論 修士 文

テキストデータにおける予定変更情報 および影響の獲得

Obtaining Information of Schedule Changes and Their Impact in Text Data

指導教員

喜連川 優 教授



東京大学 情報理工学系研究科 電子情報学専攻

氏 名 48-126450 栗原 俊明

提出日

平成26年2月6日

我々人間が適切な意思決定を行う上で、未来の予定や計画について知っておくことは必要不可欠である。例えば、企業は未来の市場変化や競合他社の計画を知ることはビジネスの成功に大きく役に立つ。政府機関であれば国際情勢の変化、一般的の家族であれば週末のイベントなどを知りたいとも考えるだろう。

しかし、未来の予定や計画は変更されることが多々ある。予定の変更があった場合、我々はこれまで想定していた未来を更新する必要がある。変更になったことに気づかずに、誤った未来を想定していれば、不適切な意思決定をしてしまうことにつながる。さらに、予定間には複雑な関連性があるため、一つの予定の変更はその他の予定の実施にも影響を及ぼすことが想定される。

そこで本研究では、テキストデータから予定変更情報を獲得し、さらにその予定変更情報が与える影響を獲得する手法を提案する。まずテキストデータ内の時間表現を正規化し、未来関連情報を抽出する。その際、新聞記事データとブログデータにおける時間表現の違いや、計画的なイベントや突発的なイベントに際して時間表現の出現パターンがどのように変化するか等の考察も行う。その後、予定変更情報の獲得手法を提案する。予定変更情報に特有の手がかり表現を用いたパターンマッチにより抽出を行った後、機械学習によってフィルタリングを行う。さらに、テキストコーパスにおける前提関係を表す表現を用いた教師なし学習により、予定変更情報により影響を受ける未来関連情報を獲得する。毎日新聞コーパスを用いた実験により、本研究における提案手法が有効であることを示す。

謝辞

はじめに、指導教官である喜連川優教授に深く感謝致します。喜連川教授には、まず、基礎を大切にすることや研究を楽しむこと等、研究に対する姿勢のあり方を教えていただきました。また、研究の進め方に関する貴重な助言を賜りました。喜連川教授に師事する中で、貴重な経験を積み、多くのことを学ぶことがきました。

次に、豊田正史准教授には、研究の方向性について多くの指導をして頂いたことに感謝致します。 夜遅くにも関わらず相談にのってくださったり、締切の直前には休日にもご相談にのっていただきました。また、論文の書き方から発表の仕方まで、本当に多くのことについても教授して頂きました。

また、鍜治伸裕特任准教授、吉永直樹特任准教授を始めとする研究室のスタッフの皆様には、ミーティングなどで研究に関する丁寧かつ熱心なご指導をして頂きましたことに感謝致します。週二回のミーティングや発表練習等で多大なる時間を割いてご指導して頂き、本当に多くのことを学ぶことができました。特に、吉永直樹特任准教授には、機械学習や自然言語処理等の分野における様々な手法を個人指導していただく等お世話になりました。併せて、研究を快適に行えるよう研究室環境を支えてくださった秘書の方々に感謝いたします。

2014年2月6日

目 次

謝辞		i
第1章	はじめに	1
1.1	テキストデータにおける未来関連情報	1
1.2	本研究の目的と貢献	2
1.3	本論文の構成	2
第2章	関連研究	4
第3章	時間表現の正規化および未来関連情報の抽出	6
3.1	時間表現の正規化	6
3.2	毎日新聞とブログ間における時間表現の比較	9
3.3	文書作成時間と参照時間の関係性分析	10
	3.3.1 計画的なイベント	11
	3.3.2 突発的なイベント	13
	3.3.3 予定変更	13
第4章	予定変更情報の獲得	16
4.1	手法の概要	16
4.2	パターンマッチによる抽出	17
4.3	機械学習を用いたフィルタリング	18
4.4	実験	18
	4.4.1 パターンマッチによる抽出	18
	4.4.2 機械学習を用いたフィルタリング	19

4.5	獲得した予	定変更情報の分析	î			 	 	21
第5章	予定変更に	こよる影響の獲得						24
5.1	手法の目的]および取り組む〜	だき課題			 	 	24
5.2	提案手法 .					 	 	25
	5.2.1 前捷	是表現を含む文を月	別いた分類器	器の構築		 	 	25
	5.2.2 分类	頁器による入力の <i>></i>	スコア付け			 	 	27
5.3	実験					 	 	28
	5.3.1 提到	ととという とうとう とうしょ とうしょ とうしょ とうしょ とうしょ とうしゅ とうしゅ とうしゅ とうしゅ とうしゅ とうしゅ とうしゅ とうしゅ				 	 	28
	5.3.2 変見	E情報からの影響 獲	矆得例			 	 	29
	5.3.3 影響	響獲得数による予定 	定変更情報の)ランキ	ング.	 	 	30
第6章	おわりに							34
参考文南	,							36
発表文南	ţ.							40

図目次

3.1	文長の分布	7
3.2	各時間表現の出現数 (毎日新聞)	9
3.3	各時間表現の出現数 (ブログ)	10
3.4	文書作成時間と参照時間の時間差(年粒度)	11
3.5	文書作成時間と参照時間の時間差(月粒度)	11
3.6	文書作成時間と参照時間の関係(地デジ化)	12
3.7	文書作成時間と参照時間の関係(消費税)	13
3.8	文書作成時間と参照時間の関係(原発・原子力)	14
3.9	文書作成時間と参照時間の関係(スペースシャトル・ディスカバリー)	14
4.1	素性とする形態素数による分類性能(学習データにおける交差検定)	20
4.2	設定する分類確率閾値による各評価値の変化	21
4.3	ドキュメント作成月別の獲得文数	22
4.4	月別の獲得文数(17年平均)	22
5.1	係り受け解析を用いた学習データの作成	26

第1章 はじめに

1.1 テキストデータにおける未来関連情報

我々人間が適切な意思決定を行う上で、未来を予測することは必要不可欠である. 例えば、企業は未来の市場変化や競合他社の計画を知ることはビジネスの成功に大きく役に立つ. 政府機関であれば国際情勢の変化、一般市民であれば週末のイベントなどを知りたいとも考えるだろう.

近年多量に蓄積され続けているブログ、オンラインニュース記事、ツイッター等のテキストデータは未来を予測する上での貴重な資源である。例えば、本研究で使用しているブログデータ(約23億文)のうち1600万文以上が未来の時間点を示す表現("来年"、"2015年1月"etc.)を含んでいる。また、CIAやGoogleが出資していることでも話題となったRecorded Future 1 という企業は、Web上の未来関連情報の活用に特化したサービスを提供している。

しかし、予定や計画は変更されることが多々ある。予定の変更があった場合、我々がこれまで想定していた未来を塗り替える必要がある。変更になったことに気づかずに、誤った未来を想定していれば、不適切な意思決定をしてしまうことにつながる。そのため、我々はこれらの予定変更に関する情報を常に把握していることが望ましい。

さらに、予定間には複雑な関連性があるため、一つの予定の変更はその他の予定の実施にも影響を及ぼすことが想定される。例えば、現在2020年の東京オリンピックに向けて、競技場の建設、交通等のインフラの整備、様々な競技の試合計画等の多数の予定が行われようとしている。そのため、もし仮に2020年の東京オリンピック開催に関して重大な変更が起きた場合には、これら多数のイベントに影響が及ぶ

¹https://www.recordedfuture.com/

ことは間違いないであろう.

また,実際に予定が変更されたものを例に挙げるとするばらば,東日本大震災による原発事故で"2030年に向けたエネルギー基本計画"が白紙に戻されたが,これは日本のエネルギーに関する様々な予定に影響を与えた.例えば,原発の建設計画,再生可能エネルギーの普及,新技術開発,海外への原発輸出政策,地球温暖化対策としてのCO2削減目標等,多くの予定が影響を受けたはずである.つまり,一つの予定が変更された際に,他に影響を受ける予定をも把握することができれば,正しい未来を思い描く上での助けとなる.

1.2 本研究の目的と貢献

本研究では、テキストデータから予定変更情報を獲得し、さらにその予定変更情報が与える影響を獲得する手法を提案する.具体的な手順としては、まず、テキストデータ内の時間表現を正規化し、未来関連情報を抽出した.その後、予定変更情報に特有のパターンを用いたパターンマッチにより抽出を行った後、機械学習によってフィルタリングを行うことで、予定変更情報を獲得した.さらに、テキストコーパスにおける前提関係を表す表現を用いた教師なし学習により、予定変更情報により影響を受ける未来関連情報を獲得した.

実験では、毎日新聞コーパスを用いて提案手法の有効性を示す. ただし、未来関連情報の抽出に関しては、ブログデータも用いて実験を行い、両データの比較を行った.

1.3 本論文の構成

本論文の構成は次のとおりである。

第2章 これまでに行われた関連研究について順に俯瞰をした後,本論文の位置付けについて説明する。

- **第3章** テキストデータにおける時間表現を正規化し、未来に関連する情報を抽出する手法について説明する. 毎日新聞とブログデータを用いて実験を行い、両データにおける時間表現の違いについても議論する.
- 第4章 未来関連情報から予定の変更に関する情報を獲得する手法について説明する. 毎日新聞データを用いた実験により、提案手法の有効性を示す.
- **第5章** 予定変更情報により影響を受ける未来関連情報を獲得する手法について説明を行う.
- 第6章 全体のまとめと今後の課題について述べる。

第2章 関連研究

これまで、テキストデータ内の時間表現の活用を目指した研究は数多く行われている.また、その中でも未来の時間表現に着目し、未来予測に役立てようとした研究もいくつか行われてきた.以下で、これらの研究について順に俯瞰する.

テキストデータ内の時間表現を活用する試みとしては、Temporal Information Retrieval なる分野が存在している [22] [13] [14] [6] [4] [7] [10] [5]. これまでテキストデータの時間情報としては文書作成時間(DCT: Document Creation Time)のみの活用にとどまっていたが、本研究分野ではテキスト内の時間表現を、クエリに対してタイムライン形式で話題を俯瞰することを支援したり、類似度の高いテキストをみつけること等に活用することを目指している [26] [24] [23]. また、時間表現を活用しているわけではないが、テキストデータにおける話題の移り変わりを把握することを目的とした分野も存在している [8] [28] [9].

また、SemEval のワークショップである TempEval [18] [19] [16] では、文構造などに着目して、テキスト内のの文書作成時間、イベント、時間表現の関係性を特定する手法の研究等が行われている.

テキストデータにおける未来関連情報に関する研究を初めて行ったのは Baeza-Yates による研究 [25] であるとされており、テキストと時間情報からなるクエリに対して、ニュースアーカイブから未来関連情報を検索する future retrieval というタスクを定義した。その後 Jatowt らは、人名、地名、組織名、イベント名などのクエリに対して、Google News アーカイブから未来関連情報を検索し、それらを分類し可視化する手法の提案 [3] や、関連するイベントが発生する年度の確率分布を求める研究 [1] 等を行っている [2]. Dias ら [11] は参照時間を用いて未来関連情報をトピックごとに分類する研究等を行った。また、Kawai らはクエリとテキスト内の未

来時間の関連性の判定等の研究 [12] を行っている.

その他にも、未来関連情報の応用的な活用として、Ho ら [27] は未来関連情報の うち、地名を含むものに着目し、位置情報を利用した推薦システムに応用している. そして、Kanhabuaら [20] は、オンラインニュース記事のユーザーに対して、記事 と関連性の高い未来関連情報を提供するためのランキング手法を提案した.

上記のように、テキストデータにおける未来関連情報に関する研究はいくつか行われている。しかし、これらの研究の中で予定変更情報に焦点を置いた研究は行われていない。

第3章 時間表現の正規化および未来関 連情報の抽出

3.1 時間表現の正規化

本節では、テキストデータにおける時間表現を正規化し、未来関連情報を抽出する手法について説明する.本研究では、未来の時間点を示す時間表現を含む文を未 来関連情報として定義する。未来関連情報の一例を以下に挙げる。

• 未来関連情報の例

- "安倍晋三首相は1日夕,官邸で記者会見を行い,平成26年4月に消費税率を現行の5%から8%に引き上げることを正式に表明した."
- "国際オリンピック委員会 (IOC) は7日 (日本時間8日), 2020年 の第32回夏季オリンピック競技大会(五輪)を東京で開催することをブエノスアイレス総会で決めた."

まず、本研究で用いる実験データについて説明する。本研究では、毎日新聞データとブログデータを用いている。それぞれの期間および文数を表 3.1 に示す。また、図 3.1 に各データの文長の分布を示す。

表 3.1: 実験データ

	期間	文数
毎日新聞	1996年1月~2012年12月	15,732,878
ブログ	2005年1月~2013年6月	2,340,932,402

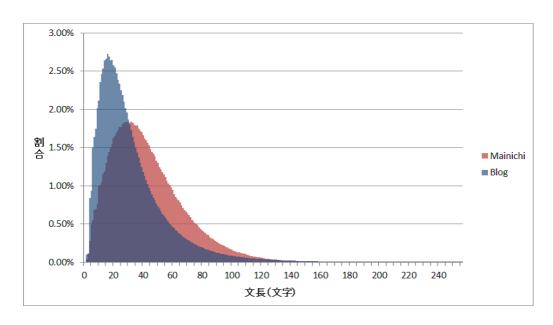


図 3.1: 文長の分布

Pustejovsky ら [15] によると、文章内の時間表現の種類には大きく分けて4種類ある.「期間」、「セット」、「日付」と「時間」である. 1 つ目の「期間」は、例えば "they have been traveling through the U.S. for three years "の" for three years "のように、時間の長さを表現する. 2つ目の「セット」は、" She goes to the gym twice a week."の"twice a week"のように時間の周期性について表現する. そして「日付」と「時間」は、" 3 p.m."や" January 25, 2010"のように、具体的な時間点を指す.

また、Omar ら [22] によると、「日付」と「時間」については、表現方法によって、「明示的表現」「相対的表現」「黙示的表現」に分けることができる。「明示的表現」は、"January 25, 2010"のように、時間軸上の1つの点を指す。「相対的表現」は、"today"や"next month"のように、その文章の文脈を考慮することで、指す日付を特定できるものである。そして"黙示的表現"は、"New Year's Day 2002"等のように、黙示的に日付を指し示しているものである。

本研究では、特定の日付を示す時間表現のうち主な明示的表現と相対的表現を正規化した。表 3.2 に正規化した表現を示す。ただし、Nは任意の数字列、"NNNN

衣 3.2: 止稅化衣块					
	絶対的表現	相対的表現			
月単位	NNNN 年(の)N 月	来年(の)N 月			
	平成N年(の)N月	今年(の)N月			
	昭和N年(の)N月	昨年(の)N月			
	NN 年(の)N 月	去年(の)N月			
		N月, 来月, 今月, 先月			
年単位	NNNN 年	来年, 今年, 昨年, 去年			
	平成 N 年, 昭和 N 年				

表 3.2: 正規化表現

年"は西暦4桁の年度表現, "NN年N月"は"98年1月"等の西暦年を省略して2桁で表しているものを指す.

相対的表現については、文書作成時間を用いて正規化を行った。また、本研究では"N月"(1月、2月 etc.)の時間表現の正規化については、文書作成時間と近い時点へ正規化した。例えば、1月に作成された文書において"12月"という表現がある場合、候補としては"前年の12月"と"同年の12月"が特に有力であると考えられるが、文書作成時間により近い"前年の12月"に正規化する等とした。

毎日新聞およびブログにおいて時間表現を正規化した結果を、それぞれ図3.2、図3.3に示す. "昨年"と"去年"については、同じ時間点を示すパターンであるため、まとめて表示した. また、青色は文書作成時間よりも"過去"、赤色は文書作成時間における"現在"、緑色は文書作成時間よりも"未来"の時間点を示しているものである. 月粒度の表現では、文書作成時間と同じ月を示す表現の場合は"現在"、年粒度の場合は、文書作成時間と同じ年を示す表現の場合に"現在"というようにしている.

この結果から、両データにおいて"N月"という表現が多いことがわかる。また、 "NN年N月"や"来年のN月"等の表現は、ブログでは毎日新聞に比べて少ない ことがわかる。

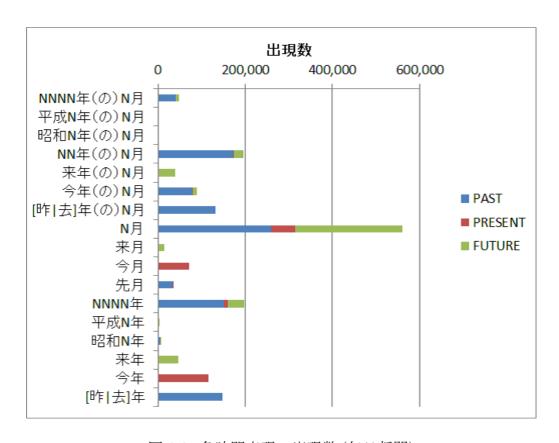


図 3.2: 各時間表現の出現数 (毎日新聞)

3.2 毎日新聞とブログ間における時間表現の比較

図 3.4 と図 3.5 は正規化した時間表現と文書作成時間との時間差について、その出現割合をプロットしたものである. 図 3.4 は年粒度のもの、図 3.5 は月粒度のものをプロットしたものである.

これらの図から,毎日新聞データとブログデータにおいて出現する時間表現の違いがわかる.年粒度,月粒度ともブログデータは毎日新聞データよりも時間差0の割合が高い.これは,ブログデータは毎日新聞データに比べて,文書が作成された日時と同じ年月の情報が多いことを示唆している.一方,毎日新聞は過去から未来までブログデータよりも幅広い情報を提供していると考えられる.



図 3.3: 各時間表現の出現数 (ブログ)

3.3 文書作成時間と参照時間の関係性分析

テキストデータは、前節で正規化を行ったテキスト内の時間情報(参照時間: Reference Time)と、テキストが作成された時間に関する情報(文書作成時間: Document Creation Time)の主に二つの時間情報を持つ。これらの二つの時間情報の関係性を視覚的に捉えることは、未来関連情報を俯瞰する上で役に立つだけでなく、テキストデータの時間情報がもつ性質の理解に役に立つ。本節では未来時間点のものに限らず、いくつかのトピックを例に、テキストデータにおける2つの時間情報の関係性を可視化した結果について紹介する。

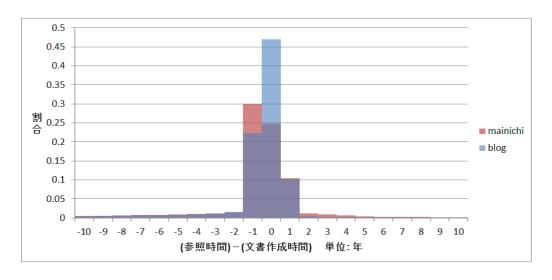


図 3.4: 文書作成時間と参照時間の時間差(年粒度)

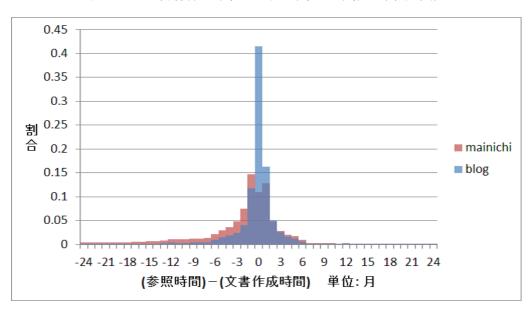


図 3.5: 文書作成時間と参照時間の時間差(月粒度)

3.3.1 計画的なイベント

図3.6は、"地上デジタル放送への移行"に関する情報の文書作成時間と参照時間の関係性をプロットしたものである。前節で正規化したもののうち、月単位の粒度のものをプロットしている。また、各棒グラフの色は同一の参照時間のものを見や

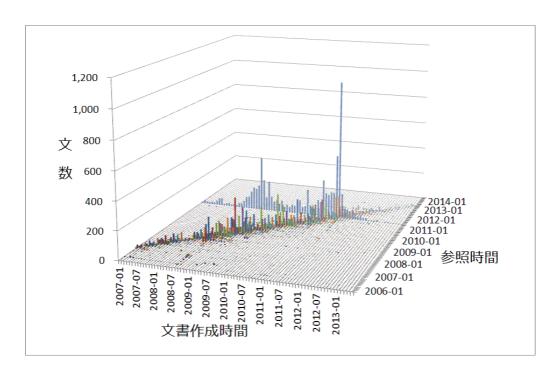


図 3.6: 文書作成時間と参照時間の関係(地デジ化)

すくするように、同一の参照時間のものを同じ色にしている.

文書作成時間と参照時間が同じまたはその前後の情報(すなわち文書が作成された月と同じ月に関する情報)は基本的によく参照されている. さらに、地上デジタル放送への完全移行月である2011年7月への参照も多く、文書作成時間の経過により、その参照数も上下している. この地上デジタル放送への移行のように、事前に計画的なイベントに関するテキスト内での時間情報は、図3.6のようになると考えられる.

次に"消費税"に関する二つの時間情報の関係性を図 3.7 に示す。2011 年末より消費税引き上げに関する議論が活発化し、「2014 年 4 月に税率 8 %、2015 年 10 月に 10 %」という法案について議論が続いた。そして 2012 年 8 月に法案が成立すると、2015 年 10 月に比べて 2014 年 4 月により焦点が当たり、引き上げの可否や準備に関する記述が増えている。

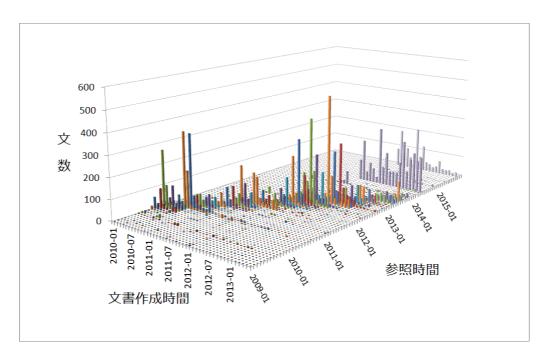


図 3.7: 文書作成時間と参照時間の関係(消費税)

3.3.2 突発的なイベント

一方,東日本大震災のような突発的に大きなイベントが発生した場合の例として, "原子力・原発"に関するものを図 3.8 に示す. 震災が発生した 2011 年 3 月は,それ以前の文書作成時間のテキストではほとんど参照されていない. しかし, 震災発生後, 2011 年 3 月という時間点への参照は,時間が経過するにつれて減少するものの他と比べて多い.

このように、テキストデータにおける時間表現の出現分布を用いることで、未来 の予定の有無やその話題の度合い、または震災や金融危機等の突発的イベントの発 生を把握できることが考えられる.

3.3.3 予定変更

また、本研究で焦点を当てている予定変更情報が、テキストデータ内の時間表現 の出現パターンに与える影響について分析するため、"スペースシャトル・ディスカ

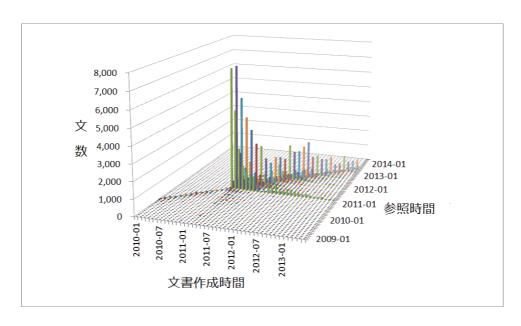


図 3.8: 文書作成時間と参照時間の関係(原発・原子力)

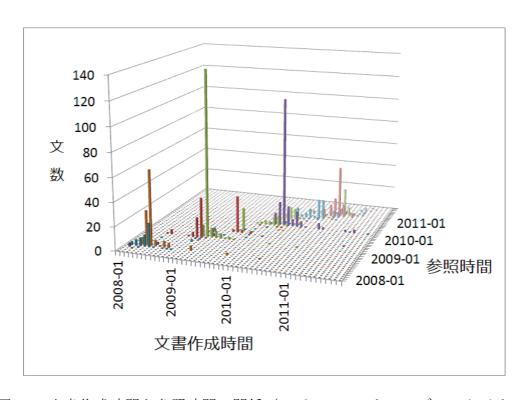


図 3.9: 文書作成時間と参照時間の関係 (スペースシャトル・ディスカバリー)

バリー"に関するものを図 3.9 に示す.この図を見ると,地上デジタル放送や消費税のものとはやや異なるプロットになっていることがわかる.例えば,一番左の塊を見てみると,水色がずっと参照されていたにも関わらず,途中からオレンジ色の参照時間にずれ込み,最終的にはオレンジ色の部分が大きくバーストしている.これは 2008 年 2 月に,2008 年 4 月に予定されていたスペースシャトルディスカバリーの打ち上げが,2008 年 5 月に延期されたことによるものである.またその他の部分ついても,同様のことが起こっていることがわかる.

第4章 予定変更情報の獲得

4.1 手法の概要

予定や計画は変更されることが多々ある.特に,災害や事故等の不測の事態が発生した際には,多くの予定が変更を余儀なくされる.予定の変更があった場合,我々はこれまでの予定情報を更新する必要がある.変更になったことに気づかずに,誤った未来を想定していれば,不適切な意思決定をしてしまうことにつながる.そのため,我々はこれらの予定変更に関する情報を常に把握していることが望ましい.

そこで、未来関連情報から予定変更情報を獲得する.本研究における予定変更情報は、以下の性質をもつ文として定義する.

予定変更情報の性質

- 未来の特定の年月日に設定されていた予定や目標を,別の日程に変更する,または,中止することを表現する.
- 情報として確定されたものである.(可能性を示唆するのみのものは含まない)

予定変更情報の一例を以下に挙げる。

● 予定変更情報の例

- "鳩山首相は米軍普天間飛行場の移設問題の結論を来年に先送りし、アフガン支援策の策定を優先させる意向を表明した。"
- "東日本大震災で大きな被害を受けた岩手,宮城,福島の3県について,7月24日に予定されていた地上デジタル放送の完全移行を最大1年間延期すると総務省が発表した."

表 4.1: パターンマッチに用いた表現

予定変更を表す名詞・動詞

延期,延長,変更,中止,前倒し,断念,保留,見送る 先送り,諦める,間に合わない,ずれ込む,ずれこむ 見合わせる,見合せる,持ち越す,持ちこす,もちこす

予定を過去化する表現

当初,もともと,予定だった,計画だった,つもりだった はずだった,が予定されていた,と発表していた に控えていた,表明していた,だったが

その他

予定より、計画より、計画〜白紙、計画〜練り直す、を待たずに 日程〜誤り、まで続投する、まで〜続けることを決めた できないことがわかった、以降も当分、以降になる見通し

本手順の流れは以下である。まず、未来関連情報全体からランダムに抽出したデータから、予定変更情報において特徴的なパターンを収集した。その後それらのパターンにマッチした文をコーパス全体から抽出し、機械学習を用いてフィルタリングを行い、獲得精度を向上させる。

4.2 パターンマッチによる抽出

毎日新聞コーパスのうち、未来の時間表現を含むもののうち 2000 文を人手でチェックし、予定変更情報に特有のパターンを収集した. 2000 文のうち 56 文が予定変更情報であった.

これらの予定変更情報内で発見されたパターンおよびそこから推測されるパターンを 40 個選択した. これらのパターンを選択する上で,日本語 WordNet 1 等も利用した. 表 4.1 にその 40 パターンを示す.主なパターンとしては"延期"や"中止"等の予定の変更を表す名詞や動詞や,"~予定だった"や"当初は~"のように予定を過去化する表現などがあった.

¹http://nlpwww.nict.go.jp/wn-ja/

4.3 機械学習を用いたフィルタリング

パターンマッチで抽出された文集合に対して機械学習を用いてフィルタリングを 行う. 本手順で除外すべき文の例を以下に挙げる.

• 不確定

- 活動期間は来年3月までの予定だが、延長の可能性もある.
- 併せて、来月31日の控訴趣意書提出期限の延期を求めた.
- 米国の利上げが11月の大統領選挙後にずれ込むとの見方が、市場関係者 の間で支配的になってきた.

● 否定

- 来年5月からの実施予定に変更がないとの立場を表明した.
- 今年12月としている工事の完了予定は変更しない.

本研究では、上記の表現を除外するには、前節でマッチしたパターンの直後の表現が大きな手がかりになることに着目した。例えばマッチパターンが"延期"の場合に、直後に"可能性"や"~しない"などの表現があれば、除外すべき情報である可能性は高い。そこで、文を $MeCab^2$ で形態素解析及び見出し語化し、マッチパターンおよびその直後のN個の形態素をその文の特徴ベクトルとした。パラメーターNについては、実験において学習データの交差検定を行い、そこで最も高いF値を示したものとする。

4.4 実験

4.4.1 パターンマッチによる抽出

コーパス全体からの表 4.1 の 40 パターンを含む文を抽出した. その結果,全体で 387,835 文あった未来関連情報のうち,17,946 文が抽出された.

²http://mecab.googlecode.com/svn/trunk/mecab/doc/index.html

表 4.2: 学習データとテストデータ

<u> </u>		/ / •	/ /
	全体	正例	負例
学習データ	200	91	109
テストデータ	200	108	92

このパターンマッチによる抽出で、どれだけの予定変更情報がとれているかを概算で評価する。全コーパスは 387,835 文で、2,000 文中 56 文が予定変更情報であったため、概算で 10、859 文存在することになる。一方、パターンマッチによる抽出後は、全体で 17,946 文あり、400 文中 199 文が正例であった。よって概算で 8,928 存在することになる。つまり、この手順による予定変更情報の損失は概算で 1,931 文(17.8%)である。

4.4.2 機械学習を用いたフィルタリング

まず、前節で抽出した文集合から後の機械学習で用いる学習データとテストデータを各々200 文ずつランダムに選択し、ラベル付した.表 4.2 にこの結果を示す.正例は予定変更情報、負例は予定変更情報でないものである.

前節で抽出した文集合に機械学習でフィルタリングを行い、精度向上を目指す.また、本研究の機械学習では機械学習を用いた手法では、パッシブアグレッシブアルゴリズムであるオンライン学習器 opal³ [21] を用いた. 学習割合や繰り返し回数等のパラメータはデフォルト値を用い、カーネルは線形カーネルとした. オプションとしては、学習データのシャッフリング (-s) とパラメータの平均化 (-a) を設定した.マッチパターン後のいくつの形態素を素性とするかについては、学習データの交差検定で最も良い評価値を示したものを採用した. 図 4.1 にその結果を示す. "Match & behind_N"は、マッチパターンとその直後 N 形態素を用いたケースを示している.交差検定としては、ランダムな 2 分割交差検定を 5 回行い、その平均値を評価値とした. その結果、Match & behind_8 が最も高い F 値を示したので、これを採用する.

³http://www.tkl.iis.u-tokyo.ac.jp/ ynaga/opal/

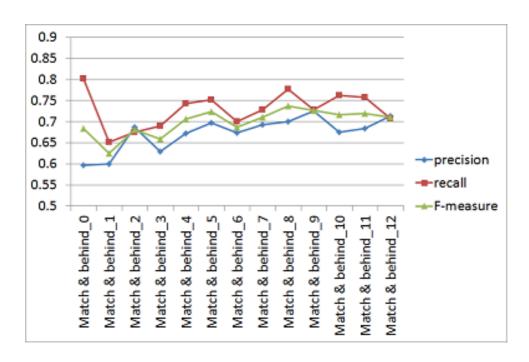


図 4.1: 素性とする形態素数による分類性能(学習データにおける交差検定)

表 4.3: 各手法の分類性能

	Precision	Recall	F-measure	概算獲得文数
All-positive	0.540	1.000	0.701	8928
BOW	0.706	0.631	0.667	5638
Match&behind_8	0.741	0.720	0.730	6432

本手順の評価を行う。図 4.3 にベースラインおよび提案手法の各評価値を示す。Recall は、パターンマッチによる抽出後の予定変更情報全体をもとに算出している。ベースラインとしては、全てを正例 と評価する場合(All-positive)と、機械学習において全品詞の Bag-of-words で行った場合(BOW)を採用した。また、オンライン学習器においては学習データの学習順序によって結果が異なるため、5 回の分類評価値の平均によって評価している。その結果、提案手法(Match & behind .8)は F 値において、All-positive より 0.029,BOW より 0.063 高い結果となった。

また、機械学習の分類において、ある閾値以上の分類確率によって正例に分類さ

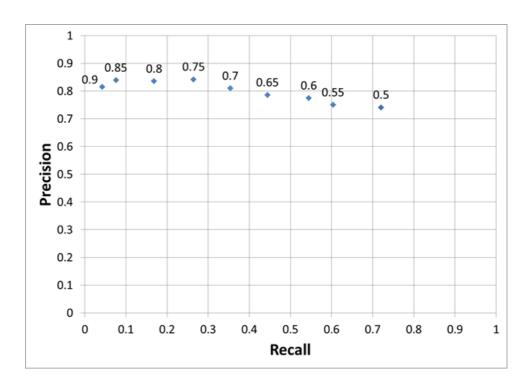


図 4.2: 設定する分類確率閾値による各評価値の変化

れているものを予定変更情報として獲得する場合の、各閾値による評価値を図 4.2 に示す。上記の場合と同様に、評価値は 5 回の分類の平均値を用いている。 閾値を上げることで、 閾値 0.75 とした場合には Precision は 0.842 まで上昇した。

4.5 獲得した予定変更情報の分析

前節までで述べた提案手法を用いて、コーパス全体から予定変更情報を獲得し、分析を行った。図 4.3 は獲得した予定変更情報を、ドキュメント作成月別に分け、その文数をプロットしたものである。この図から、2011年3月から数ヶ月間の間、他の月よりも文数が多いことがわかる。これは東日本大震災という事故が発生したことで、多くの予定や計画が日程変更や中止されたことによると考えられる。

また、図4.4は1月から12月の各月での文数の17年間の平均をとったものである. この図から、12月は他の月に比べて予定変更情報が多いことわかる.これは、12月

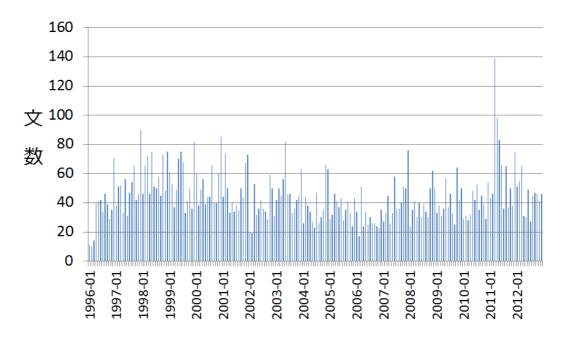


図 4.3: ドキュメント作成月別の獲得文数

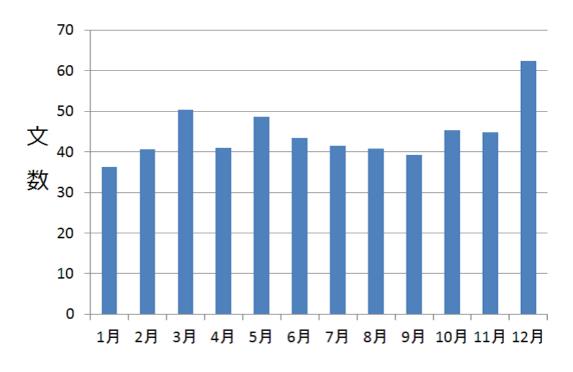


図 4.4: 月別の獲得文数 (17年平均)

は例年、特別国会や臨時国会の会期末になること、次年度の予算編成などが行われること、予定などの期限として設定されやすいことなどが理由として考えられる.

第5章 予定変更による影響の獲得

5.1 手法の目的および取り組むべき課題

本節では、ある予定が変更された際の影響をより詳細に把握することを目的として、予定変更情報によって影響を受ける予定を獲得する手法を提案する. 具体的なタスクとしては、前節で獲得した予定変更情報を入力として、未来関連情報の中から影響を受けるものを出力とする.

具体的な例を用いて説明する. 例えば,2011年05月に以下のような予定変更情報が獲得されている. 菅直人首相は3月31日に、福島第1原発の事故を踏まえ、2030年までに原発を現状より14基以上増やすとした政府のエネルギー基本計画を白紙にして見直す方針を表明している。

この予定変更情報から 2030 年のエネルギー基本計画が変更になったことがわかるが、その他にも多くの予定が影響を受けたことが考えられる. 例えば、原発の建設計画、再生可能エネルギーの普及、新技術開発、海外への原発輸出政策、地球温暖化対策としての CO2 削減目標等である. このような影響を受ける予定を獲得することが本手法の目的である.

本手法の提案における主な課題を二つ挙げる.まず一つ目は、出力の多様性を確保することである.本タスクにおいてまず考えられる手法は語句の類似度を尺度として、類似度の高い情報を獲得するというものだが、それでは予定変更情報と似た語句を持った情報以外は獲得することはできない.そこで、語句類似度以外の手法が必要となる.

二つ目は、影響の"方向"を考慮した手法設計が望まれるということである。例えば、2014年4月の"消費税率引き上げ"に伴い、"鉄道の運賃値上げ"が行われる見通しである。この際、"消費税率の引き上げ"が中止されれば、"鉄道の運賃値上

表 5.1: 前提関係を表す表現

動詞の仮定形+ば(すれば、ならば etc.) に向けて、に向け、を目指し、に伴い

げ"も中止になることは十分に考えられる.一方で,"鉄道の運賃値上げ"が中止になったとしても,"消費税率引き上げ"が中止になるとは考えにくい.このように,二つのイベントの間の影響伝播関係の方向性も考える必要がある.

5.2 提案手法

5.2.1 前提表現を含む文を用いた分類器の構築

本研究では、テキストコーパス内にある"前提関係を表す表現"に着目した.本研究では、表 5.1 の表現を"前提関係を表す表現"とした.以下ではこれらの表現を前提関係表現と呼ぶこととする.前提関係表現を含む文の例を以下に挙げる.

• 前提関係を含む文の例

- "消費税率の引き上げに伴い、鉄道の運賃を値上げする."
- "オリンピックの開幕に向けて、スタジアムの建設が進む。"
- "不信任案が成立すれば、大統領は内閣総辞職か議会解散のどちらかを 迫られる."
- "原発が再稼働すれば、節電目標幅を縮小する。"
- "再生可能エネルギー拡大に向け、大規模太陽光発電や新エネ・省エネ 設備の規制を緩和する."
- "環境にやさしい再生可能エネルギー技術の開発・普及促進を目指し、学術的な意見交換を行う。"

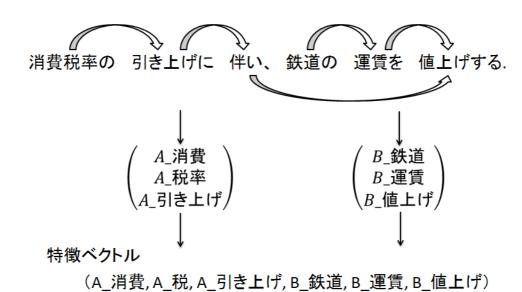


図 5.1: 係り受け解析を用いた学習データの作成

- "ソニーは放送設備分野でも世界的に圧倒的なシェアを持っており、これに放送事業が加われば、次世代テレビや衛星放送受信機などの標準規格競争で有利になる。"

乾らの研究 [29] や中島らの研究 [30] でもこのように手がかり表現に着目した研究が行われている。しかしこれらの論文では因果関係やイベント連鎖の獲得に焦点を当てており、さらに手がかりとしている表現も異なっている。Radinsky らの研究 [17] においても因果関係を表す表現に着目し、ニュースにおける未来のイベントを予測する研究が行われている。しかし、本研究で焦点を当てている予定変更による影響関係だけではなくより一般的なイベントの因果関係に着目している点や、テキストデータ内から獲得するのではなく、推論規則による予測を目的としている点で、本研究とは異なっている。

本手法では、これらの前提関係表現を含む文を用いた教師なし学習により分類器 を構築する、分類器構築の手順を以下で説明する。

係り受け解析結果から各文の特徴量を作成する手順の概要を図5.1に示す.まず, 前提関係表現を含む文の係り受け解析を行う.本研究の係り受け解析では吉永らが 開発した J.DepP 1 を用いた. 前提関係表現に係っている文節のうち,前提関係表現から 2 つ以内の係り受け関係にあるものを抽出し,形態素解析器 MeCab を用いて形態素解析した上で,それらを"前提前"の特徴量とする. ただし,本手法では動詞と名詞のみを特徴量の対象とした. また,前提関係表現が係っている文節を見つけ,その文節に係っている文節のうち 2 つ以内の係り受け関係にあるものを含めて"前提後"の特徴量とする. そして,これらを前提前と前提後で管理した上で統合し,正例の特徴量とする. 具体的には前提前の特徴量の前には" 2 点、前提後の特徴量の前には" 2 。

学習に用いる負例の作成方法については以下の手順で行う.二つの前提関係表現を含む文をランダムに選び出し、上記の手順で特徴量ベクトルを作成する.そして、一方の前提前の特徴量と他方の前提後の特徴量を統合し、それを負例の特徴量とした.これを負例が正例と同数になるまで繰り返した.

上記の手順で作成した正例と負例から,前節と同様オンライン学習器 opal を用いて分類器を構築した.本手法では,前提前と前提後の特徴量の組み合わせが大きな手がかりとなるため,2次の多項式カーネルを用いた.

また、より精度の高い前提関係を捉えため、本研究では一般的な単語("する"、 "ある"etc.) は分類器を構築する前に除外した. 具体的には、各単語の idf 値を以 下の式により算出し、idf 値が 5.5 以下のものを除外した.

$$idf_i = log(M/m_i) (5.1)$$

ただし、M はテキストコーパスの全文数、 m_i は単語 i の出現する文数である.これにより、テキストコーパス全体にある約 13 万種類ある動詞と名詞のうち、最も一般的な約 400 種類の単語を除外した.

5.2.2 分類器による入力のスコア付け

前節で構築した分類器を用いて、予定変更情報による影響を獲得する.まず、入力として、各々の予定変更情報と全ての未来関連情報のペアを作成する.ただし、本

¹http://www.tkl.iis.u-tokyo.ac.jp/ ynaga/jdepp/

研究では、予定変更情報がテキストデータ内に出現した時点で、それ以前に作成された未来関連情報から、予定の影響を獲得することを想定している。そこで、各予定変更情報について、予定変更情報が作成された月の前月までの未来関連情報のうち、予定変更情報が作成された月の翌月以降の参照時間をもつ未来関連情報を獲得対象とする。

上記条件で作られた予定変更情報と未来関連情報のペアを特徴量ベクトルに変換する.まず、予定変更情報を形態素に分割し、動詞と名詞を抽出する.そして、それらの前に "A_"とつけて管理する.次に、同様に未来関連情報の動詞と名詞を抽出し、それらの前に "B_"とつけて管理する.最後に両者を統合して一つの特徴量ベクトルとし、分類器の入力とする.ただし、前節と同様に一般的な単語については除外しておくこととする.

本研究では、各予定変更情報と未来関連情報を分類器にかけ、分類器によって算出されたマージンを各入力のスコアとした。マージンにしきい値を設け、しきい値以上のスコアをもつ入力を、変更情報により影響を受けるものとして出力する.

5.3 実験

5.3.1 提案手法の評価

2011年1月から2012年12月における予定変更情報を用いて実験を行った. 前章の手法により、本期間には1587文の予定変更情報が獲得されていた.

また、分類器を構築する際は、各予定変更情報の前月までのテキストコーパスを用いることとし、各月ごとに別々の分類器を用いた。これは、予定変更後の情報を使うことで本来なかったはずの前提条件を利用してしまうことで不公平性が生じないようにするためである。前提表現を含む文は、1996年から2012年までの間に約47万文存在した。機械学習には前節同様オンライン学習器opalを用い、多項式カーネルを用いた。

提案手法を評価する上での比較手法は、語句類似度を用いた手法とする. 本手法では、各文を MeCab で形態素に分割し、TF-IDF ベクトルに変換した後、そのコサ

表 5.2: 提案手法の評価

語句類似度	提案手法
1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1	ル 条子仏
0.45	0.60

イン類似度を各入力のスコアとした. TF-IDF の式を以下に示す.

$$tfidf_{i,j} = n_{i,j} * log(M/m_i)$$
(5.2)

 $n_{i,j}$ は単語 i の文 j における出現数,M は全文数, m_i は単語 i の出現する文数である。ただし,これまでと同様,品詞は名詞と動詞を用い,idf 値が 5.5 以下のものは除いた。

評価方法としては、提案手法と比較手法においてスコアの高いものから同数の結果を獲得対象として、その中からランダムに選んだ50個にラベル付をして、その特度を比較した。両者の獲得数は、3万(一つの予定変更情報あたり約20文)とした。

提案手法、比較手法における精度を表 5.2 に示す. ただし, 本手法の評価に焦点をおくため、正規化や予定変更情報獲得における誤りがあった場合はそれらを除き、合計 100 個になるまで評価を行った. これにより, 提案手法の方が比較手法よりも獲得性能において向上していると考えられる.

5.3.2 変更情報からの影響獲得例

本節では提案手法によって獲得された影響について、例を示す.以下の予定変更情報から獲得された未来関連情報は26個あったが、そのうちの一部を表5.3に示す.

- 入力とした予定変更情報(1)
 - 菅直人首相は3月31日に、福島第1原発の事故を踏まえ、2030年までに 原発を現状より14基以上増やすとした政府のエネルギー基本計画を白紙 にして見直す方針を表明している.

同様に、以下の予定変更情報から獲得された未来関連情報の一部を表 5.4 に示す。

表 5.3: 獲得された影響の例 (エネルギー基本計画の白紙)

世界のエネルギー需要が 2030 年に現在の 50%増となるとの予測も紹介したうえで、原発は二酸化炭素排出量が少ないことなどから、原子力は温暖化とエネルギー問題の「中核的解決手段の一つとなりうる」と位置づけた。

同法案は温室効果ガス削減中期目標を「2020年までに90年比で25%削減」と 定め、国内排出量取引▽地球温暖化対策税▽再生可能エネルギーの全量固定価 格買い取り制度を3本柱として施策の検討や実施を盛り込んでいた。

国内の原発関連メーカーは、老朽化した原発の建て替え需要が見込める 2030 年ごろまで輸出でつなぐしかない情勢で、経済産業省も総合資源エネルギー調査会の原子力部会で輸出政策を議論する見通しだ。

同委員会は、2020年代を目標に行われる最終処分地の選定過程で、処分の実施主体となる新たな認可法人「原子力発電環境整備機構」が、関係地方自治体に事前に十分な情報を公開し自治体と意思疎通を図るなど選定作業の透明性と公正さの確保を求める付帯決議を行った。

二酸化炭素排出量の削減促進を目的に、2400億円の増税となる地球温暖化対 策税を来年10月から段階的に導入。

- 入力とした予定変更情報(2)
 - 東日本大震災で大きな被害を受けた岩手,宮城,福島の3県について,7 月24日に予定されていた地上デジタル放送の完全移行を最大1年間延期 すると総務省が発表した.

5.3.3 影響獲得数による予定変更情報のランキング

提案手法によって獲得された影響の数最も多かった予定変更情報の上位 10 文を表 5.5 に示す. ランキングの上位に,地球温暖化に関する変更情報がいくつか見られるが,これは CO2 削減目標等の多くが未来の時間 ("2020年","2030年"etc.)とともに語られることから,この結果は妥当であると言える.また,エネルギー関連に関しても同様のことが言える.

表 5.4: 獲得された影響の例 (エネルギー基本計画の白紙)

原口一博総務相は16日、共同通信本社で開かれた放送協議会総会で講演し、約1年後に迫った地上デジタル放送への完全移行について「延期の選択肢は今はまったく考えていない」と、予定通り来年7月24日にアナログ放送を終了する考えを強調した

総務省や放送局、家電メーカーなどでつくる「地上デジタル推進全国会議」は 1日、ケーブルテレビ事業者が地上デジタル放送をアナログ信号に変換して送 信し、アナログテレビでも視聴できるようにする暫定措置「デジアナ変換」の 期間を、15年3月末と決めた。

地上テレビ放送のデジタル化に伴う周波数の大規模再編で、総務省は14日、2011年7月に停波して空きが出るアナログ放送の周波数を、携帯電話などの通信サービスやITS、携帯端末向け放送などに割り当てることを決めた。

総務省の情報通信審議会の情報通信政策部会は23日、2011年7月24日に地上アナログテレビ放送が終了し地上デジタル放送に移行することに伴い、生活保護世帯に地デジ対応の専用チューナーを現物給付する答申をまとめた。

受信可能世帯は、当初は放送局からの電波が届く範囲に限られるが、中継局を順次増設し、11年7月までには、東北全域で受信可能になる予定。

増田寛也総務相は9日、2011年7月24日に現行の地上アナログ放送が停止し、地上デジタル放送に完全に移行するのに合わせ、難視聴対策や経済的弱者支援など総額約2000億円の総合対策を09年度以降約3カ年で行うと発表した。

次に、この結果から本手法における今後の課題を考察する。表中の上位のものののいくつかでは、実際に延期や中止になったのは文の一部であるものがある。例えば"民主、自民、公明の~"の文では、延期になったのは所得税と相続税である。しかし、その後に自動車取得税、重量税など延期とは関係のない単語も含んでおり、本手法ではこれらも手掛かりとして影響を獲得してしまっている。これらによって精度が低下していることが考えられる。予定変更文において変更と関係のある単語を特定することで、手法の向上が図れると考える。

さらに、今後の課題としては、国や意思決定主体を特定することが挙げられる. 例えば、"イタリアのモンティ政権は~"の文では、日本の政権での変更と同じよう に影響を獲得してしまっていると考えられる. 国が違えば、影響の大きさも異なっ

表 5.5: 影響獲得数によるランキング

獲得数	予定変更情報
897	南アフリカで開かれた国連気候変動枠組み条約第17回締約国会議は、来年
001	末に期限切れを迎える京都議定書を延長する一方、20年に中国や米国など
	すべての国が参加して温室効果ガスを削減する新体制への道筋を付けた。
693	民主、自民、公明の3党合意で今年末に結論を先送りした所得税、相続税の
0.55	最高税率の見直しや、14年4月の消費税率8%への引き上げまでに結論を出
	すとした自動車取得税・重量税の存廃、増税の際の影響が大きい住宅取得の
	負担軽減策などが焦点となる。
531	国家公務員給与を平均7・8%削減する法案をはじめ、野田佳彦首相が増税の
991	前提に掲げた歳出削減と税外収入の確保策の多くが、来年の通常国会以降に
	先送りされた。
455	会期内の8月中に法案が成立すれば、2カ月後の10月から給与カットが実現
400	し、11 年度だけで約 1500 億円が復興財源に充てられるはずだった。
388	このほか、地球環境保全のために国際社会が達成すべき新たな数値目標「持
300	続可能な開発目標 を2015年までに策定することが盛り込まれたが、対象
	となる項目の具体化は先送りされた。
361	イタリアのモンティ政権は18日、これまで2013年としていた財政収支均衡
301	達成の目標年を、2年先送りし15年と修正する経済財政見通しを閣議決定し
	た。
337	11月の臨時国会で小選挙区の「0増5減」の関連法は成立したものの、衆院
331	の抜本的な定数削減は来年の通常国会会期末までに結論を出すとして先送り
	しており、議員自らのリストラを求める声も高まりそうだ。
325	前原氏は訪中時に楊潔外相らと会談し、来年の日中国交正常化40周年に向
020	けた日中関係改善策や、中国が一方的に延期を通告した東シナ海ガス田共同
	開発の条約締結交渉の再開問題、朝鮮半島情勢などについて協議する考え。
292	野田佳彦首相は19日の参院社会保障と税の一体改革特別委員会で、民主党
202	がマニフェストに掲げ、政府が2月、来年の国会への提出を閣議決定してい
	た新年金制度関連法案について、来年の国会提出を見送る意向を示した。
291	一方、関電は、停止中の原発の再稼働時期が未定であることや、今後の販売
	電力量の見極めが困難であることなどから、13年3月期通期の業績予想の公
	表を引き続き見送った。
	71 - 41 - 120 - 120

てくるため、それらを考慮することが必要であると考える.

第6章 おわりに

本論文では、テキストデータから予定変更情報を獲得し、さらにその予定変更情報が与える影響を獲得する手法を提案した.

まず、テキストデータ内の時間表現を正規化し、未来関連情報を抽出した.毎日新聞データとブログデータにおいて、新聞に特に多く見られる表現がある等、出現する時間表現に違いがあることがわかった.また、文書作成時間から参照時間までの距離をプロットすることで、毎日新聞の方が過去から未来まで幅広い時間表現が出現することがわかった.さらに、文書作成時間、参照時間、文数の3次元グラフによって可視化をすることで、計画的なイベント、突発的イベント、予定変更があった時などにその時間出現パターンにどのような影響があるかについて考察をした.

次に未来関連情報から予定変更情報を獲得する手法を提案した。まず約2000文におよび毎日新聞データを人手で確認し、予定変更に特有の表現を40パターン収集した。実験により、この40パターンで約82.2%の予定変更情報をカバーできることがわかった。その後、機械学習によるフィルタリングを行った。本研究では不確定な情報や否定の情報を除去するため、予定変更に特有の表現の直後を特徴量とする手法を提案した。実験により、提案手法によりF値が向上することを示し、手法の有効性を証明した。そして、分類確率値にしきい値を設定することで、分類精度が向上することを示した。さらに、獲得した予定変更情報を分析し、東日本大震災のあった2011年3月には多数の予定変更情報があったこと、そして月別にみると12月に予定変更が多いことがわかった。

また、予定変更情報によって影響を受ける未来関連情報を獲得する手法を提案した. 語句類似度を用いる手法では、多様な情報の獲得、影響の方向性の考慮において問題があることから、本研究では、テキストコーパス内の前提条件を含む文に着

目した. 前提条件を含む文から教師なし学習を用いて分類器を構築し、各予定変更情報と未来関連情報のペアを入力として算出したマージンがしきい値以上のものを出力とした. 実験の結果、語句類似度を用いた手法よりも高い獲得性能を示し、提案手法の有効性を示した.

参考文献

- [1] A.Jatowt and C.A.Yeung. Extracting collective expectations about the future from large text collections. *CIKM'11*, pages 909–918, 2011.
- [2] A.Jatowt, H.Kawai, K.Kanazawa, K.Tanaka, K.Kunieda, and K.Yamada. Analyzing collective view of future, time-referenced events on the web. In 19th International World Wide Web, 2010.
- [3] A.Jatowt, K.Kanazawa, S.Oyama, and K.Tanaka. Supporting analysis of future-related information in news archives and the web. In Proceedings of JCDL'2009, 2009.
- [4] A.Qamra, B.Tseng, and E.Chang. Mining blog stories using community-based and temporal clustering. In proceedings of the 15th ACM International Conference on Information and Knowledge Management (CIKM'06), 2006.
- [5] B.Shaparenko, R.Caruana, J.Gehrke, and T.Joachims. Identifying temporal patterns and key players in document collections. In Proceedings of the IEEE ICDM Workshop on Temporal Data Mining: Algorithms, Theory and Applications, 2005.
- [6] D.Koen and W.Bender. Temporal augmentation of the news. *In IBM Systems Journal*, 2000.
- [7] D.Koen and W.Bender. Temporal augmentation of the news. *In IBM Systems Journal*, 2000.

- [8] D.Metzler, C.Cai, and E.H.Hovy. Structured event retrieval over microblog archives. *NAACL*, 2012.
- [9] D.Shahaf, C.Guestrin, and E.Horvitz. Trains of thought: Generating information maps. WWW '12, 2012,.
- [10] F.Schilder and C.Habel. From temporal expressions to temporal information. In Proceedings of the Workshop on Temporal and Spatial Information Processing (TASIP'01), 2001.
- [11] G.Dias, R.Campos, and A.Jorge. Future retrieval: What does the future talk about? SIGIR 2011 Workshop on Enriching Information Retrieval, 2011.
- [12] H.Kawai, A.Jatowt, K.Tanaka, K.Kunieda, and K.Yamada. Chronoseeked: Search engine for future and past events. ICUIMC 2010 SKKU, 2010.
- [13] I.Mani and G.Wilson. Robust temporal processing of news. MUC 1998, 1998.
- [14] J.Allan, R.Gupta, and V.Khandelwal. Temporal summaries of new topics. In proceedings of the 24th Annual International ACM SIGIR Conference on Reasearch and Development in Information Retrieval (SIGIR '01), 2001.
- [15] J.Pustejovsky, J.M.Castano, R.Ingria, R.Sauri, R.J.Gaizauskas, A.Setzer, G.Katz, and D.R.Radev. Timeml: Robust specification of event and temporal expressions in text. In proceedings of the AAAI Spring Symposium on New Directions in Question Answering, 2003.
- [16] J.Pustejovsky and M.Verhagen. Semeval-2010 task 13: Evaluating events, time expressions, and temporal relations. *Proceedings of the NAACL HLT Workshop* on Semantic Evaluations, pages 851–860, 2009.
- [17] K.Radinsky, S.Davidovich, and S.Markovitch. Learning causality for news events prediction. WWW2012, 2012.

- [18] M. Verhagen, R. Gaizauskas, F. Schilder, M. Hepple, G. Katz, and J. Pustejovsky. Semeval-2007 task15: Tempeval temporal relation identification. In proceedings of the Four Int. Workshop on Semantic Evaluations, 2007.
- [19] M.Verhagen, R.Gaizauskas, F.Schilder, M.Hepple, J.Moszkowicz, and J.Pustejovsky. The tempeval challenge: identifying temporal relations in text. 2009.
- [20] N.Kanhabua, R.Blanco, and M.Matthews. Ranking related news predictions. SIGIR'11, 2011.
- [21] N.Yoshinaga and M.Kitsuregawa. Kernel slicing: Scalable online training with conjunctive features. *In proceedings of COLING*, 2010.
- [22] O.Alonso, J.Strotgen, R. Baeza-Yates, and M.Gertz. Temporal information retrieval: Challenges and opportunities. TWAW'11, 2011.
- [23] O.Alonso, M.Gertz, and R.Baeza-Yates. Clustering and exploring search results using timeline constructions. CIKM'09, 2003.
- [24] O.Alonso, M.Gertz, and R.Baeza-Yates. Clustering and exploring search results using timeline constructions. *CIKM'09*, 2009.
- [25] R.Baeza-Yates. Searching the future. In Proceedings of ACM SIGIR workshop MF/IR 2005, 2005.
- [26] R.Kessler, X.Tannier, and C.Hag'ege. Finding salient dates for building thematic timelines. *ACL*, 2012.
- [27] S.Ho, M.Lieberman, P.Wang, and H.Samet. Mining future spatiotemporal events and their sentiment from online news articles for location-aware recommendation system. ACM SIGSPATIAL MobiGIS'12, 2012.

- [28] W.Cui, S.Liu, L.Tan, C.Shi, Y.Song, Z.J.Gao, X.Tong, and H.Qu. Textflow: Towards better understanding of evolving topics in text. *InfoVis*, 2011.
- [29] 乾孝司, 乾健太郎, and 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.
- [30] 中島直哉, 吉永直樹, 鍛冶伸裕, 豊田正史, and 喜連川優. 時期依存性を有する イベント連鎖の獲得. **日本データベース学会論文誌** *vol.12*, *No.1*, *pp.103-108*, 2013.

発表文献

- 1. 栗原俊明,豊田正史,喜連川優. テキストデータの未来関連情報における予定変更情報の獲得に関する研究. 情報処理学会データベースシステム研究会,研究報告データベースシステム(DBS),2013-DBS-158(12),1-7 (2013.11.19).
- 2. 栗原俊明,豊田正史,喜連川優.テキストデータにおける予定変更情報の獲得および未来に起こり得る派生事象の予測.電子情報通信学会第6回データ工学と情報マネジメントに関するフォーラム/第12回日本データベース学会年次大会(DEIM2014)(2014.03).(to appear)