

(19) 日本国特許庁(JP)

(12) 特許公報(B2)

(11) 特許番号

特許第5717103号  
(P5717103)

(45) 発行日 平成27年5月13日(2015.5.13)

(24) 登録日 平成27年3月27日(2015.3.27)

(51) Int. Cl. F 1  
**G 0 6 F 17/30 (2006.01)**  
 G 0 6 F 17/30 3 5 0 C  
 G 0 6 F 17/30 2 2 0 Z  
 G 0 6 F 17/30 3 3 0 C

請求項の数 8 (全 16 頁)

(21) 出願番号	特願2012-143175 (P2012-143175)	(73) 特許権者	000004226 日本電信電話株式会社 東京都千代田区大手町一丁目5番1号
(22) 出願日	平成24年6月26日(2012.6.26)	(73) 特許権者	304021417 国立大学法人東京工業大学 東京都目黒区大岡山2丁目12番1号
(65) 公開番号	特開2014-6802 (P2014-6802A)	(74) 代理人	110001519 特許業務法人太陽国際特許事務所
(43) 公開日	平成26年1月16日(2014.1.16)	(72) 発明者	東中 竜一郎 東京都千代田区大手町二丁目3番1号 日 本電信電話株式会社内
審査請求日	平成26年6月9日(2014.6.9)	(72) 発明者	松尾 義博 東京都千代田区大手町二丁目3番1号 日 本電信電話株式会社内

最終頁に続く

(54) 【発明の名称】 文書間関係推定装置、方法、及びプログラム

(57) 【特許請求の範囲】

【請求項1】

時間情報及びユーザ情報が各々付与された複数の文書からなる文書群における全ての文書間の各々について、文書間の内容の関連性と、文書間の応答関係、及び文書間における何れか一方の文書が同一のユーザ情報が付与された文書のうちの先頭または最後の文書であるか否かを示す特徴を抽出する特徴抽出手段と、

前記文書群における全ての文書間の各々について、前記特徴抽出手段によって抽出された前記文書間に対する前記特徴と、前記特徴に基づいて文書間の関係を推定するための予め学習された推定モデルとに基づいて、前記文書間の関係を推定する関係推定手段と、

を含む文書間関係推定装置。

10

【請求項2】

文書群の各文書について、前記文書が他の文書と関係がある場合、前記文書と関係がある前記他の文書が1つである制約を生成する制約生成手段を更に含み、

前記関係推定手段は、前記文書群における全ての文書間の各々について、前記特徴抽出手段によって抽出された前記文書間に対する前記特徴と、前記制約生成手段によって生成された前記制約と、前記特徴及び前記制約に基づいて文書間の関係を推定するための予め学習された推定モデルとに基づいて、前記文書間の関係を推定する請求項1記載の文書間関係推定装置。

【請求項3】

前記特徴は、文書間の特徴を示す述語であって、

20

前記関係推定手段は、

前記推定モデルとして、前記特徴を示す述語又は前記推定される前記文書間の関係を示す潜在述語を用いて記述された各論理式を用いて構築されるマルコフ論理ネットワーク（MLN：Markov Logic Network）における各論理式の重みを用いて、前記マルコフ論理ネットワークにより、前記制約生成手段によって生成された前記文書間に対する前記制約を記述した論理式を満足し、かつ、前記特徴抽出手段によって抽出された前記文書間に対する前記特徴を示す述語に対して尤もらしい、前記文書間の関係を示す潜在述語を推定する請求項 2 記載の文書間関係推定装置。

【請求項 4】

前記文書間の特徴は、前記文書間のうちの先の文書が後の文書の返信先であること、後の文書が先の文書のユーザに対して返信していること、前記文書のユーザ情報が同じであること、先の文書が同じユーザ情報が付与された文書のうちの最後の文書であること、及び先の文書が同じユーザ情報が付与された文書のうちの最初の文書であることの少なくとも一つを含む請求項 1～請求項 3 の何れか 1 項記載の文書間関係推定装置。

10

【請求項 5】

時間情報及びユーザ情報が各々付与された複数の学習用文書からなる学習用文書群における全ての学習用文書間の各々について、前記特徴を抽出する学習用特徴抽出手段と、前記学習用文書群における全ての学習用文書間の各々について予め定められた文書間の関係と、前記学習用特徴抽出手段によって前記学習用文書群における全ての学習用文書間の各々について抽出された前記学習用文書間に対する前記特徴とに基づいて、前記推定モデルを学習する学習手段と、

20

を更に含む請求項 1～請求項 4 の何れか 1 項記載の文書間関係推定装置。

【請求項 6】

前記学習用文書群における各学習用文書について、前記学習用文書が他の学習用文書と関係がある場合、前記学習用文書と関係がある前記他の学習用文書が 1 つである制約を生成する学習用制約生成手段を更に含み、

前記学習手段は、前記学習用文書群における全ての学習用文書間の各々について予め定められた文書間の関係と、前記学習用特徴抽出手段によって前記学習用文書群における全ての学習用文書間の各々について抽出された前記学習用文書間に対する前記特徴と、前記学習用制約生成手段によって生成された各学習用文書に対する前記制約とに基づいて、前記推定モデルを学習する請求項 5 記載の文書間関係推定装置。

30

【請求項 7】

特徴抽出手段によって、時間情報及びユーザ情報が各々付与された複数の文書からなる文書群における全ての文書間の各々について、文書間の内容の関連性と、文書間の応答関係、及び文書間における何れか一方の文書が同一のユーザ情報が付与された文書のうちの先頭または最後の文書であるか否かを示す特徴を抽出し、

関係推定手段によって、前記文書群における全ての文書間の各々について、前記特徴抽出手段によって抽出された前記文書間に対する前記特徴と、前記特徴に基づいて文書間の関係を推定するための予め学習された推定モデルとに基づいて、前記文書間の関係を推定する

40

文書間関係推定方法。

【請求項 8】

コンピュータを、請求項 1～請求項 6 の何れか 1 項記載の文書間関係推定装置を構成する各手段として機能させるためのプログラム。

【発明の詳細な説明】

【技術分野】

【0001】

本発明は、文書間関係推定装置、方法、及びプログラムに係り、特に、文書間の関係を推定する文書間関係推定装置、方法、及びプログラムに関する。

【背景技術】

50

## 【0002】

インターネット上では様々なユーザが発言を投稿しており、たとえば、掲示板やマイクロブログサービスでは、ユーザが日々発言を投稿し、やりとりをしている。このようなデータは非常に膨大であり、構造化されていないため、効率的に閲覧することが難しい。そこで、発言間の関係を同定し、構造化する手法が提案されている。たとえば、非特許文献1では、QAサイトに投稿された発言間の関係性をマルコフロジックネットワーク(MLN)と呼ばれる教師あり学習の手法で関係づけを行っている。関係性としては、「類似」(発言同士が同様の内容を保持している)や「包含」(片方の発言がもう片方の内容を完全に含み、新たな内容も含んでいる)などである。

## 【0003】

非特許文献1では、発言の特徴、および、発言間の特徴から、発言間の関係を同定している。ここで用いられているのは発言の内容に基づく特徴である。具体的には、特徴として、どちらが長い、連続する発言かどうか、発言間の投稿間隔、反意語となる単語対が発言間にあるか、同じURLを含むかどうか、同じ固有表現を含むか否か、括弧で囲まれた同じ表現を含むかどうか、異なる固有名詞を含むかどうか、単語のコサイン類似度、名詞の包含度を用い、関係性の正解データから、関係性推定のモデルをMLNによって学習している。MLNについては、非特許文献2に詳述されている。MLNは、確率的に推論を行う仕組みとして、近年注目されている。

## 【0004】

MLNでは重み付きの述語を扱うことができ、このため、かならず成り立つわけではないような関係も論理的な関係と同時に扱うことができる。現実的な、おおそ成り立つ関係について、重みを学習によって決定し、推定に役立てることが可能な学習手法である。

## 【先行技術文献】

## 【非特許文献】

## 【0005】

【非特許文献1】Hikaru Yokono; Takaaki Hasegawa; Genichiro Kikui; Manabu Okumura, Identification of relations between answers with global constraints for Community-based Question Answering services, Proc. IJCNLP, 2011年.

【非特許文献2】吉川克正、浅原正幸、松本裕治、「Markov Logic による日本語述語項構造解析」、情報処理学会研究報告(NL-199)、2010年.

## 【発明の概要】

## 【発明が解決しようとする課題】

## 【0006】

インターネット上の発言には、発言対象のユーザ名などが入っていたり、掲示板やマイクロブログサービスによっては、発言者のユーザIDや発言同士が応答関係にあるかどうかといった関係性がシステムによって付与されていたりすることが多い。このような情報は、関係性の同定に有効だと考えられるが、従来は用いられておらず、関係性の推定性能が低くなる、という問題があった。

## 【0007】

本発明は、上記の事情を鑑みてなされたもので、文書間の関係を精度よく推定することができる文書間関係推定装置、方法、及びプログラムを提供することを目的とする。

## 【課題を解決するための手段】

## 【0008】

上記の目的を達成するために本発明に係る文書間関係推定装置は、時間情報及びユーザ情報が各々付与された複数の文書からなる文書群における全ての文書間の各々について、文書間の内容の関連性と、文書間の応答関係、及び文書間における何れか一方の文書が同一のユーザ情報が付与された文書のうちの先頭または最後の文書であるか否かを示す特徴を抽出する特徴抽出手段と、前記文書群における全ての文書間の各々について、前記特徴抽出手段によって抽出された前記文書間に対する前記特徴と、前記特徴に基づいて文書間の関係を推定するための予め学習された推定モデルとに基づいて、前記文書間の関係を推

10

20

30

40

50

定する関係推定手段と、を含んで構成されている。

【0009】

本発明に係る文書間関係推定装置は、文書群の各文書について、前記文書が他の文書と関係がある場合、前記文書と関係がある前記他の文書が1つである制約を生成する制約生成手段を更に含み、前記関係推定手段は、前記文書群における全ての文書間の各々について、前記特徴抽出手段によって抽出された前記文書間に対する前記特徴と、前記制約生成手段によって生成された前記制約と、前記特徴及び前記制約に基づいて文書間の関係を推定するための予め学習された推定モデルとに基づいて、前記文書間の関係を推定することができる。

【0010】

本発明に係る文書間関係推定方法は、特徴抽出手段によって、時間情報及びユーザ情報が各々付与された複数の文書からなる文書群における全ての文書間の各々について、文書間の内容の関連性と、文書間の応答関係、及び文書間における何れか一方の文書が同一のユーザ情報が付与された文書のうちの先頭または最後の文書であるか否かを示す特徴を抽出し、関係推定手段によって、前記文書群における全ての文書間の各々について、前記特徴抽出手段によって抽出された前記文書間に対する前記特徴と、前記特徴に基づいて文書間の関係を推定するための予め学習された推定モデルとに基づいて、前記文書間の関係を推定する。

【0011】

本発明に係るプログラムは、コンピュータを、上記の文書間関係推定装置の各手段として機能させるためのプログラムである。

【発明の効果】

【0012】

以上説明したように、本発明の文書間関係推定装置、方法、及びプログラムによれば、文書間の各々について、文書間の内容の関連性と、文書間の応答関係、及び何れか一方の文書が同一のユーザ情報が付与された文書のうちの先頭または最後の文書であるか否かを示す特徴を抽出し、予め学習された推定モデルに基づいて、文書間の関係を推定することにより、文書間の関係を精度よく推定することができる、という効果が得られる。

【図面の簡単な説明】

【0013】

【図1】本発明の実施の形態に係る発言間関係推定装置の構成を示す概略図である。

【図2】入力されるツイート集合の一例を示す図である。

【図3】本発明の実施の形態に係る発言間関係推定装置におけるモデル学習処理ルーチンの内容を示すフローチャートである。

【図4】本発明の実施の形態に係る発言間関係推定装置における発言間関係推定処理ルーチンの内容を示すフローチャートである。

【図5】実験結果を示す図である。

【発明を実施するための形態】

【0014】

以下、図面を参照して本発明の実施の形態を詳細に説明する。

【0015】

<システム構成>

図1に示すように、本発明の実施の形態に係る発言間関係推定装置100は、発言を示すテキストデータの集合が入力され、各発言間の関係を出力する。1つの発言は1つ以上の文からなるテキストデータである。この発言間関係推定装置100は、CPUと、RAMと、後述するモデル学習処理ルーチン及び発言間関係推定処理ルーチンを実行するためのプログラムを記憶したROMとを備えたコンピュータで構成され、機能的には次に示すように構成されている。図1に示すように、発言間関係推定装置100は、入力部10と、演算部20と、出力部30とを備えている。

【0016】

入力部 10 は、入力された発言の集合を受け付ける。本実施の形態の例では、発言の集合として、マイクロブログサービスの一つであるツイッター（R）の発言集合を用いる。ツイッター（R）では、ユーザが日々膨大な数の発言（ツイートと呼ばれる）を行っており、構造化が望まれるデータのの一つである。ここでは、ツイートの集合を togetter（R）（<http://togetter.com/>）と呼ばれるサービスから収集した。togetter（R）は、個人が自身のお気に入りのツイートを「まとめ」として登録することのできるサービスで、一定のトピックに関するツイートが雑多に集められている。

【0017】

また、学習データとして入力された発言の集合と共に、入力部 10 は、入力された発言間の関係を受け付ける。入力される発言間の関係は、手動で付与されたものであり、使用する関係性は、例えば、「関係あり」である。

10

【0018】

また、使用する関係性として詳細なものを用いてもよく、例えば、「賛成」及び「反対」、並びに、「矛盾」、「類似」、「演繹」及び「帰納」である。賛成は同意を表し、反対は、同意していないことを表す。矛盾は、異なる内容を述べていることを表し、類似は、発言に同意をして、類似した内容を述べていることを表す。演繹は、発言の内容をもとにして、推論による内容を発言したり、新たな情報を加え議論の展開することを表し、帰納は、発言の内容をまとめた内容を発言していることを表す。

【0019】

本実施の形態の具体例として、togetter（R）における14のまとめページから、14のツイート集合を収集した。また、実験のため、これらの集合をそのまま用いるのではなく、まとめのトピックとは直接関係のないと思われるツイートは削除し、リプライ先が含まれていないツイートがあった場合には、twitter（R）から改めて取得するという処理を行った。図 2 は、あるツイート集合からの抜粋である。IDはツイートのID、ユーザ名はツイートをしたユーザの名前、返信先は、ツイッターによって付与される返信先情報であり、in\_reply\_toから取得できる。本文はツイートの発言そのものであり、リンク先IDは、in\_reply\_to からでは判定できないが、内容から確認できる、発言対象を表すIDである。賛否は、関係性のある発言との間の関係としての「賛成」「反対」のいずれかであり、関係は、関係性のある発言との間の関係としての「矛盾」「類似」「演繹」「帰納」のいずれかである。リンク先ID、賛否、関係は、学習データとして人手で付与した発言間の関係の一例である。

20

30

【0020】

また、入力部 10 は、上記の関係性がツイート間にあるかないか、あるとしたらどの関係かを推定するために入力された、発言の集合を受け付ける。

【0021】

演算部 20 は、発言集合データベース 21、発言間関係データベース 22、特徴量生成部 25、モデル学習部 26、モデル記憶部 27、入力発言集合データベース 28、特徴量生成部 29、及び関係同定部 31 を備えている。なお、関係同定部 31 が、関係推定手段の一例である。

【0022】

発言集合データベース 21 は、入力部 10 により受け付けた学習データとしての発言の集合を記憶する。発言間関係データベース 22 は、入力部 10 により受け付けた学習データとしての発言間の関係を記憶する。

40

【0023】

特徴量生成部 25 は、形態素解析部 251、固有表現抽出部 252、発言間特徴生成部 253、及び発言間制約生成部 254 を備え、発言集合データベース 21 に記憶されている発言の全ペアについてペア間の特徴を生成する。本実施の形態の例では、発言間の特徴は一階述語論理として表現される。特徴量は、述語として表される。

【0024】

形態素解析部 251 は、各発言に対して形態素解析を行う。形態素解析の手法は、従来

50

既知の手法を用いればよく、入力を形態素に分割できるものであれば何でもよい。固有表現抽出部252は、各発言から固有表現を抽出する。固有表現抽出の手法は、従来既知の手法を用いればよく、入力について、固有表現を抽出できるものであれば何でも良い。本実施の形態の例では、どちらについてもCaboChaを用いている。

【0025】

発言間特徴生成部253は、ツイートの発言間(ツイート*i*とツイート*j*)の各々について、後述する「in\_reply\_to」、「reply」、「sameuser」、「latestutt」、及び「firstutt」を含む述語を、発言間の特徴(素性とも呼ぶ)として生成する。

【0026】

「in\_reply\_to」: ツイート*j*にリプライ先のツイートのID=*i*が指定されている場合、特徴として、述語in\_reply\_to(*i*, *j*)を生成する。たとえば、先の例における、ツイート2\_userB\_1はin\_reply\_toにID 1\_userA\_1が指定されているため、述語in\_reply\_to(1(1\_userA\_1), 2(2\_userB\_1))が特徴として生成される。

10

【0027】

「reply」: ツイート*j*が@...という形でツイート*i*のユーザー名に言及している場合、reply(*i*, *j*)を生成する。たとえば、ツイート*j* = 44\_userB\_10の中で@で言及されているユーザーのツイートのIDを用いて、述語reply(29(29\_userX\_17), 44(44\_userB\_10))、述語reply(33(33\_userX\_19), 44(44\_userB\_10))といった特徴を生成する。

【0028】

「sameuser」: ツイート*j*とツイート*i*が同一ユーザーの場合、述語sameuser(*i*, *j*)を生成する。たとえば、2\_userB\_1と4\_userB\_2のように、同じユーザーのツイート同士であれば、述語sameuser(2(2\_userB\_1), 4(4\_userB\_2))を生成する。

20

【0029】

「latestutt」: ツイート*j*とツイート*i*の間にツイート*i*と同じユーザーのツイートがない場合、つまり、相手の一番最近のツイートに対してリンクがある場合、述語latestutt(*i*, *j*)を生成する。たとえば、ツイート2のユーザーの発言が、ツイート3の時点で最も最近のそのユーザーの発言である場合、述語latestutt(2(2\_userB\_1), 3(3\_userC\_1))を生成する。すなわち、会話中で誰かの発言に対して応答を行う場合、相手の新しい発言を無視して、以前の発言に対して応答することは少ないだろう、という状況を表す意図で、ツイートが発言された時点で、応答相手の最新の発言に対しての応答かどうかを表す述語を生成する。

30

【0030】

「firstutt」: ツイート*j*がリンクをもつツイート*i*が所定の区間内でそのユーザーの最初の発言である場合、述語firstutt(*i*, *j*)を生成する。firstuttは、相手のユーザーの最初のツイートとの組み合わせであることを表す述語である。たとえば、ツイート*i*のIDが「i\_ユーザー名\_1」となっている場合、述語firstutt(*i*, *j*)を生成する。質問や問題提起など、会話の発端となるようなツイートが広く様々なユーザーから言及される場合がtogether(*R*)上で多く見られるため、together(*R*)のまとめ上でそれぞれのユーザーの初めての発言はより多くのユーザーから応答されやすいだろう、という意図でこの述語を生成する。なお、所定の区間内というのは、応答関係を判断するツイート集合内、例えば、together(*R*)のまとめ全体の範囲を表わす。

40

【0031】

なお、述語「in\_reply\_to」が、先の文書が後の文書の返信先であることを示す特徴の一例であり、述語「reply」が、後の文書が先の文書のユーザーに対して返信していることを示す特徴の一例であり、述語「sameuser」が、文書のユーザー情報が同じであることを示す特徴の一例である。また、述語「latestutt」が、先の文書が同じユーザー情報が付与された文書のうちの最後の文書であることを示す特徴の一例であり、述語「firstutt」が、先の文書が同じユーザー情報が付与された文書のうちの最初の文書であることを示す特徴の一例である。

【0032】

50

また、発言間特徴生成部 2 5 3 は、従来法で素性として用いられた、発言ペアのどちらが長い、発言ペアが連続する発言かどうか、発言間の投稿間隔、反意語となる単語対が発言間にあるか、発言ペアが同じ URL を含むかどうか、発言ペアが同じ固有表現を含むか否か、発言ペアが括弧で囲まれた同じ表現を含むかどうか、発言ペアが異なる固有名詞を含むかどうか、発言間における単語ベクトルのコサイン類似度、発言間の名詞の包含度の各々を表わす述語を、発言間の特徴として、発言間の各々について生成する。

## 【 0 0 3 3 】

固有表現は、人名、製品名、施設名、地名、時間表現、数値表現を扱っており、個々の固有表現ごとの素性ではなく、固有表現の種類毎に固有表現を含むかどうかを表わす述語を生成する。反意語に関する述語は、予め準備した反意語リストに従って、一方が、もう一方の反意語を含むかを表す述語を生成する。たとえば、好評に対して悪評と不評、重んじるに対して軽んじる、夏至に対して冬至、などが反意語のリストである。発言間の間隔については離散値の特徴であり、出現順の差が、5以下、10以下、15以下、20以下、30以下、40以下、50以下、それ以上であるかどうかを表わす述語を生成する。名詞の包含度はツイート  $i$ 、 $j$  間の両方向に対して定義され、ツイート  $i$  のツイート  $j$  に対する包含度は、(共通する名詞の異なり数) / (ツイート  $i$  に出現した名詞の数) であり、発言間の間隔と同様に、離散値の特徴として述語を生成する。単語ベクトルのコサイン類似度は、それぞれのツイートについて単語 unigram、単語 bigram それぞれのベクトルを作成し、そのベクトル間のコサイン類似度を計算し、発言間の間隔と同様に、離散値 (例えば、0.1 刻みの離散値) の特徴として述語を生成する。

## 【 0 0 3 4 】

発言間制約生成部 2 5 4 は、発言間の制約を示す論理式を生成する。ツイッターでは、一つの発言が短く、複数の発言に一度に回答することは少ない。よって、発言が他の発言と関係がある場合、関係がある他の発言が 1 つである、という制約を表す論理式を、各発言について生成する。例えば、以下に示す論理式が生成される。

## 【 0 0 3 5 】

```
for Id i
```

```
  if tweet(i) : |Id k:has_aa_relation(k,i)|<=1;
```

上記のように制約の論理式が記述される。これで、has\_aa\_relation( $k, i$ ) を満たす  $k$  が最大一つとなる。

## 【 0 0 3 6 】

また、発言間制約生成部 2 5 4 は、上記の制約に加え、基本的な制約として、推移律と呼ばれる、ツイート  $i$  と  $j$  にある関係  $R$  が成り立ち、ツイート  $j$  と  $k$  に同じく  $R$  が成り立つ場合、ツイート  $i$  と  $k$  にも同様の関係  $R$  が成り立つという制約を生成する。たとえば、あるツイートに類似している 2 つのツイート  $j, k$  は類似しているという制約を生成する。

## 【 0 0 3 7 】

例えば、以下に示す論理式である。

## 【 0 0 3 8 】

```
aa_relation(i,j,"類似") & aa_relation(j,k,"類似") => aa_relation(i,k,"類似")
```

上記の論理式は、 $i$  と  $j$  が類似、 $j$  と  $k$  が類似ならば  $i$  と  $k$  も類似であることを表わす。

## 【 0 0 3 9 】

上記の制約を表わす論理式があれば、他のツイート間の関係から別のツイート間の関係を推定することができる。

## 【 0 0 4 0 】

また、推移律ではない制約として、以下に示すような論理式を生成してもよい。

## 【 0 0 4 1 】

```
tweet(i) & !tweet(j) => !has_aa_relation(i,j)
```

上記の論理式は、ツイートとツイートでない定数の間には応答関係は存在しないことを表わす。

## 【 0 0 4 2 】

10

20

30

40

50

上記の制約を表わす論理式があれば、現実的に意味をもたない状況を出力から除くことができる。

【 0 0 4 3 】

モデル学習部 2 6 は、特徴量生成部 2 5 が出力した特徴（観測可能な述語の集合）及び発言間関係データベース 2 2 に記憶された発言間の関係（ラベル）を示す潜在述語の集合から、各述語及び潜在述語を用いて記述された各論理式の重みを学習する。学習には M L N を用いる。M L N では、ラベル間に対して予め記述した制約を表わす論理式、及び特徴量とラベルとの間に対して予め記述した制約を表わす論理式を満たしつつ、観測可能な述語について行われた推論結果が、正解データとして与えられた述語の集合に近づくように、各論理式の重みが決定される。ここで学習された論理式の重みの集合がモデルとなり、重みの集合を表わすモデルが、モデル記憶部 2 7 に記憶される。なお、ラベルは、発言間関係データベース 2 2 に記憶された発言間の関係である。

10

【 0 0 4 4 】

M L N の学習について簡単に説明する。M L N は、述語の集合 X に対してある確率を与える。この確率を最大化する述語の集合（潜在述語を含む）が、M L N における推論結果となる。M L N は、実際には一階述語論理式の集合として表される。論理式には、違反を許容する重み付きの論理式と、違反を許容しない論理式を混在させることができ、制約を表わす論理式は、違反を許容しない論理式として記述される。

【 0 0 4 5 】

述語の集合 X に対する確率は、述語を論理式中の変数に代入することで成立する（真となる）論理式すべてについて重みの和をとり、exponential を取り、正規化したものである。例えば、以下の（ 1 ）式で表される。

20

【 0 0 4 6 】

【数 1】

$$p(X) = \frac{1}{Z} \exp\left(\sum_i w_i n_i(X)\right) \dots (1)$$

30

【 0 0 4 7 】

ここで、Z は正規化のための定数であり、 $w_i$  が i 番目の論理式に対応する重みである。 $n_i$  は、ある述語の集合 X 上で真をとる論理式、例えば、`tweet(i) & tweet(j) & in_reply_to(i, j)` のそれぞれについて、カウントが 1 増える。論理式の一部のみが真でもカウントはされない。この場合、論理式 `tweet(i) & tweet(j) & in_reply_to(i, j)` が真になる回数  $n_i$  は、X 中に、この論理式が真になる  $i$  と  $j$  の組が存在する数と同じとなる。

【 0 0 4 8 】

なお、モデル学習部 2 6 は、素性として生成された観測可能な述語の集合、発言間関係データベース 2 2 に記憶された発言間の関係を表わす潜在述語（ラベル）を用いた論理式を、以下のように生成しておく。

40

【 0 0 4 9 】

例えば、ツイート  $i$ 、 $j$ 、 $k$  があるとき、以下のような述語が観測可能な述語として与えられる。

【 0 0 5 0 】

`tweet(i)` （ $i$  はツイートである）  
`tweet(j) tweet(k) in_reply_to(i, j)` （ $j$  が  $i$  に対してリプライしている）  
`in_reply_to(j, k)` （ $k$  が  $j$  に対してリプライしている）  
`has_span(i, j, "1-5")` （ $i$  と  $j$  の間隔が 1-5 の間）

【 0 0 5 1 】

加えて、`in_reply_to` 属性が付いていて、`has_aa_relation` が成り立っているならば（

50



真であれば)、 $\text{tweet}(i) \ \& \ \text{tweet}(j) \ \& \ \text{in\_reply\_to}(i, j) \Rightarrow \text{has\_aa\_relation}(i, j)$  という論理式を用意する。

【 0 0 5 2 】

また、上記のような論理式は、以下に説明するように、述語に対する論理式のテンプレートに従って、重みと共に用意される。

【 0 0 5 3 】

例えば、ツイート*i*とツイート*j*の間に $\text{in\_reply\_to}$ 属性がついていて、*i*、*j*間に $\text{has\_aa\_relation}$  が成り立っているならば(真であれば)、 $w_{\text{inreplyto}}$  という重みを得ることを示す以下のような記述が、人手によって与えられる。

【 0 0 5 4 】

$\text{if } \text{tweet}(i) \ \& \ \text{tweet}(j) \ \& \ \text{in\_reply\_to}(i, j) \ \text{add}[\text{has\_aa\_relation}(i, j)] * w_{\text{inreplyto}}$

【 0 0 5 5 】

これらは「 $\text{tweet}(i) \ \& \ \text{tweet}(j) \ \& \ \text{in\_reply\_to}(i, j)$ 」ならば「 $\text{has\_aa\_relation}(i, j)$ 」という意味の論理式がすべての*i*と*j*にあてはまるツイートに対して用意される。また、それぞれの重み $w_{\text{inreplyto}}$ がMLNに含まれる。

【 0 0 5 6 】

また、MLNの学習では、 $\text{has\_aa\_relation}(i, j)$ 、および、 $\text{has\_aa\_relation}(j, k)$ となったとすると、テンプレートで記述した論理式のうち、成立した論理式の重みを $w_{\text{inreplyto}} * 2 + w_{\text{hasspan}}$ のように計算し、確率を計算する。(  $w_{\text{inreplyto}}$  は*i*、*j*と*j*、*k*で二度真になっているため)このような、素性として与えられていない述語は潜在述語とよばれ、これが、いわゆる分類問題で言う出力するラベルにあたる。

【 0 0 5 7 】

本実施の形態の例では、正例(正解データの潜在述語を含む述語集合)と負例(それ以外の述語集合)間のマージンを最大化する学習を行って、各述語に対する論理式の重みを求める。なお、尤度を最大化する学習を行ってもよい。

【 0 0 5 8 】

入力発言集合データベース28は、入力部10により受け付けた推定対象データとしての発言の集合を記憶する。入力されたツイート集合、具体的には、ツイッターから取得できる情報である、ID、発言内容、 $\text{in\_reply\_to}$ 属性、及びユーザ名からなるデータの集合が、入力発言集合データベース28に記憶される。

【 0 0 5 9 】

特徴量生成部29は、上記の特徴量生成部25と同様に、形態素解析部251、固有表現抽出部252、発言間特徴生成部253、及び発言間制約生成部254を備え、入力発言集合データベース28に記憶されている発言の全ペアについてペア間の特徴を生成する。

【 0 0 6 0 】

関係同定部31は、入力発言集合データベース28に記憶されている発言のペアの各々について、生成された特徴(観察可能な述語集合)に対して、モデル記憶部27に記憶されたモデルを用いて、尤もらしい関係性ラベル(潜在述語)の集合を得る。具体的には、上記(1)式中の尤度を最大化する述語集合を得る。これらが推定結果であり、同定された関係である。

【 0 0 6 1 】

関係同定部31により推定された潜在述語が表わす発言の各ペアの関係性を、出力部30により出力する。

【 0 0 6 2 】

< 発言間関係推定装置の作用 >

次に、本実施の形態に係る発言間関係推定装置100の作用について説明する。まず、発言の集合と、各発言間に対して手動で付与した発言間の関係を示すラベルの集合とが発言間関係推定装置100に入力されると、発言間関係推定装置100によって、入力された、発言の集合が、発言集合データベース21へ格納され、入力された、発言間の関係を

10

20

30

40

50

示すラベルが、発言間関係データベース 22 に格納される。

【0063】

そして、発言間関係推定装置 100 によって、図 3 に示すモデル学習処理ルーチンが実行される。

【0064】

まず、ステップ S101 において、発言集合データベース 21 に格納された各発言に対して、形態素解析処理を行う。次のステップ S102 では、発言集合データベース 21 に格納された各発言から、固有表現を抽出する。

【0065】

そして、ステップ S103 において、発言集合データベース 21 に格納された発言の全ペアの各々に対して、発言間の特徴である述語を生成する。次のステップ S104 では、予め定められた制約を表わす論理式を生成する。

【0066】

ステップ S105 において、発言間関係データベース 22 に格納された発言の各ペアの関係を示すラベルを用いて、潜在述語を生成し、発言集合データベース 21 に格納された発言の全ペアの各々に対して観察可能な述語の集合及び潜在述語からなる学習データを生成する。

【0067】

そして、ステップ S106 において、上記ステップ S105 で生成した学習データに基づいて、上記ステップ S103 で生成した制約を表わす論理式を満足し、かつ、上記(1)式で表される確率が最大となる、発言の各ペアの関係を示す潜在述語及び各述語を用いた各論理式の重みを学習する。次のステップ S107 では、上記ステップ S106 で学習された各論理式の重みを、モデル記憶部 27 に格納して、モデル学習処理ルーチンを終了する。

【0068】

そして、発言間の関係を推定する推定対象の発言集合が発言間関係推定装置 100 に入力されると、発言間関係推定装置 100 によって、入力された発言集合が、入力発言集合データベース 28 へ格納される。

【0069】

そして、発言間関係推定装置 100 によって、図 4 に示す発言間関係推定処理ルーチンが実行される。

【0070】

ステップ S111 において、入力発言集合データベース 28 に格納された各発言に対して、形態素解析処理を行う。次のステップ S112 では、入力発言集合データベース 28 に格納された各発言から、固有表現を抽出する。

【0071】

そして、ステップ S113 において、入力発言集合データベース 28 に格納された発言の全ペアの各々に対して、発言間の特徴である述語を生成する。次のステップ S114 では、予め定められた制約を表わす論理式を生成する。

【0072】

ステップ S115 において、発言の全ペアの各々に対して、上記ステップ S113 で生成された発言間の特徴である述語の集合、上記ステップ S114 で生成された制約の論理式、及び学習されたモデル(各論理式の重み)に基づいて、当該ペアに対する潜在述語を推定することにより、当該ペアの関係を推定する。

【0073】

そして、ステップ S116 において、上記ステップ S115 で推定された各ペアの関係を出力部 30 により出力して、発言間関係推定処理ルーチンを終了する。

【0074】

<実施例>

学習データを元に MLN によってモデルを学習し、交差検定によって、上記の実施の形

10

20

30

40

50

態で説明した手法の評価を行った。ここでは、三つの条件で比較した。3つの条件は、ツイッター向け素性有り&制約有り、ツイッター向け素性無し&制約あり、ツイッター向け素性無し&ツイッター向け制約無しである。ツイッター向け素性とは、上記で説明した「in\_reply\_to」、「reply」、「sameuser」、「latestutt」、及び「firstutt」を含む述語である。ツイッター向け制約とは、上記の実施の形態で説明した、発言は関係を持つが一つであるという制約である。なお、ここでは、発言間に関係があるかどうか(ツイート間について、has\_aa\_relationが真か)を推定することを行う。その他の潜在ラベルについても同様の推定は可能である。

【0075】

図5に実験結果を示す。ここで、Recallは再現率を表し、関係性があるツイートペアのうち、いくつを正しく関係があると判定できたかを表す。Precisionは適合率を表し、学習モデルに基づいて関係性があると推定したツイートペアのうち、いくつが実際に正しかったかを表す。F1はRecallとPrecisionの調和平均である。この値が高ければ、正確に漏れなく関係があることを推定できていると言える。F1値に着目すると、ツイッター向けの素性および制約を入れることで、性能が改善されていることが分かる。これにより、本実施の形態で提案する発言間の特徴の有効性が示された。

【0076】

以上説明したように、本実施の形態に係る発言間関係推定装置によれば、発言(ツイート)間の各々について、発言間の内容の関連性を示す述語と、発言間の応答関係を示す述語と、先の発言が同一のユーザ情報が付与された発言のうちの最初の発言であることを示す述語と、先の発言が同一のユーザ情報が付与された発言のうちの最新の発言であることを示す述語とを、発言間の特徴として抽出すると共に、発言と関係がある他の発言が1つである制約を示す論理式を生成し、抽出された発言間の特徴と、生成された制約の論理式と、予め学習されたモデルとに基づいて、マルコフ論理ネットワークにより、発言間の関係を示す潜在述語を推定することにより、発言間の関係を精度よく推定することができる。

【0077】

また、発言対象のユーザ名、発言者のユーザID、システムが付与する発言同士の応答関係を、発言間の関係性の同定に用いることで、関係性の同定精度を向上させる。発言間の関係性の同定精度が改善し、インターネット上の膨大な発言を高精度で構造化できるようになる。発言が高精度に構造化できれば、膨大な情報から効率的に内容を閲覧したり、情報を抽出したりすることが可能となる。

【0078】

なお、本発明は、上述した実施形態に限定されるものではなく、この発明の要旨を逸脱しない範囲内で様々な変形や応用が可能である。

【0079】

例えば、発言以外の文書データ間の関係を推定することに、本発明を適用してもよい。

【0080】

また、潜在述語として、関係があることを示す潜在述語を用いて、発言間に関係があるか否かを推定する場合を例に説明したが、詳細な関係性を示す潜在述語を用いてもよい。この場合には、例えば、「賛成」及び「反対」、並びに、「矛盾」、「類似」、「演繹」及び「帰納」の各々を示す潜在述語を追加して、発言間の関係性を推定するようにしてもよい。

【0081】

また、本願明細書中において、プログラムが予めインストールされている実施形態として説明したが、当該プログラムを、コンピュータ読み取り可能な記録媒体に格納して提供することも可能である。

【符号の説明】

【0082】

10 入力部

10

20

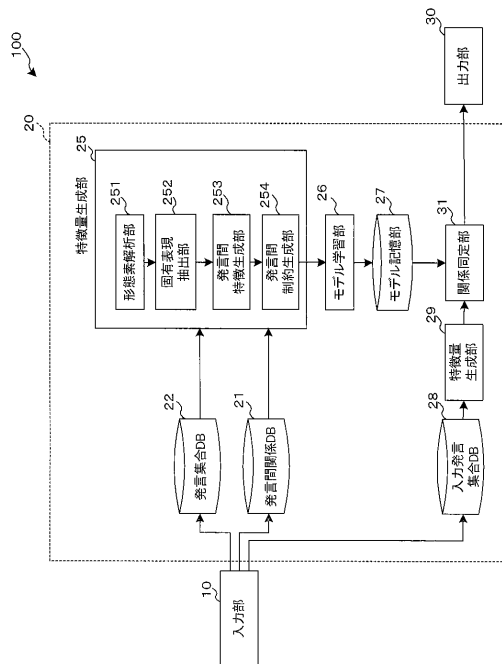
30

40

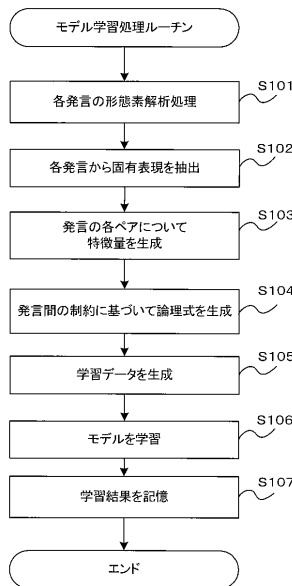
50

- 2 0 演算部
- 2 1 発言集合データベース
- 2 2 発言間関係データベース
- 2 5 特徴量生成部
- 2 6 モデル学習部
- 2 7 モデル記憶部
- 2 8 入力発言集合データベース
- 2 9 特徴量生成部
- 3 0 出力部
- 3 1 関係同定部
- 1 0 0 発言間関係推定装置
- 2 5 1 形態素解析部
- 2 5 2 固有表現抽出部
- 2 5 3 発言間特徴生成部
- 2 5 4 発言間制約生成部

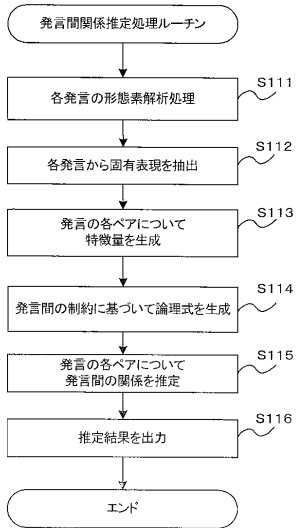
【図 1】



【図 3】



【図4】



【 図 2 】

ID	ユーザー名	返信先	本文	リンク先 ID	横否関係
1. user A _1	user A		個人の書庫をデータ化するサービスなんじゃないかと。でも書籍のデータもこの会社に蓄積しますよね。RT @ user E 安い！本 1 冊 100 円でスキヤンして電子化するサービス <a href="http://bit.ly/cEwJw6">http://bit.ly/cEwJw6</a>		
2. user B _1	user B	1. user A _1	@ user A PDF をお渡しした後にスキヤン漏れや、届いてないなどのトラブル回避のため、短時間は保管します。でも、権利関係の問題が生じるので基本的に蓄積はしません。RT 個人の書庫をデータ化するサービスなんじゃないかと。でも書籍のデータもこの会社に蓄積しますよね。	1. user A _1	反対矛盾
3. user C _1	user C		OCR で文字も読み取っていたら iPhone でも読みやすそう。RT @ user E 安い！本 1 冊 100 円でスキヤンして電子化するサービス <a href="http://www.bookscan.co.jp/">http://www.bookscan.co.jp/</a>		
4. user B _2	user B		読みやすいですね。OCR も 100% の認識率ではないですが、読めます。RT @ user C : OCR で文字も読み取っていたら iPhone でも読みやすそう。RT @ user E 安い！本 1 冊 100 円でスキヤンして	3. user C _1	横成演繹
5. user D _1	user D		これ、商売としてなりたつの？ RT @ user E 安い！本 1 冊 100 円でスキヤンして電子化するサービス <a href="http://bit.ly/cEwJw6">http://bit.ly/cEwJw6</a>		

【 図 5 】

	Recall	Precision	F1
ツイッター向け素性有り&制約有り	0.047	0.043	<b>0.045</b>
ツイッター向け素性なし&制約有り	0.035	0.033	0.034
ツイッター 向け素性無し&ツイッター向け制約無し	1.000	0.020	0.039

## フロントページの続き

- (72)発明者 森田 一  
東京都目黒区大岡山 2 - 1 2 - 1 国立大学法人東京工業大学内
- (72)発明者 奥村 学  
東京都目黒区大岡山 2 - 1 2 - 1 国立大学法人東京工業大学内

審査官 野崎 大進

- (56)参考文献 特開 2 0 1 1 - 1 8 0 9 8 8 ( J P , A )  
特開 2 0 0 7 - 2 8 7 1 3 4 ( J P , A )  
米国特許出願公開第 2 0 0 7 / 0 2 3 3 4 6 5 ( U S , A 1 )  
徳永 泰浩 他, チャット対話における発言間の継続関係と応答関係の同定, 自然言語処理, 日本  
言語処理学会, 2 0 0 5 年 1 月 1 0 日, Vol.12, No.1, pp.79-105.  
吉川 克正 他, 機械学習手法による結合推論を利用した時間的順序関係推定, 第 7 3 回人工知能  
基本問題研究会資料, 日本, 社団法人人工知能学会, 2 0 0 9 年 3 月 6 日, SIG-FPAI-A804-  
12, pp.61-67.  
宮部 泰成 他, 文書横断文間関係の特定, 言語処理学会第 1 2 回年次大会発表論文集, 日本, 言  
語処理学会, 2 0 0 6 年 3 月 1 3 日, pp.496-499.  
小西 卓哉 他, 統計的言語特性を考慮した評判情報のトピックモデリング, 第 3 回データ工学と  
情報マネジメントに関するフォーラム 論文集, 日本, 電子情報通信学会データ工学専門委員会  
, 2 0 1 1 年 7 月 2 7 日, Vol.DEIM2011, No.A8-2, pp.1-8.

- (58)調査した分野(Int.Cl., DB名)  
G 0 6 F 1 7 / 3 0  
J S T P l u s ( J D r e a m I I I )