごとの時間関係に関する偏りを取得することを試みる。例を挙げて説明すると、「笑う」という事象は、何か要因があった上で起きる事象であり、その要因の事象に対し時間関係は後に偏りやすく、また、「見ながら笑う」「言いながら笑う」など他の事象と同時に起こりやすいといった偏りが存在する。これを数値化することによって、事象そのものが他の事象とどのような時間関係を持ちやすいかを表す素性とする。

まず、時間関係を判断する手がかりとなるフレーズを表 4 のように定める。次に Twitter 上から事象間に係り受け関係があり、また事象表現間にフレーズを含む文章を *Label*(事象 1, 事象 2) と取得する。例えば、「弁当を食べる前に手を洗った」という文章は *Before*(食べる, 洗う) となる。次に、ある事象  $X$  と対に獲得された全ての事象  $y$  について出現数の総和 ( $Sum_B(X), Sum_A(X), Sum_O(X)$ ) を取り、時間関係の前後の偏り (Latent temporal value:  $LTV_{before-after}$ ) 及び同時の偏り ( $LTV_{overlap}$ ) を求め、この値を素性として用いる。

$$Sum_B(X) = \sum_{y \in D} (Before(X, y) + After(y, X))$$

$$Sum_A(X) = \sum_{y \in D} (Before(y, X) + After(X, y))$$

$$Sum_O(X) = \sum_{y \in D} (Overlap(y, X) + Overlap(X, y))$$

$$LTV_{before-after}(X) = \frac{Sum_A(X)}{Sum_B(X) + Sum_A(X)}$$

$$LTV_{overlap}(X) = \frac{Sum_O(X)}{Sum_B(X) + Sum_A(X) + Sum_O(X)}$$

$LTV_{before-after}$  は 0~1 の値をとり、0 に近づくほど他の事象より先に起きやすいことを表し、逆に 1 に近づくほど他の事象よりも後に起きやすいことを表す。また、 $LTV_{overlap}$  についても 0~1 の値をとり、0 に近づくほど他の事象とは同時に起こりにくいことを表し、1 に近づくほど同時に起きやすいことを表す。

表 4 のフレーズを元に Twitter から約 30 万事象対 (約 1 万種の事象) を獲得し、BCCWJ-TimeBank の事象との被覆率は 9 割以上となった。

## 5 実験と結果

英語の時間的順序関係推定に用いられる素性の日本語での有効性の確認、また新たに提案した素性の有効性の確認のために実験を行った。実験には表 2 のデータに対して 10 分割交差検定を行い、各ラベルの F 値とその micro 平均により評価した。

表 5 は表 3 の中で表層的な素性である 1~9 を用いたモデルに、意味の素性 (SEM) と頻度情報 (LTV)、その両方を加えたモデルの結果である。また、表におけるラベルは A(AFTER), B(BEFORE), O(OVERLAP), V(VAGUE), micro(micro 平均) である。すべてのモデルにおいて VAGUE の精度が非常に低いが、これは事例数が 41 件と少なく学習ができなかったためと考えられる。最も性能が良かったのは表層的な素性 1~9 に意味の素性 SEM を加えたモデルであった。事象が表す意味を捉えることが一定の精度向上に繋がったと考えられる。頻度の素性 LTV を加えたモデルについては、AFTER の F 値が上昇し、事象対の前後関係の識別に有効に働いているといえる。しかしながら、意味役割の素性と組み合わせても最高の性能を示すことはできなかった。これは、時間的順序関係推定において事象の頻度による偏りよりも構文的な特徴や語尾の変化などといった素性のほうが有効に働くためだと考えられる。

そこで、意味役割と頻度情報に関する素性について、4.2 節の例文のような構文的な特徴や語尾の変化などの表層的特徴からでは時間的順序関係推定が困難な場合について有効に作用するか検証した。実験対象は BCCWJ-TimeBank のうち、事象対が第一事象の

表 5: 提案素性の評価結果

素性	A	B	O	V	micro
1-9	0.524	0.774	0.561	0.0	0.655
1-9+SEM	0.537	0.782	0.567	0.083	<b>0.663</b>
1-9+LTV	0.556	0.772	0.562	0.0	0.657
1-9+SEM+LTV	0.537	0.776	0.577	0.0	0.661

表 6: 中立形・テ形接続を対象とした評価結果

素性	A	B	O	V	micro
1-9	0.0	0.757	0.432	0.0	0.636
1-9+SEM	0.0	0.792	0.500	0.0	0.679
1-9+LTV	0.0	0.765	0.462	0.0	0.657
1-9+SEM+LTV	0.0	0.792	0.511	0.0	<b>0.682</b>

中立形またはテ形によって接続しているもの 310 件とする。この 310 件の事例の内訳は、AFTER:5 件、BEFORE:210 件、OVERLAP:88 件、VAGUE:7 件と BEFORE と OVERLAP の比率が全体の 96% を占めているが、表層構造が似通っているため、表層的な特徴からでは BEFORE と OVERLAP の判別が難しい事例が多い。結果は表 6 に示す通り、表層的な素性みのモデルに比べて、意味役割と頻度情報両方を追加したモデルが BEFORE・OVERLAP 共に F 値で高い精度を出した。この結果から、文構造など表層的な特徴が似ていることから時間的順序関係推定が難しいケースにおいては、意味役割・頻度情報共に識別に有効に働くといえる。

最後に、本稿で提案したモデルと日本語での先行研究である稲田らのモデルの比較を行う。稲田らとは実験対象とするデータの事例数が異なる点から直接比較できないが、正答率での比較を行うと、我々のモデルが 66.87% に対して稲田らは 68.18% と及ばなかった。しかしながら、本稿で提案した素性を稲田らのモデルと組み合わせることによってさらなる精度向上が期待できると考えられる。

## 6 おわりに

本稿では、時間的順序関係推定において、文の表層から獲得できない事象の意味役割や時間的順序関係への偏りを新たな素性として提案した。その結果、大幅な精度向上には繋がらなかったが、表層的な特徴のみでは推定が困難なケースにおいては有効に働くことが確認できた。

今後は、エラー分析を詳細に行い、日本語における

独特な素性なども新たに作成していきたい。また、今回用いた BCCWJ-TimeBank は事例数が少なかったため、別の文書に対しての評価実験や精度調査なども今後の課題として挙げられる。

## 謝辞

本研究は JSPS26730176 の助成を受けたものです。

## 参考文献

- [1] James Pustejovsky, Jose Castano, Robert Ingria, Reser Sauri, Robert Gaizauskas, Andrea Setzer and Graham Katz. The timebank corpus. In *Proceedings of Corpus Linguistics 2003*, pp. 647–656, 2003.
- [2] Naushad UzZaman, Hector Llorens, Leon Derczynski, James Allen, Marc Verhagen, and James Pustejovsky. Semeval-2013 task 1: Tempeval-3: Evaluating time expressions, events, and temporal relations. In *Second Joint Conference on Lexical and Computational Semantics (\*SEM), Volume 2: Proceedings of the Seventh International Workshop on Semantic Evaluation (SemEval 2013)*, pp. 1–9, 2013.
- [3] Marc Verhagen, Roser Sauri, Tommaso Caselli, and James Pustejovsky. Semeval-2010 task 13: Tempeval-2. In *Proceedings of the 5th International Workshop on Semantic Evaluation (SemEval 2010)*, pp. 57–62, 2010.
- [4] Kenton Lee, Yoav Artzi, Jesse Dodge and Luke Zettlemoyer. Context-dependent Semantic Parsing for Time Expressions. In *Proceedings of the Conference of the Association for Computational Linguistics*, pp. 1437–1447, 2014.
- [5] 小西光, 浅原正幸, 前川喜久雄. 『現代日本語書き言葉均衡コーパス』に対する時間情報アノテーション. 自然言語処理, Vol. 20, No. 2, pp. 201–222, 2013.
- [6] 保田祥, 小西光, 浅原正幸, 今田水穂, 前川喜久雄. 『現代日本語書き言葉均衡コーパス』に対する時間表現・事象表現間の時間的順序関係アノテーション. 自然言語処理, Vol. 20, No. 5, pp. 657–682, 2013.
- [7] Steven Bethard. Cleartk-timeml: A minimalist approach to tempeval 2013. In *Second Joint Conference on Lexical and Computational Semantics (\*SEM), Volume 2: Proceedings of the Seventh International Workshop on Semantic Evaluation (SemEval 2013)*, pp. 10–14, 2013.
- [8] 吉川克正, 浅原正幸, 飯田龍. BCCWJ-TimeBank を対象とした時間的順序関係の推定. 言語処理学会第 20 回年次大会発表論文集, pp. 1103–1106, 2014.
- [9] 稲田和明, 松林優一郎, 乾健太郎. 日本語文書内で表現される事象間の時間的な順序関係の推定. 情報処理学会論文誌 Vol.56 No.10, pp. 2054–2071, 2015.
- [10] 高木宏伸, 嶋田和孝. 事象間の接続関係に基づく時間的順序関係推定. 言語処理学会第 21 回年次大会発表論文集, pp. 55–58, 2015.
- [11] 竹内孔一, 石原靖弘, 竹内奈央, 述語項構造シソーラスによる述語と名詞の構造化. 人工知能学会全国大会, 215-OS-08b-1, 2014.