

文書内に現れる因果関係の出現特性調査

乾 孝司[†] 奥村 学[†]

概要：本稿では、文書内に現れる因果関係の出現特性を調査した結果について述べる。まず、因果関係情報を一定量の文書集合に対して付与し、因果関係タグ付きコーパスを作成する。その後、コーパスに付与された情報を元に、文書内での因果関係の出現傾向を定量的に調査する。因果関係は定義すること自体が困難な概念であり、従来、因果関係の有無の判断には、人間の主観に頼らざるをえなかった。本研究では、因果関係情報を付与する際、言語テンプレートに基づく判断基準を採用し、得られる因果関係情報からタグ付け作業者の主観をできる限り排除することを試みる。出現特性の調査項目として、(i) 手がかり標識の有無、(ii) 出来事表現の統語カテゴリ、(iii) 出来事表現の文内位置、の3つの観点を取り上げ、因果関係が多様な言語表現形態で、かつ、多様な位置に現れることを示す。

キーワード：因果関係，コーパス

Characteristics of In-text Causal Relations

INUI Takashi OKUMURA Manabu

Abstract : In this paper, we report the results of our investigation of the characteristics of in-text causal relations. First, we designed causal relation tags. With our designed tag set, annotators annotated a set of Japanese text documents. Then, using the annotated corpus, we investigated the causal relation instances from three viewpoints: (i) cue phrase markers, (ii) part-of-speech information, and (iii) position in sentences. Our quantitative study shows that what amount of causal relation instances are present, where these relation instances are present, and which types of linguistic expressions are used for expressing these relation instances in text.

Keywords : causal relation , corpus

1 はじめに

因果関係に関する知識は、質疑応答システムや対話システムなど、幅広い自然言語アプリケーションにとって重要な知識のひとつである。近年では、大規模な文書集合から自動的に因果関係知識を獲得する手法が既に幾つか提案されている（例えば、文献 [12, 3, 13, 4, 5]）。

大規模な文書集合を利用して、因果関係知識を高い精度でかつ、効率よく自動獲得するには、文書内において因果関係がどのように現れるかという、因果関係の出現特性をあらかじめ把握しておくことが望ましいと考えられる。しかしながら、現在、知識獲得の方法論に関する研究開発が進展するその一方で、文書内に現れる因果関係の出現特性に関しては、これまでにまとまった知見が得られているとは言い難く、未知な点が多い。

本研究では、このような背景を踏まえ、文書内に現れる因果関係の出現特性を調査することを目標としている。まず、一定量の日本語文書集合を対象として、因果関係情報を注釈付け、因果関係タグ付きコーパスを作成

する。そして、コーパスに付与された情報を元に、文書内での因果関係の出現特性を定量的に調査する。調査の項目として、本稿では特に、(i) 因果関係の手がかり標識の有無、(ii) 出来事表現の統語カテゴリ、(iii) 出来事表現の文内位置の3つの観点に注目する。

上記の観点からの調査を実現するために本研究では、注釈付け対象の表現形式に関する制約は原則として設けない。例えば、(i) の観点から因果関係の出現特性を調査するために、例文 (1a) のような、明示的な手がかり標識を伴う形式に加え、(1b) のような、手がかり標識を伴わない非明示的な表現形式も注釈付けの対象に含める。また、(ii)、(iii) の観点から出現特性を調査するために、(1c) のように、出来事が名詞句として表現される形式、また、(1d) のように、文をまたいで2つの出来事が表現される形式も注釈付け対象に含める。例文 (1) に示すような、幅広い表現形式を対象として因果関係の出現傾向を定量的に調査した目立った既存研究はこれまでに存在しておらず、文書内での因果関係の出現特性を幅広く調査することが本研究の特徴となる。

- (1) a. 大雨が降ったため、川が増水した。（明示的）
b. 大雨が降り、川が増水した。（非明示的）
c. 大雨で川が増水した。（出来事が名詞句）
d. 大雨が降った。川が増水した。（文をまたぐ）

[†] 東京工業大学 精密工学研究所
Precision and Intelligence Laboratory, Tokyo Institute of Technology
連絡先: tinui@lr.pi.titech.ac.jp

本研究では新しい因果関係の判断基準を提案，採用する*（詳細は3.2節）．因果関係は定義すること自体が困難な概念であり，従来から，因果関係の有無を判断するには，人間の主観に頼らざるをえなかった．本研究では，因果関係情報を注釈付ける際，言語テンプレートに基づく判断基準を採用し，得られる因果関係情報からタグ付け作業の主観をできる限り排除することを試みる．言語テンプレートに基づく判断では，作業間で共通の，言語的な判断の拠り所が与えられる．各作業者は，定められた規則に従って，因果関係の有無を判断したい2つの出来事を含んだ言語表現を作成する．そして，その言語表現が構文，意味的に適格かどうかを判断することを通して，注目している2つの出来事間に因果関係があるかないかを判断する．

以下，2節で関連研究について述べ，3節から6節で，作成した因果関係タグ付きコーパスについて述べる．その後，7節，8節で，因果関係の出現傾向に関する定量的調査の結果を述べ，9節で本稿をまとめる．

2 関連研究

因果関係を明示的に表す接続表現に注目した研究として，Liu [7]の研究がある．Liuは，日本語の因果関係を表す接続表現の意味，用法を解明することを目的とし，“だから”など，5つの代表的な接続表現の意味，用法を分析，整理した．この研究から得られた知見を元に，因果関係を表す接続表現の用法の共通点や相違点を把握することができる．しかしながら，Liuの研究では，取り上げられた接続表現が文書内にどれぐらいの頻度で出現するかという数量的な考察はなされていない．

文書集合から因果関係知識を自動獲得する研究 [12, 3, 13, 4, 5] では，例えば，佐藤ら [12] が EDR コーパス [14] 中のテキストや Web 上のテキストデータを対象として，それらの文書に含まれる因果関係の数を評価している．Girju [3] も英語テキストを対象として同様の評価を行っている．しかしながら，因果関係知識の自動獲得に焦点を当てた諸研究では，因果関係の総数のみを評価しており，本研究で取り扱う調査項目のような，より詳細な観点からの分析は行われていない．

文書内の修辭構造を体系的に説明する修辭構造理論 [8] では，修辭関係のひとつとして因果関係に相当する関係がある．Marcu [9] は，Brown コーパスから抽出したテキストを対象にして，修辭構造解析に利用できる手がかり表現（因果関係では，“because” など）を網羅的に収集し，各表現の出現頻度分布を調査した．Marcu の調査は接続詞などに代表される手がかり表現を網羅的に収集しているが，手がかり表現に注目しているため，手がかり表現を伴わない非明示的な因果関係は調査対象とはなっていない．一方，本稿での我々の目的の一つは，先の例文 (1a) のような手がかり表現を伴う場合と，(1b) のような手がかり表現を伴わない場合の両者における，因果関係の出現割合を調査することである．

* 提案する因果関係の判断基準について，その基本的なアイデアは，文献 [5] で因果関係を定義する際に用いた手法と同様のものである．ただし，本稿での判断基準は，3.2 節でも述べるように，因果関係の必然性の強さに関する指標を導入している点で文献 [5] での定義とは異なる．

Altenberg [2] は，英語の会話書き起こしデータと文書データを対象にして，因果関係の出現傾向を調査している．Altenberg の調査は，文体や統語情報，出来事表現の出現順序など，幾つかの観点から因果関係の出現傾向を分析しており，興味深い結果を報告している．西澤ら [11] も，日本語の会話書き起こしデータを対象にして，因果関係をもつ出来事表現の順序や距離といった観点から因果関係の出現傾向を調査している．しかしながら，Altenberg，西澤らの調査はいずれも，先の Marcu の調査と同様に，明示的な表現のみを扱っており，非明示的な因果関係は調査対象から除外されている．西澤らの調査研究は，日本語を対象としており，調査項目も我々の注目点と重なる部分がある．しかしながら，彼らが会話データを扱っているのに対し，本稿では，会話データとは性質の異なる文書データを取り扱う．

3 注釈付ける因果関係情報

3.1 因果関係タグ

因果関係は，2つの出来事間の何らかの依存関係であると一般的には捉えられている．本稿では，主体の意志的な行為（例えば，〈傘をさす〉[†]）と，（単なる状態も含んだ）行為以外のすべての非意志的な事態（例えば，〈雨が降る〉）を合わせて出来事と呼び，それらの出来事間にある因果関係を取り扱う．

因果関係に関する情報は，*head*（主辞要素），*mod*（修飾要素），*causal_rel*（因果関係）の3種類の基本タグを用いて注釈付ける．例えば，次の例文 (2a) には (2b) と (2c) の2つの出来事が含まれ，これらの出来事間に因果関係が存在すると考えられる．この例文 (2a) へのタグ付与結果は図 1 のようになる．ここで，(2) 中の e_1 は因果関係の前件（原因側）出来事であること， e_2 は後件（結果側）出来事であることを示し，図 1 中の各タグの接尾数字は，それぞれの出来事 (e_1 と e_2) の接尾数字に対応する．また，スラッシュ記号は文節区切りを示す．

- (2) a. そして、遠方からの観光客がGWに入って増える。
 b. $e_1 = \langle \text{GWに入る} \rangle$
 c. $e_2 = \langle \text{遠方からの観光客が増える} \rangle$

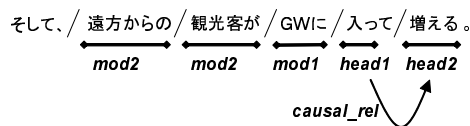


図 1 例文 (2a) へのタグ付与の例

本研究では，出来事は1つの主辞要素に0個以上の修飾要素が付随して構成されると捉え，それぞれの出来事を，1つの *head* タグと0個以上の *mod* タグを組み合わせて注釈付ける．タグ付け作業者は，まず，出来事の主辞となる要素に *head* タグを付与し，出来事を構成する要素が *head* タグの箇所以外にもある場合は，それぞれ

[†] 出来事が表された言語表現と，表現手法とは独立に人によって認識された出来事とを区別して表記するために，人によって認識された出来事を三角括弧で囲って表記する．例えば，後に示される例文 (2a) は，文書内に含まれる言語表現であることを表すが，(2b) や (2c) は，(2a) から認識された出来事であることを表す．

の箇所に *mod* タグを付与することで出来事を注釈付ける。この時、*head* タグと *mod* タグは基本的には文節単位で付与される。このようにして2つの出来事へタグを付与した後、それぞれの *head* タグの間にリンク情報として *causal_rel* タグを付与することで、当該の出来事間に因果関係があることを注釈付ける。ここで、本稿における出来事の主辞とは、ある出来事の意味内容を特定する最も中心的な構成要素を指す。大抵の場合、本稿でいう意味の主辞と一般に句構造文法で定義される構造的な主辞は等しくなるが、両者が異なるケースも存在する。例えば、「昨日買った傘」のような名詞句表現から〈傘を買う〉という出来事を認めてタグを付与する場合には、名詞句の主辞である「傘」ではなく、「買った(買う)」を出来事の意味の主辞として取り扱う。

head タグの情報は、7節で述べる出現特性調査で利用される。具体的には、文内における出来事表現の出現位置に関する調査では、*head* タグを含む文節の位置を出来事表現の位置とみなす。また、出来事表現の統語カテゴリに関する調査では、*head* タグが付与された単語の品詞情報を利用して出来事表現のカテゴリを決定する。

もし、文書内において、因果関係を発見する手がかりとなる表現が存在する場合は、基本タグ以外の付加情報として、該当箇所に *marker* タグを付与する。この情報は、7節において、因果関係の手がかり標識の有無の割合を調査する際に利用される。

3.2 言語テンプレートに基づく判断基準

因果関係は、定義すること自体が困難な概念であり、従来から、因果関係の有無を判断するには、人間の主観に頼らざるをえなかった。本研究では、タグ付け作業者の主観をできる限り排除するために、言語テンプレートに基づく因果関係の判断基準を提案、採用する。

言語テンプレートに基づく判断では、作業者間で共通の、言語テンプレートという言語的な判断の拠り所が与えられる。言語テンプレートの例を図2に示す。テンプレート中の鉤括弧はスロットを表し、実際に判断に利用される際は、因果関係の有無を判断したい2つの出来事間の言語表現がスロットにそれぞれ代入される(テンプレート中の *adv* については後述する)。

『 e_1 』(という)状態になれば、それに伴い、
adv 『 e_2 』(という)状態になる。

図2 言語テンプレートの例

因果関係の有無の判断は次の手順で行なう。まず、事前に言語テンプレートを作成、用意する。今回は、図2を含め、表1に示した18通りのテンプレートを用意した。次に、因果関係の e_1 と e_2 の候補となる2つの出来事表現を文書から選択、抽出し、それらを言語テンプレートの各スロットに代入する。この時、それぞれの言語テンプレートについて、ひとつの文(以下、テンプレート文)が完成する。次に、完成したテンプレート文が構文的に適切であり、かつ、その文意が意味的に適切であるかどうかを考える。もし、あるテンプレート文が構文的、意味的に適切であると判断されれば、スロットに代入した2つの出来事間には因果関係があると判断し、3.1節の記述に従ってタグを付与する。

2つの出来事のうち、どちらが e_1 でどちらが e_2 となるかは、言語テンプレート上で指定されており、それぞれの出来事表現がテンプレートのどちらのスロットに代入されるかによって決定される。もし、18通りのすべてのテンプレート文がいずれにおいても適切と判断されない場合、スロットに代入した2つの出来事間には因果関係がないと判断する。

各スロットに代入される出来事表現は、例えば、「GWに入る」や「遠方からの観光客が増える」のような単文形式を想定している。もし、対象文書から抽出した、因果関係の有無を判断したい出来事表現の主辞要素が活用していれば、基本形に戻した状態でスロットに代入する(「GWに入っ」「GWに入る」)。また、出来事表現の主辞要素が名詞句(NP)として表現されている場合は、次のいずれかの書き換え操作を施した結果をスロットに代入する。

- NP NP + する(「停電」「停電する」)
- NP NP + が起こる(「地震」「地震が起こる」)
- NP NP + になる(「大雨」「大雨になる」)
- 動詞の名詞化 動詞の基本形(「疲れ」「疲れる」)

次に、*adv* について説明する。言語テンプレート中の *adv* は、テンプレートを使用する際に、「しばしば」「大抵」「常に」あるいは「」のいずれかの単語に置き換えられる。「」は説明のために導入した擬似的な記号であり、「」に置き換えることは、*adv* 部分を削除したテンプレートを使用することに等しい。

ここで、上記の3つの副詞は、必然性の強い2つの出来事と共に現れた時のみ語用論的に適切な文として認識される性質をもつ。この性質を利用し、「蓋然」「偶然」という、必然性の強さに関する2つ属性値を *causal_rel* タグに付加することを考える。

「蓋然」 因果関係に立つ出来事間で、必ずそうなると決まっていること、それ以外になりようのないこと。あるいは、ある程度確実であること。

「偶然」 通常、ある2つの出来事の間には因果関係はなく、予測していないことが起こること。

具体的には、18通りの言語テンプレートのそれぞれについて、*adv* を上記の3つの副詞のいずれかで置き換えたテンプレートを用いて因果関係があると判断される場合には、「蓋然」の属性値を付加する。*adv* を3つの副詞のいずれかで置き換えたテンプレートを用いると因果関係があると判断されないが、「」に置き換えた(すなわち、*adv* を削除した)テンプレートでは因果関係があると判断される場合には、「偶然」の属性値を付加する。

周知の通り、因果関係という概念自体が、現在でも哲学を含めた学問領域における議論の対象となっており、多様な事象間関係がその範疇下におかれる。ただ、最終的に獲得される因果関係知識を工学的な応用システムに適用することを念頭においた場合、注釈付けされた因果関係事例の中で、「蓋然」の強さをもつ事例は特に有用性が高いと期待している。

3.3 注釈付けの対象と優先規則

理想的には、文書内のあらゆる出来事に対して、因果関係が含まれるか否かの判断を試みるべきである。し

表 1 言語テンプレート

id	the linguistic templates
1	『 e_{1a} (という)ことが起こるその結果として、 adv 『 e_{2a} (という)ことが起こる。
2	『 e_{1a} (という)状態になれば、それに伴い、 adv 『 e_{2a} (という)状態になる。
3	『 e_{1a} (という)状態になれば、それに伴い、 adv 『 e_{2a} (という)状況になる。
4	『 e_{1a} (という)状態であると、 adv 『 e_{2a} (という)状態である。
5	『 e_{1a} (という)状態であると、 adv 『 e_{2a} (という)状況である。
6	『 e_{1a} (という)ことをする結果、 『 e_{2a} (という)ことが adv 起こる。
7	『 e_{1a} (という)ことをすると、 adv 『 e_{2a} (という)状態になる。
8	『 e_{1a} (という)ことをすると、 adv 『 e_{2a} (という)状況になる。
9	『 e_{1a} (という)ことをすると、 adv 『 e_{2a} (という)状態を保つ。
10	『 e_{2a} (という)ことをするのは、 adv 『 e_{1a} (という)状態の時である。
11	『 e_{2a} (という)ことをするのは、 adv 『 e_{1a} (という)状況の時である。
12	『 e_{1a} (という)状態になる場合、 adv 『 e_{2a} (という)ことをする。
13	『 e_{1a} (という)状況になる場合、 adv 『 e_{2a} (という)ことをする。
14	『 e_{1a} (という)状態では、 adv 『 e_{2a} (という)ことをする。
15	『 e_{1a} (という)状況では、 adv 『 e_{2a} (という)ことをする。
16	『 e_{1a} (という)ことが起こらなければ、 adv 『 e_{2a} (という)ことができない。
17	X が 『 e_{2a} (という)ことを実現する手段として、 adv X が 『 e_{1a} (という)ことを行なう。
18	X が 『 e_{1a} (という)ことをすることによって、 adv X が 『 e_{2a} (という)ことができる。

かしながら、任意の2つの出来事について、それらが文書内で現れる位置が離れるに従い、それらの間に因果関係が成立する確率は減少すると考えられる。そこで今回は、隣接する2文以内に2つの出来事の主辞要素が含まれる場合のみ、因果関係の有無を判断する。この制約は主辞要素に *head* タグ付与する場合のみ適用する。日本語では、照応現象や省略現象が頻繁に文書内に生じ、しばしば主辞要素とその修飾要素が異なる文において表現される。そこで、修飾要素については、それが現れる文の位置に関わらず、*mod* タグの付与を試みる。

実世界においては、 e_1 に対して複数の e_2 が存在したり、複数の因果関係が複雑に絡み合っって認識されるケースも存在する。しかしながら今回は、タグ付与作業にかかる作業者の負荷を軽減することを優先し、ひとつの e_1 にはやはりひとつの e_2 のみを割り当てる。注釈付ける出来事対を選択する際は、下記の優先規則を順に適用して選択する。

1. 文書内で近い位置にある出来事対を優先する。
2. 「蓋然」関係となる出来事対を優先する。

4 対象テキスト

注釈付け対象テキストとして、毎日新聞 1995 年版 [15] から抽出した記事を採用した。因果関係に関する情報を正確に注釈付けるには、記事の意味内容を、実際に注釈付けを行う作業者が十分に理解する必要がある。そこで、予備調査の時点で内容の理解が比較的容易であった社会面の記事を、注釈付け対象とした。また、予備調査において、記事の長さ按比例して、注目箇所と周辺文脈との関連性が強くなり、注釈付け作業が難しくなる傾向が観察された。この結果を踏まえ、社会面の中でも記事あたり 10 文以内で構成されている比較的短い記事を 750 記事 (3912 文) 抽出し、注釈付け対象に選んだ。

5 注釈付け作業の流れ

注釈付け作業への訓練期間の後、約 1 ヶ月をかけ、上記の 750 記事に注釈付けを行った。作業者は言語学の素養をもつ 2 名に筆者の 1 人を加えた計 3 名 (以下、

A ~ C) である。注釈付けの訓練期間中は作業者間での相談を許可していたが、本作業中は相談をもち、各作業者が、互いに独立に、750 記事に注釈付けを行った。以下に作業の流れを示す。

1. 注釈付け対象テキストが記事単位で作業者に提示される。記事内の文はあらかじめ文節単位に分割されている。また、接続助詞、動詞など、因果関係の発見に寄与しそうな語句があらかじめ作業画面上でハイライト表示され、作業者に注意を促す。各作業者は、因果関係を含む箇所を探しながら作業画面上の記事を読み、3.2 節で述べた因果関係の判断基準に合格する出来事対を発見次第、タグを付与する。なお、注釈付け作業は、専用のタグ付与インタフェースを介して行われ、作業者はキーボードとマウスを使った簡単な操作でタグの付与が可能である。
2. 一定数の記事 (今回は 30 記事) に対して注釈付け作業を完了した時点で、各作業者は修正過程に移行する。修正過程では、まず、タグが付与された文字列部分が記事から自動的に抽出される (抽出結果の例として、表 3 参照)。各作業者は、抽出結果に対し、付与されたタグの範囲や種類を確認する。誤りがあれば修正を施し、修正箇所が無くなった時点で新たな記事群 (30 記事) への注釈付けに進む。

すべての作業者が 750 記事への注釈付けを終えた後、明らかに注釈付けが不適当であると思われる事例について、機械的な修正、削除処理を施す。例えば、 e_1 と e_2 が対になっておらず、片方の出来事だけが注釈付けされている場合は、該当するタグを機械的に削除した。次節以降では、この自動修正削除処理の後に得られたデータに対して、評価ならびに考察を行う。なお、一般には、メタ作業による人手の修正過程を設けることがある。しかしながら、現在の因果関係の判断基準では、メタ作業者の主観が入る余地が残されているため、メタ作業による修正過程は設けない。

6 因果関係タグ付きコーパスの概要

6.1 総数

作業ごとに認められた因果関係の総数を表 2 に示す．括弧内は記事あたりの平均因果関係数である．また，表 3 に認められた因果関係を例示する．

表 2 因果関係の総数

A	2014 (2.7)
B	1587 (2.1)
C	1048 (1.4)

今回の注釈付け作業では，作業間で同一の言語テンプレートに基づく判断基準を採用した．これにより，作業ごとの因果関係数は同程度数になると期待したが，結果は大きく異なっていた．最も多くタグを付与した作業員 A と最も少なくタグを付与した作業員 C では，約 2 倍の差がある．この差に関する考察については，8 節で後述する．

6.2 作業員間の一致度

次に，作業員間の一致度を求める．ここで，作業員間で「判断が一致している」とは，異なる作業員によってタグが付与された 2 つの因果関係事例が，次に定める一致関係を満たすことをいう．

- 2 つの因果関係事例 x と y がある．ここで， x は e_{1x} と e_{2x} から構成され， y は e_{1y} と e_{2y} から構成される．さらに， e_{1x} は $head_{1x}$ をその主辞要素とし，同様に， e_{2x} ， e_{1y} ， e_{2y} は，それぞれ， $head_{2x}$ ， $head_{1y}$ ， $head_{2y}$ をその主辞要素とする．もし， $head_{1x}$ と $head_{1y}$ が同一の文節内にあり，かつ， $head_{2x}$ と $head_{2y}$ が同一の文節内にある時， x と y は一致関係にある．

作業員（被験者）間の判断の一致度を評価する尺度として， κ 統計量が知られている．しかしながら本稿では，次の理由から，一致度を κ 統計量によって評価せず，単に事例の一致数を数えることで評価する．

κ 統計量では，複数の被験者が，同じ判断対象について判断を試行することを前提としている．一方で，今回の作業は，判断対象（因果関係の有無を判断する出来事対）があらかじめ作業員に与えられているわけではない．作業員には記事単位のデータが与えられ，各作業員は与えられた記事を読みながら，判断対象（出来事対）を走査，特定し，その後，因果関係の有無の判断を試みる．このように，判断対象は作業員が発見的に特定するため，判断対象が被験者間で異なっている可能性が十分高い．この点において，本稿のデータは， κ 統計量で一致度を測定するのに適していない．

表 4 に作業員間の一致数の集計結果を示す．表 4 左の“1”は，該当する作業員がタグを付与した，すなわち因果関係を認めたことを示し，“0”はタグを付与していないことを示す．つまり，第 1 行目から第 3 行目までは，各作業員が単独でタグを付与した事例数であり，第 4 行目から第 7 行目までは，“1”となっている 2 人以上の作業員の判断が一致してタグを付与した事例数であ

表 3 認められた因果関係（“転落”，“負う”を含む例）

	<i>mod1</i>	<i>head1</i>	<i>mod2</i>	<i>head2</i>
中学校舎から		転落する		死亡
6階から		転落する		意識不明
川に		転落		助け上げ
		転落	胸などを	打つ
軽乗用車と		殴る	けがを	負う
郵便物が		衝突	打撲傷を	負う
顔や手などに		爆発する	重傷を	負う
火傷を		負う		重傷
重傷を		負う		休職する

表 4 作業員間の一致数

A	B	C	混合	蓋然	偶然
1	0	0	921	632	535
0	1	0	487	487	255
0	0	1	187	134	207
1	1	0	372	230	90
1	0	1	133	92	77
0	1	1	140	107	83
1	1	1	588	270	64

る．例えば，第 4 行目（“110”）の“混合”の列では，作業員 A と作業員 B の判断は一致しており，A と B は共通した文節内に *head* タグを付与しているが，作業員 C はその箇所にタグを付与していない事例が 372 件あったことを表す．

表 4 の“混合”の列に，必然性の強さを区別しないで求めた一致数の集計結果を示す．作業員全員で判断が一致した事例は 588 件あり，2 人以上で判断が一致した事例は計 1233 件（= 372 + 133 + 140 + 588）あった．

次に，必然性の強さを区別して一致度を求める．そのために，先の一致関係の定義を修正し，必然性の強さを考慮した一致関係を定める．具体的には，先の定義に加え，必然性の強さの属性値も一致している場合のみ，一致関係にある，すなわち「判断が一致している」とする．例として，必然性の強さを区別しない場合に“111”に該当する事例において，A のみが「蓋然」の強さ，B と C が「偶然」の強さを属性値としてタグを付与したとしよう．この時，事例は，「蓋然」の強さに注目した場合は“100”に該当するとして集計され，「偶然」の強さに注目した場合は“011”に該当するとして集計される．

表 4 の右 2 列に，必然性の強さを区別し，「蓋然」と「偶然」のそれぞれの必然性の強さに注目した事例の一致数の集計結果を示す．表から，「蓋然」の強さをもつ因果関係事例に比べて「偶然」の強さをもつ事例において，作業員間の判断の一致度が低いことがわかる．この傾向は，作業員全員一致（“111”）の場合に特に顕著であり，「蓋然」では作業員全員の判断が一致した事例が 270 件あったが，「偶然」では 64 件しかない．「蓋然」の強さに注目して集計した事例集合と「偶然」の強さに注目して集計した事例集合の 2 群の集合について，2 人以上の作業員が一致する事例の比率を母比率の差の検定によって検定した．帰無仮説を「母比率の差が $d\%$ である」とし

て検定を実施したところ、 $d \leq 7$ では p 値 ≤ 0.00805 となり、有意水準 1% で仮説が棄却された。ここから、「蓋然」の強さをもつ事例の方が、2人以上の作業者が一致する比率が統計的に高いことが確認できる。一般に、作業間の一貫度が高い事例ほど、その事例の信頼性は高いと仮定できる。この仮定および、表 4 の結果を踏まえると、「偶然」の強さを属性値としてもつ事例集合に比べて、「蓋然」の強さを属性値としてもつ事例集合には、より信頼性の高い事例が集中していると期待できる。

7 出現特性調査

本節では、これまでに述べてきた因果関係タグ付きコーパスを元に、文書内に現れる因果関係の出現傾向を調査した結果について述べる。また、調査結果を踏まえて、因果関係知識の自動獲得の際に考慮すべき点について述べる。

調査結果の信頼性を高めるために、以降の議論では、「蓋然」の強さをもち、かつ、2人以上の作業者から因果関係があると判断された 699 件 (= 230+92+107+270) を利用する。調査項目は、(i) 因果関係の手がかり標識の有無、(ii) 出来事表現の統語カテゴリ、(iii) 出来事表現の文内位置の 3 点である。

7.1 手がかり標識の有無

文書内に現れる因果関係が、手がかりとなる標識を伴う割合について調査した。今回、注釈付け作業時に、*head*, *mod*, *causal_rel* の基本タグ以外の付加情報として、因果関係タグを付与する手がかりとなった語句に *marker* タグを合わせて付与した。ここでは、699 件の因果関係事例に対し、*marker* タグが一緒に付与されているかどうか、その割合を調べた。

表 5 手がかり標識の有無

あり	219
なし	480

結果を表 5 に示す。また、いずれかの作業者が *marker* タグを付与した手がかり標識を表 6 に列挙する。手がかり標識を一切伴わずに、因果関係をもつ 2 つの出来事が表現されることに関しては、従来からも言及がある (例えば、文献 [1])。表 5 から、このことを定量的に確認した。取り扱うデータに依存する部分もあるが、今回の場合、手がかり標識が存在する事例はおおよそ 3 割に満たなかった。

この結果を因果関係知識の自動獲得の観点から考察する。従来から開発されてきた因果関係知識の自動獲得の手法 [12, 4, 5] は、手がかり標識に基づく手法が中心である。しかしながら、今回の表 5 の結果から、高い被覆率で因果関係知識を獲得するには、手がかり標識を伴う場合に加え、手がかり標識を伴わない場合をうまく取り扱える手法の確立が必要であるといえる[‡]。

7.2 出来事表現の統語カテゴリ

次に、因果関係をもつ出来事の表現形式について調査した。具体的には、因果関係をもつ 2 つの出来事 e_1 と e_2 のそれぞれについて、それらが動詞句 (VP) と名詞

[‡] 手がかり標識に依存しない自動獲得手法として、例えば、鳥澤の手法 [13] がある。

表 6 *marker* タグが付与された手がかり標識

	頻度
ため	120
で	35
結果	5
ので	5
と	5
場合	4
ば	4
ことから	4
から	3
理由で 目的で 影響で より ように	2
ようとして ところが	
背景には 続いて 事故で 事件で	
取り入れようと 際に 際 限り	1
れると み による ており せようと	
ことで うとした ことによって	

表 7 出来事表現の統語カテゴリ

カテゴリ	該当例	e_1	e_2
VP	動詞-自立 (“焼く”)	365	412
	形容詞-自立 (“難しい”)		
NP	名詞-サ変接続 (“停電”)	322	269
	名詞-一般 (“火災”)		
	その他 (“うっとり”)		

句 (NP) のどちらの形式で表現される傾向にあるかを調べた。統語カテゴリの決定には、主辞要素の末尾形態素の品詞に注目し、動詞が形容詞であれば動詞句、名詞であれば名詞句とした。品詞情報は ChaSen [10] の解析結果を利用した。

結果を表 7 に示す。表 7 から、動詞句を形成する出来事表現が過半数を占めていることがわかる。しかしながら、名詞句を形成する事例も決して少なくはない。この傾向は、 e_1 と e_2 のどちらの出来事についても確認できる。

従来の因果関係知識の自動獲得手法では、動詞句に注目する傾向があったのに対し、今回の結果は、今後、名詞句への対応も必要であることを示唆している。名詞句は通常、“傘”や“川”といった静的なモノを表現するために使用されることが多い。そのため、名詞句を対象とした知識自動獲得を実現するには、文書中の名詞句がモノではなく、“火災”のように出来事を表現していることを自動判定する機構が必要となる。

7.3 出来事表現の出現位置

7.3.1 個々の出来事表現の出現位置

因果関係をもつ出来事が文内において表現される位置について調査した。具体的には、まず、記事内でタグが付与された元の文に対し、文末文節を根 (深さ = 0) とする係り受け木を考える。そして、出来事表現の位置をその主辞要素で代表させ、主辞要素を含む文節が位置している根からの深さを調査した。係り受け情報は CaboCha [6] の解析結果を利用した。

図 3 に e_1 の位置の分布を示す。また、図 4 に e_2 の位置の分布を示す。各図中の “f” は、ある深さでの頻度を示し、“c” はその深さまでの累積頻度を示す。図 4 の e_2 の位置分布について見ると、出来事が動詞句で形

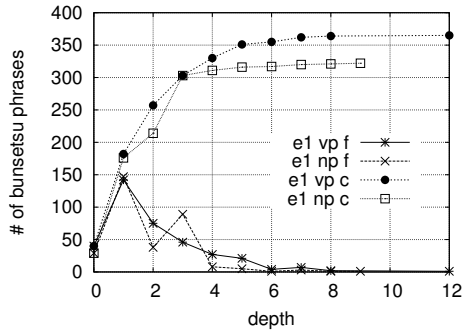


図3 出来事表現の出現位置 (e_1)

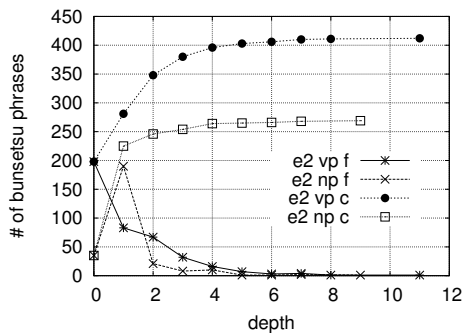


図4 出来事表現の出現位置 (e_2)

成されていれば、その出来事は文末に集中して位置し、名詞句で形成されていれば、深さ = 1、すなわち、文末に係っている文節に集中して位置することがわかる。また、動詞句と名詞句、いずれの統語カテゴリーにおいても、深さ = 2 までに 8 割以上の事例が集中していることがわかる。同様に、図 3 の e_1 についても、深さ = 3 までに 8 割以上の事例が集中していることが確認できる。

以上の結果から、因果関係をもつ出来事（の主辞）を発見的に探査する場合、比較的浅い位置のみを走査するだけで、大部分の事例が発見できることがわかる。

7.3.2 出来事表現間の相対的な出現位置

次に、因果関係をもつ e_1 と e_2 の相対的な出現位置関係を調査した。 e_1 と e_2 の位置は先程と同様、それぞれの主辞を含む文節の位置で代表させる。7.3.1 節では、個々の主辞が位置する深さに注目したが、ここでは個々の主辞の深さの差に注目する。

結果を表 8 に示す。表中の“ $e_1 \Rightarrow e_2$ ”は、 e_1 が e_2 よりも係り受け木の深い位置にある事を示し、“ $e_2 \Rightarrow e_1$ ”は、 e_1 が e_2 よりも浅い位置にある事を示す。“no dep”は、 e_1 から係り受け関係を辿った際、先祖あるいは子孫の位置に e_2 が存在しない事を示す。

まず、 e_1 と e_2 の出現順序について調査する。表 8 の“ $e_1 \Rightarrow e_2$ ”と“ $e_2 \Rightarrow e_1$ ”に該当する事例数の比較から、 e_1 は e_2 よりも深い位置に位置する傾向が強いことがわかる。つまり、原因を述べた後に続いて結果が述べられることが多い。しかしながら、両者の出来事の位置が逆転する事例も存在する。事例分析から、 e_1 が、 e_2 の目的に対する手段を表現する場合に、“ $e_2 \Rightarrow e_1$ ”の順序になる傾向が観察された。“ $e_2 \Rightarrow e_1$ ”に該当する事例を含む元文を表 9 に例示する。ただし表 9 において、ここでは出現順序に焦点を当てているため、議論に直接関係

表 8 出来事表現間の相対的な出現位置

		$e_1 \Rightarrow e_2$	$e_2 \Rightarrow e_1$
文内	深さの差 = 1	259	15
	= 2	152	23
	> 2	33	4
no dep		72	
文間			141

表 9 $e_2 \Rightarrow e_1$ となる例

市への葬儀費用の返還を求める住民監査請求を行う
相互理解を深める目的で世界民族芸能祭を開催
冥福を祈って、黙とうした
目立つよう真っ赤に塗る
期待にお応えすべく努力する

のない修飾要素などは実際の元文から除去している。

続いて、出来事が位置する深さの差について調査する。直感的には、 e_1 の主辞要素が e_2 の主辞要素に直接係っている場合が多いと予想され、実際の結果も 259 件と最も多かった。しかしながら、それ以外の位置関係にある場合も少なからず存在していることが明らかになった。深さの差 = 2 となる事例には、「濡れたシートで足を滑らせる」の「濡れる」と「滑る」のように、連体修飾関係を挟む場合が多く含まれていた。“no dep”に該当する事例には、抽出された元文中に並列関係が存在している場合が多く、現在の二項関係という枠組みでは捉え難い事例が多く集まっていた。

8 考察：判断の網羅性の検証

先に表 2 で見たように、今回の注釈付け作業では、付与したタグの総数が作業員間で大きく異なる。本稿の最後にここでは、この作業員間のタグ総数の相違について、各作業員の判断の網羅性という観点から考察する。

本稿で述べた注釈付け作業では、作業員に記事が提示され、各作業員は、そこから発見的に因果関係が成立する出来事対を特定する。この結果、因果関係タグが付与されていない出来事対には、次の 2 通りの可能性が残ることになる。

- 該当箇所の出来事対には因果関係がないと判断され、タグが付与されなかった。
- 該当箇所の出来事対に関して因果関係の有無を判断しておらず、タグが付与されなかった。

1 点目の可能性は、各作業員の因果関係の認識そのものに起因する。一方、2 点目の可能性は、作業員の判断が網羅性に欠けていることから生じており、注釈付け作業の設計やインターフェースに関わる問題である。特に、今回の注釈付け作業では、手がかり標識を伴わない箇所への注釈付けを要求しており、2 点目に該当する箇所が先行研究と比較しても多いと考えられる。

以上を踏まえ、上記の 2 通りの可能性のうち、2 点目に該当する事例の割合を検証するために、以下の手続きから構成される追作業を実施した。

1. e_1 と e_2 の候補となる出来事対を作業員に提示する。
2. 各作業員は、提示された出来事対に関して因果関係の有無を判断する。

表 10 本作業と追作業

		本作業	
		1	0
追作業	1	<i>a</i>	<i>b</i>
	0	<i>c</i>	<i>d</i>

表 11 作業 A の結果

			A	1	0
1			196	327	
0			10	67	

表 12 作業 B の結果

			B	1	0
1			217	345	
0			0	38	

表 13 作業 C の結果

			C	1	0
1			233	278	
0			2	87	

提示する出来事対は本作業時に A から C のいずれかの作業者がタグを付与した事例である。ある作業 X に提示される事例は、作業 X が本作業時に注釈つけた因果関係事例と、X は注釈づけなかったが、X 以外の作業者が注釈つけた因果関係事例をおおよそ 1 : 2 の割合で混ぜたもので、それらが無作為な順序で提示される。

提示された出来事対は、本作業と追作業の結果によって、表 10 のいずれかに分類される。表中の“1”はタグを付与した（因果関係があると判断した）ことを示し、“0”はタグを付与していない（因果関係がないと判断した）ことを示す。もし、追作業で提示されたすべての出来事対に対して、作業者が本作業においても見落とすことなく因果関係の有無に関する判断を行っており、かつ、その判断の再現性が保証されていれば、追作業で提示された事例は、すべて“*a*”か“*d*”のカテゴリに分類されることになる。一方、追作業では、判断対象が直接作業者に提示されるため、「該当部分の出来事対に関して因果関係の有無を判断しておらず、タグが付与されなかった」という可能性は存在しない。また、提示される事例は、本作業時に A から C のいずれかの作業者がタグを付与した事例である。以上から、本作業時に判断がなされず、見落とされていた事例の割合に比例して、集計結果には“*b*”のカテゴリに分類される事例が増加すると予想できる。

上記の追作業を 600 件の出来事対を用いて実施した。各作業者の結果を表 11 から表 13 に示す。各表から、いずれの作業においてもカテゴリ“*b*”の事例数が多いことがわかる。このことは、本作業時にタグが付与されなかった事例の多くについて、実際には判断がなされなかったことを示唆しており、また、判断の試行対象が作業者によって大きく異なっていた可能性を示している。定性的な事例分析から、手がかり標識が伴わない場合、特に、出来事対が連体修飾節を介して表現される場合や、2 つの出来事が文をまたいで表現される場合では、判断がなされない傾向が強いことを確認した。

9 おわりに

本稿では、まず、一定量のテキスト集合に対して因果関係情報を注釈付けたその過程について述べた。また、注釈付けの結果得られたコーパスを用いて、テキスト中に現れる因果関係の出現傾向を調査し、その結果を報告した。調査から、手がかり標識と伴う割合など、因果関係知識の自動獲得手法を開発するにあたり、考慮すべき点に関する有益な知見を得た。

今後の課題として、コーパスの増量と共に、3.3 節で述べた制約を解除することがある。また、今回は出来事の主辞要素（*head* タグ）の情報を元に議論を進めたが、修飾要素（*mod* タグ）の情報からも検討を進めたい。

謝辞

本研究は、21 世紀 COE プログラム「大規模知識資源の体系化と活用基盤構築」、ならびに科学研究費補助金学術創成研究「言語理解と行動制御」（課題番号：13NP0301）の支援を受けて行われた。本研究で用いたコーパスの作成にあたり、ランゲージュアの衛藤純司氏、日本システムアプリケーションの植田禎子氏、十河則子氏、奈良先端科学技術大学院大学（現在、富士通研究所）の高橋哲朗氏の諸氏から多大な協力を頂きました。諸氏の皆様に感謝いたします。

参考文献

- [1] 有田節子. 因果の言語学. 月刊言語, Vol. 25, No. 5, pp. 20–23, 1996.
- [2] B. Altenberg. Causal linking in spoken and written English. *Studia Linguistica*, Vol. 38, p. 1, 1984.
- [3] R. Girju. Automatic detection of causal relations for question answering. In *Proc. of ACL 2003, Workshop on Multilingual Summarization and Question Answering - Machine Learning and Beyond*, 2003.
- [4] T. Inui. *Acquiring Causal Knowledge from Text Using Connective Markers*. PhD thesis, Graduate School of Information Science, Nara Institute of Science and Technology, 2004.
- [5] 乾孝司, 乾健太郎, 松本裕治. 接続標識「ため」に基づく文書集合からの因果関係知識の自動獲得. 情報処理学会論文誌, Vol. 45, No. 3, pp. 919–933, 2004.
- [6] 工藤拓, 松本裕治. チャンキングの段階適用による係り受け解析. 情報処理学会論文誌, Vol. 43, No. 6, pp. 1834–1842, 2002.
- [7] Y. Liu. 現代日本語における因果関係を表す接続語の意味・用法 - 「だから」「したがって」「それゆえ(に)」「その結果」「そのため(に)」を中心に -. PhD thesis, 名古屋大学大学院国際言語文化研究科, 2004.
- [8] W. C. Mann and S. A. Thompson. Rhetorical structure theory: A theory of text organization. In *USC Information Sciences Institute, Technical Report ISI/RS-87-190*, 1987.
- [9] D. Marcu. *The Rhetorical Parsing, Summarization, and Generation of Natural Language Texts*. PhD thesis, Department of Computer Science, University of Toronto, 1997.
- [10] 松本裕治, 北内啓, 山下達雄, 平野善隆, 松田寛, 高岡一馬, 浅原正幸. 日本語形態素解析システム『茶筌』version 2.2.1 使用説明書. 奈良先端科学技術大学院大学, 2000.
- [11] 西澤信一郎, 中川裕志. 日本語の自由会話における談話構造の推定 - 因果関係を表す接続詞の場合 -. 自然言語処理 (技術資料), Vol. 4, No. 4, pp. 61–72, 1997.
- [12] 佐藤浩史, 笠原要, 松澤和光. テキスト上の表層的因果知識の獲得とその応用. 信学技報 (TL98-23), 1999.
- [13] 鳥澤健太郎. 「常識的」推論規則のコーパスからの自動抽出. 言語処理学会第 9 回年時大会, pp. 318–321, 2003.
- [14] T. Yokoi. The edr electronic dictionary. *Communications of the ACM*, Vol. 38, No. 11, pp. 42–44, 1995.
- [15] 毎日新聞社. 毎日新聞 CD-ROM 版 (1995).