

(19) 日本国特許庁(JP)

(12) 公開特許公報(A)

(11) 特許出願公開番号

特開2011-138440

(P2011-138440A)

(43) 公開日 平成23年7月14日(2011.7.14)

(51) Int.Cl.	F I	テーマコード (参考)
G06F 17/27 (2006.01)	G06F 17/27 Z	5B075
G06F 17/30 (2006.01)	G06F 17/30 320D	5B091
	G06F 17/30 170A	

審査請求 未請求 請求項の数 13 O L (全 48 頁)

(21) 出願番号	特願2009-299287 (P2009-299287)	(71) 出願人	301022471 独立行政法人情報通信研究機構 東京都小金井市貫井北町4-2-1
(22) 出願日	平成21年12月30日(2009.12.30)	(74) 代理人	100115749 弁理士 谷川 英和
		(72) 発明者	村田 真樹 東京都小金井市貫井北町4-2-1 独立行政法人情報通信研究機構内
		(72) 発明者	小島 正裕 東京都小金井市貫井北町4-2-1 独立行政法人情報通信研究機構内
		(72) 発明者	鳥澤 健太郎 東京都小金井市貫井北町4-2-1 独立行政法人情報通信研究機構内

最終頁に続く

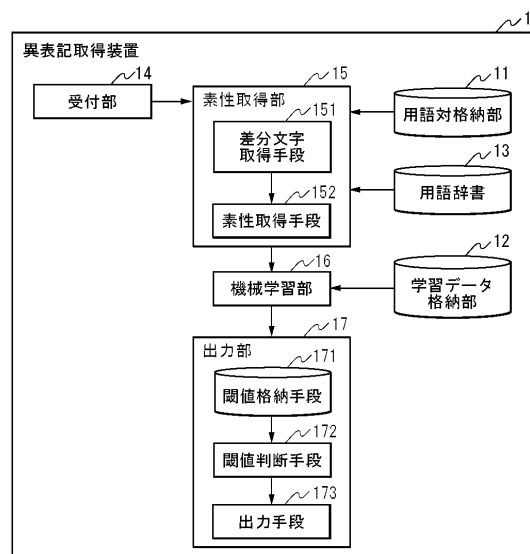
(54) 【発明の名称】 異表記取得装置、異表記取得方法、およびプログラム

(57) 【要約】

【課題】従来、十分な異表記抽出の精度が得られなかった。

【解決手段】用語対の異なる文字である編集箇所の子種に関する素性である字種関連素性、用語辞書を用いて取得された素性である辞書関連素性、用語対を構成する2つの用語の類似度を示す素性である類似度素性のうちの1以上の素性を含む複数の素性と、用語対が異表記の用語対であるかを示す情報である正負情報とを対応付けた学習データを2以上格納し、編集距離が1以上の用語対ごとに、字種関連素性、辞書関連素性、類似度素性のうちの1以上を含む複数の素性を取得する素性取得部と、用語対に対して、2以上の学習データと取得された複数の素性とを用いて、教師あり機械学習法により、各用語対が異表記の用語対であるか否かを判断する機械学習部と、判断結果を出力する出力部とを具備する異表記取得装置により、精度の高い異表記の用語対抽出ができる。

【選択図】 図1



【特許請求の範囲】

【請求項 1】

編集距離が 1 以上の用語対を 1 以上格納し得る用語対格納部と、
用語対の異なる文字である編集箇所の字種に関する素性である字種関連素性、用語辞書を用いて取得された素性である辞書関連素性、前記用語対を構成する 2 つの用語の類似度を示す素性である類似度素性のうちの一以上の素性を含む複数の素性と、前記用語対が異表記の用語対であるかを示す情報である正負情報とを対応付けた学習データを 2 以上格納し得る学習データ格納部と、
前記用語対格納部の用語対ごとに、字種関連素性、辞書関連素性、類似度素性のうちの一以上を含む複数の素性を取得する素性取得部と、
前記各用語対に対して、前記学習データ格納部の 2 以上の学習データと、前記素性取得部が取得した複数の素性とを用いて、教師あり機械学習法により、前記用語対格納部の各用語対が異表記の用語対であるか否かを判断する機械学習部と、
前記機械学習部における判断結果を出力する出力部とを具備する異表記取得装置。

10

【請求項 2】

前記字種関連素性は、
用語対が有する 2 つの用語の編集箇所の字種が異なり、かつ、当該 2 つの用語の編集箇所が数字であるか否かを示す情報であり、
前記素性取得部は、
前記用語対格納部の用語対ごとに、用語対が有する 2 つの用語の編集箇所の字種が異なり、かつ、当該 2 つの用語の編集箇所が同じ値の数字であるという条件に合致するか否かを判断し、当該判断結果を字種関連素性として取得する請求項 1 記載の異表記取得装置。

20

【請求項 3】

前記字種関連素性は、
用語対が有する 2 つの用語の編集箇所の字種がローマ字であり、かつ、当該 2 つの用語の編集箇所が大文字と小文字の違いであるか否かを示す情報であり、
前記素性取得部は、
前記用語対格納部の用語対ごとに、用語対が有する 2 つの用語の編集箇所の字種がローマ字であり、かつ、当該 2 つの用語の編集箇所が大文字と小文字の違いであるという条件に合致するか否かを判断し、当該判断結果を字種関連素性として取得する請求項 1 記載の異表記取得装置。

30

【請求項 4】

用語と、当該用語の代表表記とを有する 1 以上の用語情報を格納し得る用語辞書をさらに具備し、
前記辞書関連素性は、
用語対が有する 2 つの用語の代表表記が同一であるか否かを示す情報であり、
前記素性取得部は、
前記用語対格納部の用語対ごとに、用語対が有する 2 つの用語の代表表記を、前記用語辞書から取得し、当該取得した 2 つの代表表記が同一であるか否かを判断し、当該判断結果を辞書関連素性として取得する請求項 1 記載の異表記取得装置。

40

【請求項 5】

前記辞書関連素性は、
スタッキングアルゴリズムを使用して、前記教師あり機械学習法とは異なる分類方法、または同一の分類方法であるが学習データが異なる分類方法により、用語対が異表記の用語対であるか否かを判断した結果であり、
前記素性取得部は、
前記用語対格納部の用語対ごとに、前記教師あり機械学習法とは異なる分類方法、または同一の分類方法であるが学習データが異なる分類方法により、当該用語対が異表記の用語対であるか否かを判断し、当該判断結果を辞書関連素性として取得する請求項 1 記載の異表記取得装置。

50

【請求項 6】

用語と、当該用語の読みとを有する 1 以上の用語情報を格納し得る用語辞書をさらに具備し、

前記辞書関連素性は、

用語対が有する 2 つの用語の読みが一致するか否かを示す情報であり、

前記素性取得部は、

前記用語対格納部の用語対ごとに、前記用語辞書から前記用語対が有する 2 つの用語の読みを取得し、当該 2 つの用語の読みが一致するか否かを判断し、当該判断結果を辞書関連素性として取得する請求項 1 記載の異表記取得装置。

【請求項 7】

前記機械学習部は、

前記用語対格納部の各用語対が異表記の用語対であるか否かを判断するとともに、異表記の用語対である確度を示すスコアも取得し、

前記出力部は、

前記機械学習部が取得したスコアを出力する請求項 1 から請求項 6 いずれか記載の異表記取得装置。

【請求項 8】

前記出力部は、

スコアの閾値を格納している閾値格納手段と、

前記機械学習部が取得したスコアが前記閾値以上または前記閾値より大きいと判断する閾値判断手段と、

前記閾値判断手段が前記閾値以上または前記閾値より大きいと判断したスコアに対応する用語対を、異表記の用語対であるとの判断結果とし、当該判断結果または異表記の用語対または異表記でない用語対のいずれか 1 以上を出力する出力手段とを具備する請求項 7 記載の異表記取得装置。

【請求項 9】

用語対の異なる文字である編集箇所の子種に関する素性である字種関連素性、用語辞書を用いて取得された素性である辞書関連素性、前記用語対を構成する 2 つの用語の類似度を示す素性である類似度素性のうちの一以上の素性を含む複数の素性と、前記用語対が異表記の用語対であるかを示す情報である正負情報とを対応付けた学習データを 2 以上格納し得る学習データ格納部と、

異表記のパターンを示す第一文字列と第二文字列とを対に有する異表記パターンを 1 以上格納し得る異表記パターン格納部と、

1 以上の用語を受け付ける受付部と、

前記受付部が受け付けた 1 以上の各用語に対して、前記異表記パターン格納部の 1 以上の各異表記パターンを適用し、1 以上の用語を生成し、前記 1 以上の各用語と前記生成した用語とを有する 1 以上の異表記の候補の用語対である異表記候補用語対を生成する用語対生成部と、

前記用語対生成部が生成した 1 以上の異表記候補用語対ごとに、字種関連素性、辞書関連素性、類似度素性のうちの一以上の素性を含む複数の素性を取得する素性取得部と、

前記用語対生成部が生成した各異表記候補用語対に対して、前記学習データ格納部の 2 以上の学習データと、前記素性取得部が取得した複数の素性とを用いて、教師あり機械学習法により、前記用語対格納部の各異表記候補用語対が異表記の用語対であるか否かを判断する機械学習部と、

前記機械学習部における判断結果を出力する出力部とを具備する異表記取得装置。

【請求項 10】

編集距離が 1 の異表記の用語対を 1 以上格納し得る異表記用語対格納部と、

前記異表記用語対格納部に格納されている 1 以上の異表記の用語対の編集箇所を取得する編集箇所取得部と、

前記編集箇所取得部が取得した編集箇所から、第一文字列と第二文字列とを対に有する異

10

20

30

40

50

表記パターンを取得する異表記パターン取得部と、
前記異表記パターン取得部が取得した異表記パターンを、前記異表記パターン格納部に蓄積する異表記パターン蓄積部とをさらに具備する請求項 9 記載の異表記取得装置。

【請求項 1 1】

前記用語対の編集距離は 2 であり、

前記素性取得部は、

前記用語対の 2 つの差分文字の組を、それぞれ取得する差分文字取得手段と、

前記差分文字取得手段が取得した 2 つの差分文字を、独立に対象として、字種関連素性、辞書関連素性、類似度素性のうちの一以上を含む複数の素性を、2 組取得する素性取得手段とを具備し、

前記機械学習部は、

前記素性取得手段が取得した 2 組の複数の素性のうちの組ごとに、当該各組の複数の素性と、前記学習データ格納部の 2 以上の学習データとを用いて、教師あり機械学習法により、前記用語対格納部の各組の複数の素性が異表記の用語対に対応する素性の組であるか否かを判断し、当該 2 つの判断結果を用いて、編集距離が 2 である用語対が異表記の用語対であるか否かを判断する請求項 1 から請求項 1 0 いずれか記載の異表記取得装置。

【請求項 1 2】

記憶媒体に、

編集距離が 1 以上の用語対、および、

用語対の異なる文字である編集箇所の字種に関する素性である字種関連素性、用語辞書を用いて取得された素性である辞書関連素性、前記用語対を構成する 2 つの用語の類似度を示す素性である類似度素性のうちの一以上の素性を含む複数の素性と、前記用語対が異表記の用語対であるかを示す情報である正負情報とを対応付けた学習データを 2 以上格納しており、

素性取得部、機械学習部、および出力部により実現される異表記取得方法であって、前記素性取得部により、前記記憶媒体の用語対ごとに、字種関連素性、辞書関連素性、類似度素性のうちの一以上を含む複数の素性を取得する素性取得ステップと、

前記機械学習部により、前記各用語対に対して、前記記憶媒体の 2 以上の学習データと、前記素性取得ステップで取得された複数の素性とを用いて、教師あり機械学習法により、前記記憶媒体の各用語対が異表記の用語対であるか否かを判断する機械学習ステップと、前記出力部により、前記機械学習ステップにおける判断結果を出力する出力ステップとを具備する異表記取得方法。

【請求項 1 3】

記憶媒体に、

編集距離が 1 以上の用語対、および、

用語対の異なる文字である編集箇所の字種に関する素性である字種関連素性、用語辞書を用いて取得された素性である辞書関連素性、前記用語対を構成する 2 つの用語の類似度を示す素性である類似度素性のうちの一以上の素性を含む複数の素性と、前記用語対が異表記の用語対であるかを示す情報である正負情報とを対応付けた学習データを 2 以上格納しており、

コンピュータを、

前記記憶媒体の用語対ごとに、字種関連素性、辞書関連素性、類似度素性のうちの一以上を含む複数の素性を取得する素性取得部と、

前記各用語対に対して、前記記憶媒体の 2 以上の学習データと、前記素性取得部が取得した複数の素性とを用いて、教師あり機械学習法により、前記記憶媒体の各用語対が異表記の用語対であるか否かを判断する機械学習部と、

前記機械学習部における判断結果を出力する出力部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】

【0 0 0 1】

10

20

30

40

50

本発明は、異表記の用語対を取り出す異表記取得装置等に関するものである。

【背景技術】

【0002】

従来の異表記を取り出す技術としては、荒牧らの研究がある（非特許文献1参照）。この研究は、医療分野の専門用語の異表記の取り出しを行うものであった。なお、異表記とは、例えば「スパゲティ」に対して、「スパゲッティ」など、同義語であるが異なる表現の用語を言う。また、異表記の二つの用語を異表記対という。

【0003】

異表記対の第一の考え方は、以下である。例えば、用語対の例1（問い合わせメール、問い合わせメール）、例2（学園闘争、学園紛争）について、例1は異表記対とし、例2は、異表記対ではなく日本語同義語対とする。つまり、第一の考え方において、同一語の異形なら異表記対とし、同一語でなければ、例え意味が同等でも異表記対としない。闘争と紛争は、ほぼ同等の意味を有するが、同一の語でないので、例2は異表記対とはしない。一方、例1の「問い合わせ」「問い合わせ」は、表記は異なるが同一の語の異形と判断できるので、異表記対とする。

10

【0004】

また、異表記対の第二の考え方は、同義語も異表記とする考え方である。第二の考え方では、上記の例1だけではなく、例2（学園闘争、学園紛争）も異表記対となる。

【0005】

さらに、異表記、異表記対の考え方は、上記の考え方と類似する考え方でも良く、異表記、異表記対は広く解するものとする。

20

【0006】

また、従来技術として、機械学習法についての技術がある（例えば、非特許文献2～非特許文献4参照）

【先行技術文献】

【非特許文献】

【0007】

【非特許文献1】Eiji Aramaki, Takeshi Imai, Kengo Miyo, Kazuhiko Ohe: Orthographic Disambiguation Incorporating Transliterated Probability, International Joint Conference on Natural Language Processing (IJCNLP2008), pp.48-55, 2008.

30

【非特許文献2】村田真樹, 機械学習に基づく言語処理, 龍谷大学理工学部. 招待講演. 2004. <http://www2.nict.go.jp/x/x161/member/murata/ps/kougi-ml-siryuu-new2.pdf>

【非特許文献3】サポートベクトルマシンを用いたテンス・アスペクト・モダリティの日英翻訳, 村田真樹, 馬青, 内元清貴, 井佐原均, 電子情報通信学会言語理解とコミュニケーション研究会 NLC2000-78, 2001年.

【非特許文献4】SENSEVAL2J辞書タスクでのCRLの取り組み, 村田真樹, 内山将夫, 内元清貴, 馬青, 井佐原均, 電子情報通信学会言語理解とコミュニケーション研究会 NLC2001-40, 2001年.

【発明の概要】

【発明が解決しようとする課題】

40

【0008】

しかしながら、従来技術においては、日本語の一般の異表記を扱うものではなく、また、従来技術を日本語の一般の異表記に適用したのでは、十分な異表記抽出の精度が得られなかった。

【課題を解決するための手段】

【0009】

本第一の発明の異表記取得装置は、編集距離が1以上の用語対を1以上格納し得る用語対格納部と、用語対の異なる文字である編集箇所の子種に関する素性である字種関連素性、用語辞書を用いて取得された素性である辞書関連素性、用語対を構成する2つの用語の類似度を示す素性である類似度素性のうちの1以上の素性を含む複数の素性と、用語対が

50

異表記の用語対であるかを示す情報である正負情報とを対応付けた学習データを2以上格納し得る学習データ格納部と、用語対格納部の用語対ごとに、字種関連素性、辞書関連素性、類似度素性のうちの一以上を含む複数の素性を取得する素性取得部と、各用語対に対して、学習データ格納部の2以上の学習データと、素性取得部が取得した複数の素性を用いて、教師あり機械学習法により、用語対格納部の各用語対が異表記の用語対であるか否かを判断する機械学習部と、機械学習部における判断結果を出力する出力部とを具備する異表記取得装置である。

【0010】

かかる構成により、用語対の分野を問わず、精度の高い異表記の用語対の抽出が可能となる。

10

【0011】

また、本第二の発明の異表記取得装置は、第一の発明に対して、字種関連素性は、用語対が有する2つの用語の編集箇所の字種が異なり、かつ、2つの用語の編集箇所が数字であるか否かを示す情報であり、素性取得部は、用語対格納部の用語対ごとに、用語対が有する2つの用語の編集箇所の字種が異なり、かつ、2つの用語の編集箇所が同じ値の数字であるという条件に合致するか否かを判断し、判断結果を字種関連素性として取得する異表記取得装置である。

【0012】

かかる構成により、用語対の分野を問わず、精度の高い異表記の用語対の抽出が可能となる。

20

【0013】

また、本第三の発明の異表記取得装置は、第一の発明に対して、字種関連素性は、用語対が有する2つの用語の字種がローマ字であり、かつ、2つの用語の編集箇所が大文字と小文字の違いであるか否かを示す情報であり、素性取得部は、用語対格納部の用語対ごとに、用語対が有する2つの用語の編集箇所の字種がローマ字であり、かつ、2つの用語の編集箇所が大文字と小文字の違いであるという条件に合致するか否かを判断し、判断結果を字種関連素性として取得する異表記取得装置である。

【0014】

かかる構成により、用語対の分野を問わず、精度の高い異表記の用語対の抽出が可能となる。

30

【0015】

また、本第四の発明の異表記取得装置は、第一の発明に対して、用語と、用語の代表表記とを有する1以上の用語情報を格納し得る用語辞書をさらに具備し、辞書関連素性は、用語対が有する2つの用語の代表表記が同一であるか否かを示す情報であり、素性取得部は、用語対格納部の用語対ごとに、用語対が有する2つの用語の代表表記を、用語辞書から取得し、取得した2つの代表表記が同一であるか否かを判断し、判断結果を辞書関連素性として取得する異表記取得装置である。

【0016】

かかる構成により、用語対の分野を問わず、精度の高い異表記の用語対の抽出が可能となる。

40

【0017】

また、本第五の発明の異表記取得装置は、第一の発明に対して、辞書関連素性は、スタッキングアルゴリズムを使用して、教師あり機械学習法とは異なる分類方法、または同一の分類方法であるが学習データが異なる分類方法により、用語対が異表記の用語対であるか否かを判断した結果であり、素性取得部は、用語対格納部の用語対ごとに、教師あり機械学習法とは異なる分類方法、または同一の分類方法であるが学習データが異なる分類方法により、用語対が異表記の用語対であるか否かを判断し、判断結果を辞書関連素性として取得する異表記取得装置である。

【0018】

かかる構成により、用語対の分野を問わず、精度の高い異表記の用語対の抽出が可能と

50

なる。

【0019】

また、本第六の発明の異表記取得装置は、第一の発明に対して、用語と、用語の読みとを有する1以上の用語情報を格納し得る用語辞書をさらに具備し、辞書関連素性は、用語対が有する2つの用語の読みが一致するか否かを示す情報であり、素性取得部は、用語対格納部の用語対ごとに、用語辞書から用語対が有する2つの用語の読みを取得し、2つの用語の読みが一致するか否かを判断し、判断結果を辞書関連素性として取得する異表記取得装置である。

【0020】

かかる構成により、用語対の分野を問わず、精度の高い異表記の用語対の抽出が可能となる。

10

【0021】

また、本第七の発明の異表記取得装置は、第一から第六いずれかの発明に対して、機械学習部は、用語対格納部の各用語対が異表記の用語対であるか否かを判断するとともに、異表記の用語対である確度を示すスコアも取得し、出力部は、機械学習部が取得したスコアを出力する異表記取得装置である。

【0022】

かかる構成により、用語対の分野を問わず、さらに精度の高い異表記の用語対の抽出が可能となる。

【0023】

また、本第八の発明の異表記取得装置は、第七の発明に対して、出力部は、スコアの閾値を格納している閾値格納手段と、機械学習部が取得したスコアが閾値以上または閾値より大きいか否かを判断する閾値判断手段と、閾値判断手段が閾値以上または閾値より大きいと判断したスコアに対応する用語対を、異表記の用語対であるとの判断結果とし、判断結果または異表記の用語対または異表記でない用語対のいずれか1以上を出力する出力手段とを具備する異表記取得装置である。

20

【0024】

かかる構成により、用語対の分野を問わず、さらに精度の高い異表記の用語対の抽出が可能となる。

【0025】

また、本第九の発明の異表記取得装置は、用語対の異なる文字である編集箇所の字種に関する素性である字種関連素性、用語辞書を用いて取得された素性である辞書関連素性、用語対を構成する2つの用語の類似度を示す素性である類似度素性のうちの一以上の素性を含む複数の素性と、用語対が異表記の用語対であるかを示す情報である正負情報とを対応付けた学習データを2以上格納し得る学習データ格納部と、異表記のパターンを示す第一文字列と第二文字列とを対に有する異表記パターンを1以上格納し得る異表記パターン格納部と、1以上の用語を受け付ける受付部と、受付部が受け付けた1以上の各用語に対して、異表記パターン格納部の1以上の各異表記パターンを適用し、1以上の用語を生成し、1以上の各用語と生成した用語とを有する1以上の異表記の候補の用語対である異表記候補用語対を生成する用語対生成部と、用語対生成部が生成した1以上の異表記候補用語対ごとに、字種関連素性、辞書関連素性、類似度素性のうちの一以上の素性を含む複数の素性を取得する素性取得部と、用語対生成部が生成した各異表記候補用語対に対して、学習データ格納部の2以上の学習データと、素性取得部が取得した複数の素性とを用いて、教師あり機械学習法により、用語対格納部の各異表記候補用語対が異表記の用語対であるか否かを判断する機械学習部と、機械学習部における判断結果を出力する出力部とを具備する異表記取得装置である。

30

40

【0026】

かかる構成により、異表記の用語対の候補を自動生成できる。

【0027】

また、本第十の発明の異表記取得装置は、第九の発明に対して、編集距離が1の異表記

50

の用語対を1以上格納し得る異表記用語対格納部と、異表記用語対格納部に格納されている1以上の異表記の用語対の編集箇所を取得する編集箇所取得部と、編集箇所取得部が取得した編集箇所から、第一文字列と第二文字列とを対に有する異表記パターンを取得する異表記パターン取得部と、異表記パターン取得部が取得した異表記パターンを、異表記パターン格納部に蓄積する異表記パターン蓄積部とをさらに具備する異表記取得装置である。

【0028】

かかる構成により、異表記の用語対の候補を自動生成するための異表記パターンを自動的に取得できる。

【0029】

また、本第十一の発明の異表記取得装置は、第一から第十いずれかの発明に対して、用語対の編集距離は2であり、素性取得部は、用語対の2つの差分文字の組を、それぞれ取得する差分文字取得手段と、差分文字取得手段が取得した2つの差分文字を、独立に対象として、字種関連素性、辞書関連素性、類似度素性のうちの一以上を含む複数の素性を、2組取得する素性取得手段とを具備し、機械学習部は、素性取得手段が取得した2組の複数の素性のうちの組ごとに、各組の複数の素性と、学習データ格納部の2以上の学習データとを用いて、教師あり機械学習法により、用語対格納部の各組の複数の素性が異表記の用語対に対応する素性の組であるか否かを判断し、2つの判断結果を用いて、編集距離が2である用語対が異表記の用語対であるか否かを判断する異表記取得装置である。

【0030】

かかる構成により、編集距離が2の用語対でも、精度高く、異表記の用語対であるか否かを判断できる。

【発明の効果】

【0031】

本発明による異表記取得装置によれば、用語対の分野を問わず、精度の高い異表記の用語対の抽出が可能となる。

【図面の簡単な説明】

【0032】

【図1】本発明の実施の形態1における異表記取得装置のブロック図

【図2】同異表記取得装置の動作について説明するフローチャート

【図3】同素性取得処理の動作について説明するフローチャート

【図4】同用語辞書の例を示す図

【図5】同サポートベクトルマシン法のマージン最大化の概念を示す図

【図6】同実験で用いた編集距離が1の日本語用語対の中に、多数決により日本語異表記対であるか日本語異表記対でないかを判定した内訳を示す図

【図7】同Landisらによる一致度の評価方法を示す図

【図8】同クローズドデータとオープンデータに対して、ベースラインの手法を適用した結果を示す図

【図9】同ブートストラップ法を用いて素性が有効であるかどうかの検討をした結果を示す図

【図10】同素性の例を示す図

【図11】同提案手法を用い、大規模類似語リストから編集距離が1の日本語異表記対と分類された用語対が、種々の辞書にどの程度の割合で含まれているかの検討結果を示す図

【図12】同SVMの分類精度を示す図

【図13】同種々の辞書と用語対DBにおいて、編集距離が1の日本語異表記対であると分類された日本語用語対と、分類されなかった日本語用語対をそれぞれランダムに、5組ずつ取り出した結果を示す図

【図14】同閾値の評価基準を示す図

【図15】同再現率と適合率の比率を示す図

【図16】同ベースライン手法を用いた場合の実験結果を示す図

10

20

30

40

50

- 【図 17】同ベースライン手法を用いた場合の実験結果を示す図
- 【図 18】同ベースライン手法を用いた場合の実験結果を示す図
- 【図 19】本発明の実施の形態 2 における異表記取得装置のブロック図
- 【図 20】同異表記取得装置の動作について説明するフローチャート
- 【図 21】同異表記パターンの例を示す図
- 【図 22】上記実施の形態におけるコンピュータシステムの概観図
- 【図 23】同コンピュータシステムのブロック図
- 【発明を実施するための形態】

【0033】

以下、異表記取得装置等の実施形態について図面を参照して説明する。なお、実施の形態において同じ符号を付した構成要素は同様の動作を行うので、再度の説明を省略する場合がある。

10

【0034】

(実施の形態 1)

【0035】

本実施の形態において、編集距離が 1 または 2 以上の用語対から、異なる文字の字種に関する素性、用語辞書を用いて取得された素性、2 つの用語の類似度のうちの 1 以上の素性を含む複数の素性を取り出し、当該複数の素性を用いて、用語対が異表記の用語対であるか否かを、教師あり機械学習法により判断する異表記取得装置について説明する。

【0036】

20

図 1 は、本実施の形態における異表記取得装置 1 のブロック図である。

異表記取得装置 1 は、用語対格納部 11、学習データ格納部 12、用語辞書 13、受付部 14、素性取得部 15、機械学習部 16、出力部 17 を備える。素性取得部 15 は、差分文字取得手段 151、素性取得手段 152 を備える。出力部 17 は、閾値格納手段 171、閾値判断手段 172、出力手段 173 を備える。

【0037】

用語対格納部 11 は、編集距離が 1 または 2 以上の用語対を 1 以上格納し得る。編集距離とは、異なる文字の数である。また、用語対とは、2 つの用語である。編集距離が 2 の用語対は、異なる文字数が 2 つの用語である。なお、用語とは、通常、名詞や名詞句であるが、形容詞等の他の品詞の用語でも良い。

30

【0038】

学習データ格納部 12 は、2 以上の学習データを格納し得る。学習データは、用語対の複数の素性と正負情報とを有する。学習データは、用語対を有しても良い。用語対の複数の素性は、ここでは、字種関連素性、辞書関連素性、類似度素性のうちの 1 以上の素性を含む、とする。なお、素性とは、異表記取得装置 1 が学習する際に手掛かりとする情報のことである。

【0039】

字種関連素性とは、用語対の異なる文字である編集箇所の字種に関する素性である。字種関連素性は、例えば、用語対が有する 2 つの用語の編集箇所の字種が異なり、かつ、2 つの用語の編集箇所が同じ値の数字であるか否かを示す情報である。また、字種関連素性は、例えば、用語対が有する 2 つの用語の文字数が同数であり、かつ、2 つの用語の編集箇所の字種が異なり、かつ、2 つの用語の編集箇所が同じ値の数字であるか否かを示す情報である。また、字種関連素性は、例えば、用語対が有する 2 つの用語の文字数が同数であり、かつ、2 つの用語の編集箇所の字種がローマ字であり、かつ、2 つの用語の編集箇所が大文字と小文字の違いであるか否かを示す情報である。また、字種関連素性は、例えば、用語対が有する 2 つの用語の文字数が同数であり、かつ、2 つの用語の編集箇所の字種がローマ字であり、かつ、2 つの用語の編集箇所が大文字と小文字の違いであるか否かを示す情報である。

40

【0040】

また、辞書関連素性とは、用語辞書 13 を用いて取得された素性である。辞書関連素性は、例えば、スタッキングアルゴリズムを使用して、機械学習部 16 が利用する教師あり

50

機械学習法とは異なる分類方法、または同一の分類方法であるが学習データが異なる分類方法により、用語対が異表記の用語対であるか否かを判断した結果である。ここで、「学習データが異なる」とは、学習データの元になる用語対の集合が異なる場合、学習データが有する素性が異なる場合などがある。また、辞書関連素性は、例えば、用語対が有する2つの用語の代表表記が同一であるか否かを示す情報である。また、辞書関連素性は、例えば、用語対が有する2つの用語の読みが一致するか否かを示す情報である。また、辞書関連素性は、例えば、用語対が有する2つの用語の文字数が同数であり、かつ、2つの用語の読みが一致するか否かを示す情報である。なお、分類方法とは、異表記の用語対であるか否かの分類の方法である。また、教師あり機械学習法とは異なる分類方法とは、分類のやり方、アルゴリズムが教師あり機械学習法とは異なることである。

10

【0041】

また、類似度素性とは、用語対を構成する2つの用語の類似度を示す素性である。2つの用語の類似度は、それらの用語がWeb上でよく似た文脈に出現するかどうかの情報を利用して求める。なお、用語の類似度を取得する技術は、「風間淳一, De Saeger, Stijn, 鳥澤健太郎, 村田真樹「係り受けの確率的クラスタリングを用いた大規模類似語リストの作成」言語処理学会第15回年次大会(NLP2009)」等に記載されている。つまり、2つの用語の類似度の取得方法は公知技術である。2つの用語の類似度の算出方法は問わない。

【0042】

また、正負情報とは、用語対が異表記の用語対であるか否かを示す情報である。正負情報は、異表記の用語対であれば正例(例えば「1」)、異表記の用語対でなければ負例(例えば「0」)である。

20

【0043】

また、用語辞書とは、異表記の用語の情報を含む情報群である。用語辞書の例やデータ構造の例については後述する。

【0044】

また、他の素性として、編集箇所の文字または編集箇所の文字の周辺の文字の情報である編集箇所文字素性がある。

【0045】

用語辞書13は、1以上の用語情報を格納し得る。用語辞書13は、例えば、異表記の用語の情報を含む情報群である。用語辞書13は、異表記の2つの用語が、陽に対応付けられている必要はない。用語情報は、例えば、用語と用語の代表表記とを有する。用語情報は、例えば、用語と、用語の読みとを有する。

30

【0046】

受付部14は、ユーザからの入力を受け付ける。この入力とは、例えば、異表記取得装置1を動作させるための動作指示である。受付部14は、異表記であるか否かを判断する対象の用語対を受け付けても良い。動作指示などの入力手段は、テンキーやキーボードやマウスやメニュー画面によるもの等、何でも良い。受付部14は、テンキーやキーボード等の入力手段のデバイスドライバや、メニュー画面の制御ソフトウェア等で実現される。

【0047】

素性取得部15は、用語対格納部11の用語対ごとに、字種関連素性、辞書関連素性、類似度素性のうちの一以上を含む複数の素性を取得する。複数の素性とは、例えば、後述する68の素性である。

40

【0048】

素性取得部15は、用語対格納部11の用語対ごとに、用語対が有する2つの用語の編集箇所の字種が異なり、かつ、2つの用語の編集箇所が同じ値の数字であるという条件に合致するか否かを判断し、判断結果を字種関連素性として取得する。なお、素性取得部15は、例えば、用語対「三者会談」「三者会談」に対して、編集箇所が「3」「三」であるので、上記条件に合致する、と判断する。また、素性取得部15は、例えば、用語対「125件」「百二十五件」に対して、編集箇所が「125」「百二十五」であるので、用

50

語の文字数は同数ではないが、上記条件に合致する、と判断する。

【0049】

素性取得部15は、例えば、用語対格納部11の用語対ごとに、用語対が有する2つの用語の編集箇所の字種がローマ字であり、かつ、2つの用語の編集箇所が大文字と小文字の違いであるという条件に合致するか否かを判断し、判断結果を字種関連素性として取得する。なお、素性取得部15は、例えば、編集箇所が「A」と「a」の用語対に対して合致すると判断し、編集箇所が「A」と「b」の用語対に対して合致しないと判断する。

【0050】

素性取得部15は、例えば、用語対格納部11の用語対ごとに、機械学習部16が利用する教師あり機械学習法とは異なる分類方法、または同一の分類方法であるが学習データが異なる分類方法により、用語対が異表記の用語対であるか否かを判断し、判断結果を辞書関連素性として取得する。かかる辞書関連素性を利用する機械学習法を、スタッキングアルゴリズムによる方法という。機械学習部16が利用する教師あり機械学習法とは異なる分類方法とは、上記の機械学習法がSVMである場合、SVMとは異なる決定木などの機械学習法、後述するルールに基づく分類方法等である。

10

【0051】

スタッキングアルゴリズムは、詳細には、例えば、以下の手順による分類方法である。まず、JUMAN辞書を使って教師データを作成する。つまり、JUMAN辞書の単語の集合から、編集距離が1文字の単語対を取り出す。ここで、編集距離が1文字の単語対は、904612組、取り出せる。そのうち、代表表記が等しい単語対(25934組)を取り出す。次に、JUMAN辞書で、代表表記が等しい単語対を正例、そうでないものを負例とする。以上により、教師データを作成する。

20

【0052】

次に、その教師データを学習データとした機械学習を行う。なお、教師データは、上述した教師データに限らず、他の教師データを用いてもよい。また、機械学習の際に利用する素性は、本発明の全素性(S1からS68の素性)のうち、S54の素性を取り除いた素性を利用する。なお、機械学習の際に利用する素性は、他の素性を用いてもよい。

【0053】

そして、実際に、S54の素性を付与したいデータを、上記学習結果を利用して、分類する。分類結果において正例となったか、負例となったかの情報をS54の素性として、そのデータに付与する。

30

【0054】

そして、S54の素性が付与された学習データ(68の素性を有する)を用いて、問題となる用語対に対して、機械学習を行うことで、問題となる用語対が異表記対であるか否かを判断していく。

【0055】

スタッキングアルゴリズムによる方法では、JUMAN辞書で、代表表記が一致するか否かについて学習した結果を素性として付与できるので、実際にJUMAN辞書に記載されていない用語対に対しても、JUMAN辞書で、代表表記が一致するとされる傾向のある用語対か否かの情報を付与できることとなる。

40

【0056】

素性取得部15は、例えば、用語対格納部11の用語対ごとに、用語対が有する2つの用語の代表表記を、用語辞書13から取得し、取得した2つの代表表記が同一であるか否かを判断し、判断結果を辞書関連素性として取得する。

【0057】

素性取得部15は、例えば、用語対格納部11の用語対ごとに、用語辞書13から2つの用語の読みを取得し、2つの用語の読みが一致するか否かを判断し、判断結果を辞書関連素性として取得する。また、素性取得部15は、例えば、用語対格納部11の用語対ごとに、用語対が有する2つの用語の文字数が同数であり、かつ、用語辞書13から2つの用語の読みを取得し、2つの用語の読みが一致するか否かを判断し、判断結果を辞書関連

50

素性として取得しても良い。

【0058】

なお、上述した判断結果とは、例えば、上記条件に合致する場合の判断結果は「1」、その他の場合の判断結果は「0」などである。

【0059】

差分文字取得手段151は、編集距離が2つの用語対について、2つの差分文字の組を、それぞれ取得する。例えば、編集距離が2つの用語対が、(1)「できる」「出来る」(2)「理解できる」「できる」(3)「IX(ローマ数字の9)」「9」である場合を考える。(1)は両方の用語対が同じ文字数である場合である。(2)はどちらか一方の用語の文字数がもう片方の用語の文字数より2つ多いまたは、少ない場合である。(3)はどちらか一方の用語の文字数がもう片方の用語の文字数より1つ多いまたは、少ない場合である。(1)の場合、差分文字取得手段151は、「できる」および「出来る」の用語に対して、前方から後方に1, 2, 3・・・と文字に番号をつけ、それぞれの用語で同じ文字番号を持ち、違う文字である「で」「出」と「き」「来」が差分文字であるとして、「で」「出」と「き」「来」の2組の差分文字の組を取得する。(2)の場合、差分文字取得手段151は、「理」「」と「解」「」(「」はNULLである)の2組の差分文字の組を取得する。(3)の場合、差分文字取得手段151は、「I」「9」と「X」「」の2組の差分文字、または「I」「」と「X」「9」の2組の差分文字を取得する。

10

【0060】

また、差分文字取得手段151は、編集距離が1の用語対について、差分文字の組を1組取得する。例えば、編集距離が1の用語対が、(1)「ご苦労」「御苦労」(2)「Firefox」「FireFox」(3)「肝炎ウイルス」「肝炎ウイルス」(4)「文学史上」「文学史」(5)「咲き分け」「咲分け」(6)「クロゼット」「クローゼット」(7)「大人・子供」「大人子供」(8)「第1位」「第一位」である場合、差分文字取得手段151は、それぞれ(1)「ご」「御」(2)「f」「F」(3)「イ」「イ」(4)「上」「」(5)「き」「」(6)「」「ー」(7)「・」「」(8)「1」「ー」を取得する。

20

【0061】

さらに、差分文字取得手段151は、編集距離が3以上の用語対について、3組以上の差分文字の組を取得する。例えば、編集距離が4の用語対が、「1025位」「千二十五位」である場合、差分文字取得手段151は、「1」「千」、「0」「二」、「2」「十」「5」「五」という4組の差分文字を取得する。ここで、差分文字とは、2つの用語の異なる文字である。

30

【0062】

素性取得手段152は、差分文字取得手段151が取得した2つの差分文字を、独立に対象として、字種関連素性、辞書関連素性、類似度素性のうちの一以上を含む複数の素性を、2組取得する。例えば、編集距離が2つの用語対が、(1)「できる」「出来る」(2)「理解できる」「できる」(3)「IX」「9」である場合を考える。(1)の用語対について、素性取得手段152は、「で」「出」と「き」「来」の2組の差分文字の組のそれぞれを対象に素性の抽出を行い、それぞれ差分文字から抽出した素性は、別のものと考え、2種類のテストデータを作成する。素性取得手段152は、例えば、用語対が有する2つの用語の編集箇所の字種が異なり、かつ、2つの用語の編集箇所が同じ値の数字であるか否かを示す字種関連素性について、「で」「出」の編集箇所が同じ値の数字でない」と判断し、当該字種関連素性「0」を取得する。また、素性取得手段152は、例えば、用語辞書13から2つの用語の読みを取得し、2つの用語の読みが一致するか否かを示す辞書関連素性「1」を取得する。素性取得手段152は、用語辞書13から「出」の読み「で」を取得し、「で」と「出」の読みが一致すると判断する。また、素性取得手段152は、例えば、差分文字「で」「出」に対して、差分文字(編集箇所)の前後の文字という素性について、前の文字の素性「」(なし)、後の文字の素性「き」と「来」を取得する。また素性取得手段152は、例えば、差分文字「き」「来」に対して、差分文字の

40

50

前後の文字という素性について、前の文字の素性「出」と「で」、後の文字の素性「る」を取得する。かかる処理により、別の差分文字も素性に含めることとなる。

【0063】

また、(2)の用語対について、素性取得手段152は、(1)と同様に、「理」「」と「解」「」の2組の差分文字の組のそれぞれを対象に素性の抽出を行い、それぞれ差分文字から抽出した素性は、別のものと考え、2種類のテストデータを作成する。

【0064】

さらに、(3)の用語対について、素性取得手段152は、(1)(2)と同様に、例えば、「I」「9」と「X」「」の2組の差分文字の組のそれぞれを対象に素性の抽出を行い、それぞれ差分文字から抽出した素性は、別のものと考え、2種類のテストデータを作成する。

10

【0065】

また、素性取得手段152は、差分文字取得手段151が取得した1組以上の差分文字を用いて、字種関連素性、辞書関連素性、類似度素性のうちの一以上を含む複数の素性を取得する。なお、字種関連素性、辞書関連素性、類似度素性などの素性を取得する具体的な方法は後述する。

【0066】

機械学習部16は、各用語対に対して、学習データ格納部12の2以上の学習データと、素性取得部15が取得した複数の素性を用いて、教師あり機械学習法により、用語対格納部11の各用語対が異表記の用語対であるか否かを判断する。

20

【0067】

機械学習部16は、用語対格納部11の各用語対が異表記の用語対であるか否かを判断するとともに、異表記の用語対である確度を示すスコアも取得しても良い。

【0068】

機械学習部16は、素性取得手段152が取得した2組の複数の素性のうちの組ごとに、各組の複数の素性と、学習データ格納部12の2以上の学習データを用いて、教師あり機械学習法により、用語対格納部11の各組の複数の素性が異表記の用語対に対応する素性の組であるか否かを判断し、2つの判断結果を用いて、編集距離が2である用語対が異表記の用語対であるか否かを判断する。

【0069】

教師あり機械学習法のアルゴリズムは問わない。教師あり機械学習法とは、例えば、サポートベクターマシン(SVM)などである。SVMは、「<http://chasen.org/~taku/software/TinySVM/>」「<http://ja.wikipedia.org/wiki/%E3%82%B5%E3%83%9D%E3%83%BC%E3%83%88%E3%83%99%E3%82%AF%E3%82%BF%E3%83%BC%E3%83%9E%E3%82%B7%E3%83%B3>」(平成21年12月12日検索)などに記載されている。なお、教師あり機械学習法の詳細は、後述する。

30

【0070】

また、上記の、2つの判断結果を用いてとは、2つとも正例とされた場合に異表記の用語対としても良いし、2つとも負例とされた場合に異表記の用語対ではないとしても良いし、2つのスコアのうちのスコアが0に近い方のスコアを採用して、採用したスコアが正の場合は正例(異表記の用語)、負の場合は負例(異表記の用語でない)と判断しても良いし、スコアの絶対値が大きい方のスコアを採用して、採用したスコアが正の場合は正例(異表記の用語)、負の場合は負例(異表記の用語でない)と判断しても良い。また、2つのスコアのうちの、小さい方のスコアを取得し、当該小さい方のスコアが正の場合は正例、負の場合は負例と判断しても良い。つまり、2つの判断結果の使い方は問わない。なお、上記の(2)の場合(どちらか一方の用語の文字数がもう片方の用語の文字数より2つ多いまたは、少ない場合)、大規模類似語リストの中から、約1万5千のタグ付けを行った結果、このパターンの2文字差分データには、異表記対であると判定する用語対はなかった。

40

【0071】

50

さらに、2組の差分文字の組（例えば、「I」「9」と「X」「」、または「I」「」と「X」「9」）、つまり2つの問題（問題1、問題2）ができる場合、それぞれの差分文字を対象に素性の抽出を行い、それぞれ差分文字から抽出した素性は、別のものと考え、4種類のテストデータを作成する。そして、2つの問題ごとに、算出したスコアが0に近い方を取得し、問題ごとのスコアのうちの、絶対値が高いスコアを当該問題のスコアとし、スコアが正の場合は正例、負の場合は負例と判断しても良い。例えば、編集距離が2の用語対が(3)「IX」「9」である場合、問題「I」「9」と「X」「」、および「I」「」と「X」ができる。そして、機械学習部16は、「I」「9」と「X」「」のスコアの小さい方を取得し、また、「I」「」と「X」「9」のスコアの小さい方を取得し、2つの取得されたスコアのうちの、値が大きい方を「IX」「9」の用語対におけるスコアとする。そして、機械学習部16は、当該スコアが正の場合は正例、負の場合は負例と判断しても良い。なお、例えば、機械学習部16は、「I」「9」と「X」「」のスコアが0に近い方を取得し、また、「I」「」と「X」「9」のスコアが0に近い方を取得し、2つの取得されたスコアのうちの、絶対値が大きい方を「IX」「9」の用語対におけるスコアとしても良い。つまり、4種類のテストデータの判断結果を如何に用いてスコアを算出するかは問わない。

10

20

30

40

50

【0072】

出力部17は、機械学習部16における判断結果を出力する。また、出力部17は、機械学習部16が取得したスコアを出力しても良い。判断結果とは、各用語対が異表記の用語対であるか否かを示す情報、または異表記の1以上の用語対、または異表記でない1以上の用語対などである。また、出力部17は、判断結果とスコアの両方を出力しても良いし、一方を出力しても良い。

【0073】

また、出力とは、ディスプレイへの表示、プロジェクターを用いた投影、プリンタへの印字、音出力、外部の装置への送信、記録媒体への蓄積、他の処理装置や他のプログラムなどへの処理結果の引渡しなどを含む概念である。

【0074】

閾値格納手段171は、スコアの閾値を格納している。

【0075】

閾値判断手段172は、機械学習部16が取得したスコアが閾値以上または閾値より大きいかなんかを判断する。

【0076】

出力手段173は、閾値判断手段172が閾値以上または閾値より大きいと判断したスコアに対応する用語対を、異表記の用語対であるとの判断結果とし、判断結果または異表記の用語対または異表記でない用語対のいずれか1以上を出力する。

【0077】

用語対格納部11、学習データ格納部12、用語辞書13、および閾値格納手段171は、不揮発性の記録媒体が好適であるが、揮発性の記録媒体でも実現可能である。

【0078】

用語対格納部11、学習データ格納部12、および用語辞書13に格納されている情報が記憶される過程は問わない。

【0079】

素性取得部15、機械学習部16、閾値判断手段172は、通常、MPUやメモリ等から実現され得る。素性取得部15等の処理手順は、通常、ソフトウェアで実現され、当該ソフトウェアはROM等の記録媒体に記録されている。但し、ハードウェア（専用回路）で実現しても良い。

【0080】

出力部17は、ディスプレイやスピーカ等出力デバイスを含むと考えても含まないと考えても良い。出力部17は、出力デバイスのドライバソフトまたは、出力デバイスのドライバソフトと出力デバイス等で実現され得る。

【 0 0 8 1 】

次に、異表記取得装置 1 の動作について図 2 のフローチャートを用いて説明する。ここでは、異表記取得装置 1 は、編集距離が 1 の用語対に対して、異表記の用語対であるか否かを判断することとする。

【 0 0 8 2 】

(ステップ S 2 0 1) 受付部 1 4 は、動作開始の指示を受け付けたか否かを判断する。指示を受け付ければステップ S 2 0 2 に行き、受け付けなければステップ S 2 0 1 に戻る。

【 0 0 8 3 】

(ステップ S 2 0 2) 素性取得部 1 5 は、カウンタ i に 1 を代入する。

10

【 0 0 8 4 】

(ステップ S 2 0 3) 素性取得部 1 5 は、 i 番目の用語対が用語対格納部 1 1 に存在するか否かを判断する。 i 番目の用語対が存在すればステップ S 2 0 4 に行き、存在しなければ処理を終了する。

【 0 0 8 5 】

(ステップ S 2 0 4) 素性取得部 1 5 は、用語対格納部 1 1 から、 i 番目の用語対を読み出す。

【 0 0 8 6 】

(ステップ S 2 0 5) 素性取得部 1 5 は、 i 番目の用語対の素性を取得する処理を行う。素性取得処理について、図 3 のフローチャートを用いて説明する。

20

【 0 0 8 7 】

(ステップ S 2 0 6) 機械学習部 1 6 は、ステップ S 2 0 5 で取得された複数の素性と、学習データ格納部 1 2 の 2 以上の学習データとを用いて、教師あり機械学習を行い、スコアを取得する。

【 0 0 8 8 】

(ステップ S 2 0 7) 出力部 1 7 は、ステップ S 2 0 6 で取得されたスコアを用いて、 i 番目の用語対は異表記の用語対であるか否かを判断する。例えば、出力部 1 7 を構成する閾値判断手段 1 7 2 は、閾値格納手段 1 7 1 から閾値を読み出し、ステップ S 2 0 6 で取得されたスコアが閾値より大きいまたは閾値以上であれば、 i 番目の用語対は異表記の用語対であると判断し、スコアが閾値以下または閾値より小さい場合は、 i 番目の用語対は異表記の用語対でない、と判断する。

30

【 0 0 8 9 】

(ステップ S 2 0 8) 出力部 1 7 は、ステップ S 2 0 8 での判断結果が、異表記の用語対であればステップ S 2 0 9 に行き、異表記の用語対でなければステップ S 2 1 0 に行く。

【 0 0 9 0 】

(ステップ S 2 0 9) 出力部 1 7 は、 i 番目の用語対を異表記の用語対であるとして出力する。

【 0 0 9 1 】

(ステップ S 2 1 0) 素性取得部 1 5 は、カウンタ i を 1 , インクリメントする。ステップ S 2 0 3 に戻る。

40

【 0 0 9 2 】

次に、ステップ S 2 0 5 の素性取得処理について、図 3 のフローチャートを用いて説明する。

【 0 0 9 3 】

(ステップ S 3 0 1) 素性取得部 1 5 を構成する差分文字取得手段 1 5 1 は、2 つの用語の編集箇所を取得する。

【 0 0 9 4 】

(ステップ S 3 0 2) 素性取得部 1 5 の素性取得手段 1 5 2 は、ステップ S 3 0 1 で取得された編集箇所を用いて、字種関連素性を取得する。字種関連素性の具体的な取得方法

50

については後述する。

【0095】

(ステップS303)素性取得手段152は、用語辞書13を用いて、辞書関連素性を取得する。辞書関連素性の具体的な取得方法については後述する。

【0096】

(ステップS304)素性取得手段152は、2つの用語の類似度を取得する。この類似度は、類似度素性である。

【0097】

(ステップS305)素性取得手段152は、その他、予め決められた素性を取得する。その他の予め決められた素性の例は、後述する。

【0098】

(ステップS306)素性取得手段152は、スタッキングアルゴリズムを使用して、ステップS302からステップS305において取得した複数の素性を用いて、ステップS206における教師あり機械学習法とは異なる分類方法により、用語対が異表記の用語対であるか否かを判断し、その判断結果を取得する。

【0099】

以下、本実施の形態における異表記取得装置1の具体的な動作について説明する。

【0100】

今、用語辞書13は、例えば、図4に示すような構造を有する、とする。図4において、一用語の情報は、レコードになっている。各レコードは、「用語」「読み」「品詞」「代表表記」「カテゴリ」「ドメイン」の属性値を有する。用語辞書13は、例えば、JUMAN辞書(「<http://nlp.kuee.kyoto-u.ac.jp/nl-resource/juman.html>」参照[平成21年12月13日検索])である。また、用語辞書13は、例えば、日本語ワードネット辞書(<http://nlpwww.nict.go.jp/wn-ja/index.ja.html>参照[平成21年12月13日検索])や、異体字辞書や、EDR電子化辞書(http://www2.nict.go.jp/r/r312/EDR/J_index.html参照[平成21年12月13日検索])である。なお、異体字辞書とは、異体字の対を有する辞書である。異体字とは、読み方や用い方が同じでも字形に異なる部分のある字体のことである。旧字体と新字体がある漢字などに多く見られ、例えば「沢」と「澤」は異体字の関係にある。なお、異体字辞書は、異体字ではないが、異体字のように代替可能な漢字の対を有しても良い。さらに、用語辞書13は、異体字辞書とは別に、異体字のように代替可能な漢字の対を有する辞書を有しても良い。

【0101】

また、学習データ格納部12に格納されている学習データが有する複数の素性、および素性取得部15が取得する複数の素性は、ここでは、68種類である、とする。以下に、68の素性(S1からS68)について説明する。また、以下、用語対の具体例として、用語対「ショウウインドウ」「ショーウインドウ」を用いて、素性を例示する。

【0102】

S1は、「一目の表記の編集箇所」であり、上記具体例では、「ウ」である。素性S1を取得する場合、差分文字取得手段151は、用語対の構成する2つの用語を1文字ずつずらしながら文字を比較し、編集箇所を得る。例えば、差分文字取得手段151は、「ショウウインドウ」の1文字目「シ」と、「ショーウインドウ」の1文字目「シ」とから比較し、同一と判断し、2文字目も同一と判断し、3文字目「ウ」と「ー」とが異なると判断し、一目の表記の編集箇所「ウ」と二目の表記の編集箇所「ー」を取得する。

【0103】

S2は、「二目の表記の編集箇所」であり、上記具体例では、「ー」である。一目の表記とは用語対を構成する一目の用語(例えば、「ショウウインドウ」)であり、二目の表記とは用語対を構成する二目の用語(例えば、「ショーウインドウ」)である。

【0104】

S3は「編集箇所の前方の1文字」であり、上記具体例では、「ヨ」である。

10

20

30

40

50

- 【0105】
S4は「編集箇所の後方の1文字」であり、上記具体例では、「ウ」である。
- 【0106】
S5は、「編集箇所の前方の連続する2文字」であり、上記具体例では、「ショ」である。
- 【0107】
S6は、「編集箇所の前方の連続する3文字」であり、上記具体例では、「ショ」である。
- 【0108】
S7は、「編集箇所の前方2文字目の文字」であり、上記具体例では、「シ」である。 10
- 【0109】
S8は「編集箇所の前方3文字目の文字」であり、上記具体例では、「(del)」である。(del)とは、文字が無いことを示す。
- 【0110】
S9は「編集箇所の後方の2文字」であり、上記具体例では、「ウイ」である。
- 【0111】
S10は「編集箇所の後方の3文字」であり、上記具体例では、「ウイン」である。
- 【0112】
S11は「編集箇所の後方2文字目の文字」であり、上記具体例では、「イ」である。
- 【0113】
S12は「編集箇所の後方3文字目の文字」であり、上記具体例では、「ン」である。 20
- 【0114】
S13は「'S1の情報 - S2の情報'とした文字列」であり、上記具体例では、「ウー」である。
- 【0115】
S14は「'S3の情報 - S13の情報'とした文字列」であり、上記具体例では、「ヨ - ウ - ー」である。
- 【0116】
S15は「'S5の情報 - S13の情報'とした文字列」であり、上記具体例では、「シヨ - ウ - ー」である。 30
- 【0117】
S16は「'S6の情報 - S13の情報'とした文字列」であり、上記具体例では、「シヨ - ウ - ー」である。
- 【0118】
S17は「'S13の情報 - S4の情報'」であり、上記具体例では、「ウ - ー - ウ」である。
- 【0119】
S18は「'S3の情報 - S13の情報 - S4の情報'」であり、上記具体例では、「ヨ - ウ - ー - ウ」である。
- 【0120】
S19は「'S5の情報 - S13の情報 - S4の情報'とした文字列」であり、上記具体例では、「シヨ - ウ - ー - ウ」である。 40
- 【0121】
S20は「'S6の情報 - S13の情報 - S4の情報'とした文字列」であり、上記具体例では、「シヨ - ウ - ー - ウ」である。
- 【0122】
S21は「'S13の情報 - S7の情報'とした文字列」であり、上記具体例では、「ウ - ー - シ」である。
- 【0123】
S22は「'S3の情報 - S13の情報 - S7の情報'とした文字列」であり、上記具体 50

例では、「ヨ - ウ - - - シ」である。

【0124】

S23は「'S5の情報 - S13の情報 - S7の情報'とした文字列」であり、上記具体例では、「シヨ - ウ - - - シ」である。

【0125】

S24は「'S6の情報 - S13の情報 - S7の情報'とした文字列」であり、上記具体例では、「シヨ - ウ - - - シ」である。

【0126】

S25は「'S13の情報 - S8の情報'とした文字列」であり、上記具体例では、「ウ - - - (del)」である。

10

【0127】

S26は「'S3の情報 - S13の情報 - S8の情報'とした文字列」であり、上記具体例では、「ヨ - ウ - - - (del)」である。

【0128】

S27は「'S5の情報 - S13の情報 - S8の情報'とした文字列」であり、上記具体例では、「シヨ - ウ - - - (del)」である。

【0129】

S28は「'S6 の情報 - S13の情報 - S8の情報'とした文字列」であり、上記具体例では、「シヨ - ウ - - - (del)」である。なお、2つの用語が与えられ、編集箇所が判断できれば、単なる文字列の処理（操作）により、素性取得手段152は、S3からS28の素性を取得できる。

20

【0130】

S29は「S1の字種」であり、上記具体例では、「カタカナ」である。文字を与えられた場合、当該文字の字種（漢字、ひらがな、カタカナ、アルファベット等）を取得する技術は公知技術である。

【0131】

S30は「S2の字種」であり、上記具体例では、「カタカナ」である。

【0132】

S31は「S3の字種」であり、上記具体例では、「カタカナ」である。

【0133】

S32は「S4の字種」であり、上記具体例では、「カタカナ」である。

30

【0134】

S33は「S13の字種」であり、上記具体例では、「カタカナ」である。

【0135】

S34は「S14の字種」であり、上記具体例では、「カタカナ」である。

【0136】

S35は「S17の字種」であり、上記具体例では、「カタカナ」である。

【0137】

S36は「S18の字種」であり、上記具体例では、「カタカナ」である。

【0138】

S37は「S1の品詞」であり、上記具体例では、「名詞」である。ここで、文字の品詞は、その文字（ここではS1）が属している用語の品詞である。例えば、用語に対して、形態素解析をかけ、用語を単語に区切り、品詞情報を取得する。そして、文字の品詞は、当該取得した品詞情報が示す品詞である。

40

【0139】

S38は「S2の品詞」であり、上記具体例では、「名詞」である。

【0140】

S39は「S3の品詞」であり、上記具体例では、「名詞」である。

【0141】

S40は「S4の品詞」であり、上記具体例では、「名詞」である。

50

- 【 0 1 4 2 】
S 4 1 は「 S 1 3 の品詞 」であり、上記具体例では、「名詞」である。
- 【 0 1 4 3 】
S 4 2 は「 S 1 4 の品詞 」であり、上記具体例では、「名詞」である。
- 【 0 1 4 4 】
S 4 3 は「 S 1 7 の品詞 」であり、上記具体例では、「名詞」である。
- 【 0 1 4 5 】
S 4 4 は「 S 1 8 の品詞 」であり、上記具体例では、「名詞」である。
- 【 0 1 4 6 】
S 4 5 は「 S 1 の品詞と位置情報 」であり、上記具体例では、「名詞 , 3 」である。こ
こで「 3 」は、3文字目であることを示す。 10
- 【 0 1 4 7 】
S 4 6 は「 S 2 の品詞と位置情報 」であり、上記具体例では、「名詞 , 3 」である。
- 【 0 1 4 8 】
S 4 7 は「 S 3 の品詞と位置情報 」であり、上記具体例では、「名詞 , 2 」である。
- 【 0 1 4 9 】
S 4 8 は「 S 4 の品詞と位置情報 」であり、上記具体例では、「名詞 , 4 」である。
- 【 0 1 5 0 】
S 4 9 は「 S 1 3 の品詞と位置情報 」であり、上記具体例では、「名詞 , 3 」である。
- 【 0 1 5 1 】
S 5 0 は「 S 1 4 の品詞と位置情報 」であり、上記具体例では、「名詞 , 2 」である。 20
- 【 0 1 5 2 】
S 5 1 は「 S 1 7 の品詞と位置情報 」であり、上記具体例では、「名詞 , 6 」である。
- 【 0 1 5 3 】
S 5 2 は「 S 1 8 の品詞と位置情報 」であり、上記具体例では、「名詞 , 3 」である。
- 【 0 1 5 4 】
S 5 3 は「日本語用語対の類似度」であり、上記具体例では、例えば、0 . 9 である。
- 【 0 1 5 5 】
S 5 4 は「スタッキングアルゴリズムを使用して、日本語用語対の J U M A N 辞書の代
表表記が一致するかどうか」を示す情報であり、上記具体例では、「 1 」である。つま
り、ここでは、機械学習部 1 6 が利用する教師あり機械学習法とは異なる分類方法は、用語
対を構成する 2 つの用語の、J U M A N 辞書における代表表記が一致するか否かにより分
類する以下の方法である。まず、J U M A N 辞書の単語の集合から、編集距離が 1 文字の
単語対を取り出す。ここで、編集距離が 1 文字の単語対は、9 0 4 6 1 2 組、取り出せる
。そのうち、代表表記が等しい単語対 (2 5 9 3 4 組) を取り出す。次に、J U M A N 辞
書で、代表表記が等しい単語対を正例、そうでないものを負例とする。以上により、教師
データを作成する。次に、その教師データを学習データとした機械学習を行う。なお、教
師データは、上述した教師データに限らず、他の教師データを用いてもよい。また、機
械学習の際に利用する素性は、本発明の全素性 (S 1 から S 6 8 の素性) のうち、S 5 4 の
素性を取り除いた素性を利用する。なお、機械学習の際に利用する素性は、他の素性を用
いてもよい。そして、実際に、S 5 4 の素性を付与したいデータを、上記学習結果を利用
して、分類する。分類結果において正例となったか、負例となったかの情報を S 5 4 の素
性として、そのデータに付与する。スタッキングアルゴリズムによる方法では、J U M A N
辞書で、代表表記が一致するか否かについて学習した結果を素性として付与できるので
、実際に J U M A N 辞書に記載されていない用語対に対しても、J U M A N 辞書で、代表
表記が一致するとされる傾向のある用語対が否かの情報を付与できることとなる。 30
- 【 0 1 5 6 】
S 5 5 は「日本語用語対の文字数が同数で編集箇所が両方とも数字の場合であり、同じ
値が違ふ値かどうか」であり、上記具体例では、「 0 」である。なお、「 2 次キャッシュ
」と「二次キャッシュ」の用語対の場合、「一週間あたり」「 1 週間あたり」の用語対の 40

場合は、S55の素性は「1」となる。なお、文字数が同数である条件をはずし、S55は、「日本語用語対の編集箇所が両方とも数字の場合であり、同じ値か違う値かどうか」が好適である。

【0157】

S56は「日本語用語対の文字数が同数で編集箇所が両方ともひらがなの場合であり、同じ音声か違う音声かどうか」であり、上記具体例では、「0」である。なお、「おかあちゃん」「おかあちゃん」の用語対の場合は、S56の素性は「1」となる。なお、文字数が同数である条件をはずし、S56は、「日本語用語対の編集箇所が両方ともひらがなの場合であり、同じ音声か違う音声かどうか」が好適である。

【0158】

S57は「日本語用語対の文字数が同数で編集箇所が両方ともカタカナの場合であり、同じ音声か違う音声かどうか」であり、上記具体例では、「1」である。なお、「オリーブ・オイル」「オリーブ・オイル」の用語対の場合、「ウインドウ」「ウインドウ」の用語対の場合も、S57の素性は「1」となる。なお、文字数が同数である条件をはずし、S57は、「日本語用語対の編集箇所が両方ともカタカナの場合であり、同じ音声か違う音声かどうか」が好適である。

【0159】

S58は「日本語用語対の文字数が同数で編集箇所が両方ともローマ字の場合であり、大文字と小文字の違いだけかどうか」であり、上記具体例では、「0」である。なお、「300kbps」「300Kbps」の用語対の場合、「Windows上」「windows上」の用語対の場合は、S58の素性は「1」となる。なお、文字数が同数である条件をはずし、S58は、「日本語用語対の編集箇所が両方ともローマ字の場合であり、大文字と小文字の違いだけかどうか」が好適である。また、「Windows」は登録商標です。

【0160】

S59は「日本語用語対の文字数が同数で一方の編集箇所に濁点をつけるともう一方の編集箇所になるかどうか」であり、上記具体例では、「0」である。なお、「触れるくらい」「触れるくらい」の用語対の場合、「飲むくらい」「飲むくらい」の用語対の場合は、S59の素性は「1」となる。なお、S59は、文字数が同数である条件をはずし、「日本語用語対の一方の編集箇所に濁点をつけるともう一方の編集箇所になるかどうか」が好適である。

【0161】

S60は「日本語用語対の文字数が同数で一方の編集箇所に半濁点をつけるともう一方の編集箇所になるかどうか」であり、上記具体例では、「0」である。なお、S60は、文字数が同数である条件をはずし、「日本語用語対の一方の編集箇所に半濁点をつけるともう一方の編集箇所になるかどうか」が好適である。

【0162】

S61は「編集箇所が日本語用語対の一方にしかなく、その編集箇所が'化'、'系'、'類'、'型'、'形'、'氏'、'ー'、'・'かどうか」であり、上記具体例では、「0」である。なお、「サーバ」「サーバー」の用語対の場合、「ハンセン病患者」「ハンセン氏病患者」の用語対の場合、「日本語パッチ」「日本語化パッチ」の用語対の場合、「30種類ほど」「30種ほど」の用語対の場合は、S61の素性は「1」となる。

【0163】

S62は「編集箇所が日本語用語対の一方にしかなく、その編集箇所の用語が日本語用語対の最後の文字と一致するかどうか」であり、上記具体例では、「0」である。なお、「妊娠・授乳中」「妊娠中・授乳中」の用語対の場合、「国産・輸入車」「国産車・輸入車」の用語対の場合は、S62の素性は「1」となる。

【0164】

S63は「編集箇所が日本語用語対の一方にしかなく、その編集箇所が桁数をあらわす用語かどうか（例えば、「千」「万」など）」であり、上記具体例では、「0」である。なお

10

20

30

40

50

、「2万5000人」「25000人」の用語対の場合、「1万6500円」「16500円」の用語対の場合は、S63の素性は「1」となる。

【0165】

S64は「日本語用語対のJUMAN辞書の定義されている代表表記が一致するかどうか」であり、上記具体例（用語対「ショウウインドウ」「ショーウインドウ」）では、例えば、「1」である。なお、素性取得手段152は、用語対を構成する各用語の代表表記を、用語辞書13から取得し、比較することにより、素性を取得する。

【0166】

S65は「日本語用語対が日本語ワードネット辞書に類義語対として定義されているかどうか」であり、上記具体例（用語対「ショウウインドウ」「ショーウインドウ」）では、例えば、「1」である。素性取得手段152は、用語対を構成する各用語をキーとして、日本語ワードネット辞書を検索し、類義語対として定義されているか否かを判断する。本処理は、通常の検索処理である。

10

【0167】

S66は「日本語用語対の編集箇所が異体字辞書に異体字として定義されているかどうか」であり、上記具体例（用語対「ショウウインドウ」「ショーウインドウ」）では、例えば、「0」である。異体字辞書は、2つの異体字の対の情報を有する。

【0168】

S67は「日本語用語対の文字数が同数で編集箇所が漢字とひらがなの場合であり、JUMAN辞書の読みが一致するかどうか」であり、上記具体例（用語対「ショウウインドウ」「ショーウインドウ」）では、「0」である。

20

【0169】

S68は「日本語用語対の文字数が同数で編集箇所が両方とも漢字の場合であり、JUMAN辞書の読みが一致するかどうか」であり、上記具体例（用語対「ショウウインドウ」「ショーウインドウ」）では、「0」である。

【0170】

また、上記の68の素性は、上述したように、字種関連素性、辞書関連素性、類似度素性、編集箇所文字素性などが含まれる。

【0171】

また、上記の68の素性をグループ化すると、例えば、以下のG1からG7のグループに分かれる、と考えられる。

30

【0172】

G1は、S1からS52の素性であり、編集箇所とその周辺の文字列に関する情報である編集箇所文字素性である。

【0173】

G2は、S53の素性であり、類似度素性である。

【0174】

G3は、S54の素性であり、スタッキングアルゴリズムを使用した情報である素性である。

【0175】

G4は、S55からS60の素性であり、編集箇所に関する情報である編集箇所関連素性である。

40

【0176】

G5は、S61からS63の素性であり、用語対のパターンに関する情報である用語対パターン素性である。

【0177】

G6は、S64からS66の素性であり、種々の辞書による情報である辞書関連素性である。

【0178】

G7は、S67からS68の素性であり、読みに関する情報である読み関連素性である

50

。

【0179】

そして、まず、異表記取得装置1において、予め正しい異表記の用語対のデータ(正例)を手で構築しておき、正例の用語対と、正例であることを示す正負情報(例えば、「1」と)を対応付けて、学習データ格納部12に格納しておく。また、異表記取得装置1において、予め異表記でない用語対のデータ(負例)を手で構築しておき、負例の用語対と、負例であることを示す正負情報(例えば、「0」と)を対応付けて、学習データ格納部12に格納しておく。

【0180】

次に、異表記取得装置1の素性取得部15により、各用語対の、上述した68の素性を取得し、正負情報または用語対と対応付けて、68の素性を学習データ格納部12に蓄積する。

10

【0181】

以上の処理により、学習データ格納部12の学習データが構築された。

【0182】

次に、異表記の用語対であるか否かを判断したい1以上の用語対を用語対格納部11に格納する。

【0183】

そして、ユーザは、異表記取得装置1に、動作開始の指示を入力する。すると、受付部14は、動作開始の指示を受け付ける。

20

【0184】

次に、用語対格納部11に格納されている用語対を順に、以下のように処理する。つまり、素性取得部15は、各用語対の素性を取得する処理を行う。かかる素性取得処理については説明済みである。

【0185】

次に、機械学習部16は、取得された68の素性と、学習データ格納部12の学習データを用いて、教師あり機械学習を行い、スコアを取得する。

【0186】

次に、出力部17は、取得された各用語対のスコアを用いて、各用語対は異表記の用語対であるか否かを判断する。

30

【0187】

次に、出力部17は、異表記の用語対であると判断した用語対のみ出力する。ここで、出力とは、たとえば、予め決められた記憶媒体への蓄積である。

(機械学習について)

【0188】

以下、機械学習部16が行う機械学習、および機械学習部16が行う教師あり機械学習法とは異なる分類方法(スタッキングアルゴリズムで利用)について説明する。

【0189】

まず、機械学習法とは、問題-解の組のセットを多く用意し、それで学習を行ない、どういう問題のときにどういう解になるかを学習し、その学習結果を利用して、新しい問題のときも解を推測できるようにする方法である(例えば、非特許文献2~非特許文献4参照)。

40

【0190】

どういう問題のときにどういう解になるかという、問題の状況を機械に伝える際に、素性(解析に用いる情報で問題を構成する各要素)が必要になる。問題を素性によって表現するのである。

【0191】

すなわち、機械学習の手法は、素性の集合-解の組のセットを多く用意し、それで学習を行ない、どういう素性の集合のときにどういう解になるかを学習し、その学習結果を利用して、新しい問題のときもその問題から素性の集合を取り出し、その素性の場合の解を

50

推測する方法である。

【0192】

機械学習の手法として、例えば、k近傍法、シンプルベイズ法、決定リスト法、最大エントロピー法、サポートベクトルマシン法などの手法を用いることができる。

【0193】

k近傍法は、最も類似する一つの事例のかわりに、最も類似するk個の事例を用いて、このk個の事例での多数決によって分類先(解)を求める手法である。kは、あらかじめ定める整数の数字であって、一般的に、1から9の間の奇数を用いる。

【0194】

シンプルベイズ法は、ベイズの定理にもとづいて各分類になる確率を推定し、その確率値が最も大きい分類を求める分類先とする方法である。

10

【0195】

シンプルベイズ法において、文脈bで分類aを出力する確率は、以下の数式1で与えられる。

【数1】

$$p(a|b) = \frac{p(a)}{p(b)} p(b|a)$$

$$\cong \frac{\tilde{p}(a)}{p(b)} \prod_i \tilde{p}(f_i|a)$$

20

【0196】

ただし、ここで文脈bは、あらかじめ設定しておいた素性 f_j ($F, 1 \leq j \leq k$)の集合である。 $p(b)$ は、文脈bの出現確率である。ここで、分類aに非依存であって定数のために計算しない。 $P(a)$ (ここでPはpの上部にチルダ)と $P(f_i|a)$ は、それぞれ教師データ(判断情報と同意義)から推定された確率であって、分類aの出現確率、分類aのときに素性 f_i を持つ確率を意味する。 $P(f_i|a)$ として最尤推定を行って求めた値を用いると、しばしば値がゼロとなり、数式2の2行目の式の値がゼロで分類先を決定することが困難な場合が生じる。そのため、スムージングを行う。ここでは、以下の数式2を用いてスムージングを行ったものを用いる。

30

【数2】

$$p(f_i|a) = \frac{\text{freq}(f_i, a) + 0.01 * \text{freq}(a)}{\text{freq}(a) + 0.01 * \text{freq}(a)}$$

【0197】

ただし、 $\text{freq}(f_i, a)$ は、素性 f_i を持ちかつ分類がaである事例の個数、 $\text{freq}(a)$ は、分類がaである事例の個数を意味する。

【0198】

決定リスト法は、素性と分類先の組とを規則とし、それらをあらかじめ定めた優先順序でリストに蓄えおき、検出する対象となる入力を与えられたときに、リストで優先順位の高いところから入力のデータと規則の素性とを比較し、素性が一致した規則の分類先をその入力の分類先とする方法である。

40

【0199】

決定リスト方法では、あらかじめ設定しておいた素性 f_j ($F, 1 \leq j \leq k$)のうち、いずれか一つの素性のみを文脈として各分類の確率値を求める。ある文脈bで分類aを出力する確率は以下の数式3によって与えられる。

【数3】

$$p(a|b) = p(a|f_{\max})$$

50

【 0 2 0 0 】

ただし、 f_{\max} は以下の数式 4 によって与えられる。

【 数 4 】

$$f_{\max} = \operatorname{argmax}_{f_j \in F} \max_{a_i \in A} \tilde{p}(a_i | f_j)$$

【 0 2 0 1 】

また、 $P(a_i | f_j)$ (ここで P は p の上部にチルダ) は、素性 f_j を文脈に持つ場合の分類 a_i の出現の割合である。

【 0 2 0 2 】

最大エントロピー法は、あらかじめ設定しておいた素性 f_j ($1 \leq j \leq k$) の集合を F とするとき、以下の所定の条件式 (数式 5) を満足しながらエントロピーを意味する式 (数式 6) を最大にするときの確率分布 $p(a, b)$ を求め、その確率分布にしたがって求める各分類の確率のうち、最も大きい確率値を持つ分類を求める分類先とする方法である。

【 数 5 】

$$\sum_{a \in A, b \in B} p(a, b) g_j(a, b) = \sum_{a \in A, b \in B} \tilde{p}(a, b) g_j(a, b)$$

$$\text{for } \forall f_j (1 \leq j \leq k)$$

【 数 6 】

$$H(p) = - \sum_{a \in A, b \in B} p(a, b) \log(p(a, b))$$

【 0 2 0 3 】

ただし、 A 、 B は分類と文脈の集合を意味し、 $g_j(a, b)$ は文脈 b に素性 f_j があって、なおかつ分類が a の場合 1 となり、それ以外で 0 となる関数を意味する。また、 $P(a_i | f_j)$ (ここで P は p の上部にチルダ) は、既知データでの (a, b) の出現の割合を意味する。

【 0 2 0 4 】

数式 5 は、確率 p と出力と素性の組の出現を意味する関数 g をかけることで出力と素性の組の頻度の期待値を求めることになっており、右辺の既知データにおける期待値と、左辺の求める確率分布に基づいて計算される期待値が等しいことを制約として、エントロピー最大化 (確率分布の平滑化) を行なって、出力と文脈の確率分布を求めるものとなっている。最大エントロピー法の詳細については、以下の参考文献 (1) および参考文献 (2) に記載されている。

【 0 2 0 5 】

参考文献 (1) : Eric Sven Ristad, Maximum Entropy Modeling for Natural Language, (ACL/EACL Tutorial Program, Madrid, 1997)

【 0 2 0 6 】

参考文献 (2) : Eric Sven Ristad, Maximum Entropy Modeling Toolkit, Release 1.6 beta, (<http://www.mnemonic.com/software/memt>, 1998))

【 0 2 0 7 】

サポートベクトルマシン法は、空間を超平面で分割することにより、二つの分類からなるデータを分類する手法である。

【 0 2 0 8 】

図 5 にサポートベクトルマシン法のマージン最大化の概念を示す。図 5 において、白丸は正例、黒丸は負例を意味し、実線は空間を分割する超平面を意味し、破線はマージン領域の境界を表す面を意味する。図 5 (A) は、正例と負例の間隔が狭い場合 (スモールマージン) の概念図、図 5 (B) は、正例と負例の間隔が広い場合 (ラージマージン) の概念図である。

【 0 2 0 9 】

10

20

30

40

50

このとき、二つの分類が正例と負例からなるものとする、学習データにおける正例と負例の間隔（マージン）が大きいものほどオープンデータで誤った分類をする可能性が低いと考えられ、図5（B）に示すように、このマージンを最大にする超平面を求めそれを用いて分類を行なう。

【0210】

基本的には上記のとおりであるが、通常、学習データにおいてマージンの内部領域に少数の事例が含まれてもよいとする手法の拡張や、超平面の線形の部分を非線型にする拡張（カーネル関数の導入）がなされたものが用いられる。

【0211】

この拡張された方法は、以下の識別関数（ $f(x)$ ）を用いて分類することと等価であり、その識別関数の出力値が正か負かによって二つの分類を判別することができる。

【数7】

$$f(x) = \text{sgn} \left(\sum_{i=1}^l \alpha_i y_i K(x_i, x) + b \right)$$

$$b = - \frac{\max_{i, y_j = -1} b_i + \min_{i, y_j = 1} b_i}{2}$$

$$b_i = \sum_{j=1}^l \alpha_j y_j K(x_j, x_i)$$

【0212】

ただし、 x は識別したい事例の文脈（素性の集合）を、 x_i と y_j （ $i = 1, \dots, l$, $y_j \in \{1, -1\}$ ）は学習データの文脈と分類先を意味し、関数 sgn は、

$$\text{sgn}(x) = \begin{cases} 1 & (x \geq 0) \\ -1 & (\text{otherwise}) \end{cases}$$

であり、また、各 α_j は数式8の式（8-2）と式（8-3）の制約のもと、式（8-1）を最大にする場合のものである。

【数8】

$$L(\alpha) = \sum_{i=1}^l \alpha_i - \frac{1}{2} \sum_{i,j=1}^l \alpha_i \alpha_j y_i y_j K(x_i, x_j) \quad \text{式(8-1)}$$

$$0 \leq \alpha_i \leq C \quad (i=1, \dots, l) \quad \text{式(8-2)}$$

$$\sum_{j=1}^l \alpha_j y_j = 0 \quad \text{式(8-3)}$$

【0213】

また、関数 K はカーネル関数と呼ばれ、様々なものが用いられるが、本形態では、例えば、以下の多項式（数式9）のものを用いる。

【数9】

$$K(x, y) = (x \cdot y + 1)^d$$

【0214】

数式8、数式9において、 C 、 d は実験的に設定される定数である。例えば、 C はすべての処理を通して1に固定した。また、 d は、1と2の二種類を試している。ここで、 $\alpha_i > 0$ となる x_i は、サポートベクトルと呼ばれ、通常、数式7の和をとっている部分は、この事例のみを用いて計算される。つまり、実際の解析には学習データのうちサポートベクトルと呼ばれる事例のみしか用いられない。

【0215】

10

20

30

40

50

なお、拡張されたサポートベクトルマシン法の詳細については、以下の参考文献(3)および参考文献(4)に記載されている。

【0216】

参考文献(3) : Nello Cristianini and John Shawe-Taylor, An Introduction to Support Vector Machines and other kernel-based learning methods, (Cambridge University Press, 2000)

【0217】

参考文献(4) : Taku Kudoh, Tinysvm: Support Vector machines, ([http://cl.aistnara.ac.jp/taku-ku//software/Tiny SVM/index.html](http://cl.aistnara.ac.jp/taku-ku//software/Tiny_SVM/index.html), 2000)

【0218】

サポートベクトルマシン法は、分類の数が2個のデータを扱うものである。したがって、分類の数が3個以上の事例を扱う場合には、通常、これにペアワイズ法またはワンVSレスト法などの手法を組み合わせる用いることになる。

【0219】

ペアワイズ法は、 n 個の分類を持つデータの場合に、異なる二つの分類先のあらゆるペア($n(n-1)/2$ 個)を生成し、各ペアごとにどちらがよいかを二値分類器、すなわちサポートベクトルマシン法処理モジュールで求めて、最終的に、 $n(n-1)/2$ 個の二値分類による分類先の多数決によって、分類先を求める方法である。

【0220】

ワンVSレスト法は、例えば、 a 、 b 、 c という三つの分類先があるときは、分類先 a とその他、分類先 b とその他、分類先 c とその他、という三つの組を生成し、それぞれの組についてサポートベクトルマシン法で学習処理する。そして、学習結果による推定処理において、その三つの組のサポートベクトルマシンの学習結果を利用する。推定すべき問題が、その三つのサポートベクトルマシンではどのように推定されるかを見て、その三つのサポートベクトルマシンのうち、その他でないほうの分類先であって、かつサポートベクトルマシンの分離平面から最も離れた場合のものの分類先を求める解とする方法である。例えば、ある解くべき問題が、「分類先 a とその他」の組の学習処理で作成したサポートベクトルマシンにおいて分離平面から最も離れた場合には、その解くべき問題の分類先は、 a と推定する。

【0221】

機械学習部16が推定する、解くべき問題についての、どのような解(分類先)になりやすいかの度合いの求め方は、機械学習部16が機械学習の手法として用いる様々な方法によって異なる。

【0222】

例えば、本発明の実施の形態において、機械学習部16が、機械学習の手法として k 近傍法を用いる場合、機械学習部16は、教師データの事例同士で、その事例から抽出された素性の集合のうち重複する素性の割合(同じ素性をいくつ持っているかの割合)にもとづく事例同士の類似度を定義して、前記定義した類似度と事例とを学習結果情報として学習データ格納部12に記憶しておく。

【0223】

そして、機械学習部16は、素性取得部15によって解くべき問題の素性が抽出されたときに、記憶された類似度と事例を参照して、素性取得部15によって抽出された解くべき問題の素性について、その解くべき問題の素性の類似度が高い順に k 個の事例を選択し、選択した k 個の事例での多数決によって決まった分類先を、解くべき問題の分類先(解)として推定する。すなわち、機械学習部16では、解くべき問題についての、どのような解(分類先)になりやすいかの度合いを、選択した k 個の事例での多数決の票数とする。

【0224】

また、機械学習手法として、シンプルベイズ法を用いる場合には、教師データの事例について、前記事例の解と素性の集合との組を学習データとして学習データ格納部12に記

10

20

30

40

50

憶する。そして、機械学習部 16 は、素性取得部 15 によって解くべき問題の素性が抽出されたときに、学習データ格納部 12 の判断情報の解と素性の集合との組をもとに、ベイズの定理にもとづいて素性取得部 15 で取得した解くべき問題の素性の集合の場合の各分類になる確率を算出して、その確率の値が最も大きい分類を、その解くべき問題の素性の分類（解）と推定する。すなわち、機械学習部 16 では、解くべき問題の素性の集合の場合にある解となりやすさの度合いを、各分類になる確率とする。

【0225】

また、機械学習手法として決定リスト法を用いる場合には、教師データの事例について、素性と分類先との規則を所定の優先順序で並べたリストを、予め、何らかの手段により、学習データ格納部 12 に記憶させる。そして、素性取得部 15 によって解くべき問題の素性が抽出されたときに、機械学習部 16 は、学習データ格納部 12 のリストの優先順位の高い順に、抽出された解くべき問題の素性と規則の素性とを比較し、素性が一致した規則の分類先をその解くべき問題の分類先（解）として推定する。

10

【0226】

また、機械学習手法として最大エントロピー法を使用する場合には、教師データの事例から解となりうる分類を特定し、所定の条件式を満足し、かつエントロピーを示す式を最大にするときの素性の集合と解となりうる分類の二項からなる確率分布を求めて、学習データ格納部 12 に記憶する。そして、素性取得部 15 によって解くべき問題の素性が抽出されたときに、機械学習部 16 は、学習データ格納部 12 の確率分布を利用して、抽出された解くべき問題の素性の集合についてその解となりうる分類の確率を求めて、最も大きい確率値を持つ解となりうる分類を特定し、その特定した分類をその解くべき問題の解と推定する。すなわち、機械学習部 16 では、解くべき問題の素性の集合の場合にある解となりやすさの度合いを、各分類になる確率とする。

20

【0227】

また、機械学習手法としてサポートベクトルマシン法を使用する場合には、教師データの事例から解となりうる分類を特定し、分類を正例と負例に分割して、カーネル関数を用いた所定の実行関数にしたがって事例の素性の集合を次元とする空間上で、その事例の正例と負例の間隔を最大にし、かつ正例と負例を超平面で分割する超平面を求めて学習データ格納部 12 に記憶する。そして、素性取得部 15 によって解くべき問題の素性が抽出されたときに、機械学習部 16 は、学習データ格納部 12 の超平面を利用して、解くべき問題の素性の集合が超平面で分割された空間において正例側か負例側のどちらにあるかを特定し、その特定された結果にもとづいて定まる分類を、その解くべき問題の解と推定する。すなわち、機械学習部 16 では、解くべき問題の素性の集合の場合にある解となりやすさの度合いを、分離平面からのその解くべき問題の事例への距離の大きさとする。

30

（実験結果 1）

【0228】

以下、異表記取得装置 1 の実験結果について説明する。まず、実験に利用するデータについて説明する。

【0229】

実験で用いるデータは、大規模類似語リストである。大規模類似語リストとは、検索エンジン研究基盤 T S U B A K I (<http://tsubaki.ixnlp.nii.ac.jp/se/index.cgi>参照 [平成 21 年 12 月 13 日検索]) の約 1 億ページ・60 億文のデータから 100 万語を抽出し、その 100 万語の各々の語に対して最大 500 個の類義語を類似度付きで生成したものである。この大規模類似語リストに含まれる 100 万語の日本語用語と、その日本語用語の各々の類義語の日本語用語を日本語用語対とする。

40

【0230】

そして、大規模類似語リストから、編集距離が 1 の日本語用語対をランダムに 14185 組取り出した。その取り出した日本語用語対が日本語異表記対であるか、日本語異表記対でないかのタグ付けを 3 人の評定者の多数決により行った。3 人の評定者のタグ付けがどれくらい一致しているのか、カップバ統計量 K を用いて判定する。

50

【0231】

右の用語と左の用語の2つの用語を有する日本語用語対の組み合わせが左にあるか、右にあるかにより、異なる情報になる素性がある。その素性に対応し、情報量を増やすために、日本語用語対の組み合わせを左右入れ替えたデータも用いる。つまり、本実験では、大規模類似リストから取り出した14185組に加え、合計28370組の実験データを用いる。

【0232】

また、28370組の実験データを1つのまとまったデータであるとする、実験の公正性が失われるのではないかと考え、28470組ある実験データの半分を素性の考案を行うデータとして用いる。残り半分の実験データをクロズドデータで考案された素性が、他のデータにおいても有効であるかどうかの検討を行うデータとして用いる。素性の考案を行うデータをクロズドデータと呼ぶ。検討を行うデータをオープンデータと呼ぶ。クロズドデータは10分割クロスバリデーションによる学習により精度の測定を行う。なお、10分割クロスバリデーションとは、実験対象のデータを、第一から第十の10に分割し、以下の(1)から(10)の学習を行う。(1)第一をテストデータとし、第二から第十を学習データとして、学習を行う。(2)第二をテストデータとし、第一、第三から第十を学習データとして、学習を行う。(3)第三をテストデータとし、第一、第二、第四から第十を学習データとして、学習を行う。(4)第四をテストデータとし、他を学習データとして、学習を行う。(5)第五をテストデータとし、他を学習データとして、学習を行う。(6)第六をテストデータとし、他を学習データとして、学習を行う。(7)第七をテストデータとし、他を学習データとして、学習を行う。(8)第八をテストデータとし、他を学習データとして、学習を行う。(9)第九をテストデータとし、他を学習データとして、学習を行う。(10)第十をテストデータとし、他を学習データとして、学習を行う。なお、10分割クロスバリデーションは、公知技術である。

【0233】

また、クロズドデータを学習データ(学習データ格納部12に格納されるデータ)、オープンデータをテストデータ(異表記の用語対であるか否かを判断されるデータ)とし、オープンクロスによる学習により精度の測定を行う。

【0234】

図6に、実験で用いた編集距離が1の日本語用語対の中に、多数決により日本語異表記対であるか日本語異表記対でないかを判定した内訳を示す。

【0235】

なお、カッパ統計量Kとは、K人評定者のカテゴリ評定における一致度を表す数値のことである。カッパ統計量Kの算出方法は公知であるので、説明を省略する。

【0236】

評定者の間に完全な一致があればKの値は1になる。チャンスレベルでの一致であればKの値は0である。一致度が高くなればKの値は0から1に近づく。図7に、Landsらによる一致度の評価方法を示す。本実験では、14185組の日本語用語対を対象に、3人で日本語異表記対であるか日本語異表記対でないかの2カテゴリでカッパ統計量Kを求めたところ、一致度は0.84であった。これは0.8以上の一致度であるため、ほぼ完全な一致であると評価できる。

【0237】

次に、異表記取得装置1における異表記の用語対であるか否かの判断手法が優れていることを示すために、異表記取得装置1の判断手法と比較対照となるベースライン手法について説明する。

【0238】

編集距離の小さい(例えば、編集距離が1)日本語異表記対の抽出を行う対象の日本語用語対に対して、ベースライン手法では、以下のルールを適用し、機械的に日本語異表記対であるか日本語異表記対でないかについて判定を行う。

(ルール1)文字数が同じ日本語用語対の編集箇所が同じ値の数字を表す場合、日本語異

10

20

30

40

50

表記対であると判定する。

(ルール2)文字数が同じ日本語用語対の編集箇所が同じ意味のアルファベットを表す場合、日本語異表記対であると判定する。

(ルール3)文字数が同じ日本語用語対がJUMANを使い読み方を調べることで、読み方が一致する場合、日本語異表記対であると判定する。

(ルール4)ルール1、ルール2、ルール3と一致しなかった場合、日本語異表記対でないと判定する。

【0239】

上記のルール1からルール4を適用するベースライン手法において、日本語用語対「第2版」「第二版」については、以下のように判断される。この日本語用語対における編集箇所は、「2」と「二」であり数字を表している。ルール1を適用し、同じ値を表しているため、ベースライン手法では、この日本語用語対は日本語異表記対であると判定される。

10

【0240】

日本語用語対「Tea」「tea」については、以下のように判断される。この日本語用語対における編集箇所は、「T」と「t」でありアルファベットを表している。そして、ルール2が適用され、この日本語用語対は、同じ意味の語であるためこの日本語用語対は日本語異表記対であると判定される。

【0241】

日本語用語対「誉める」「褒める」については、以下のように判断される。この日本語用語対は、「ほめる」と「ほめる」にJUMANを使い読み方を調べることで、そして、ルール3が適用され、読み方が一致し、この日本語用語対は日本語異表記対であると判定される。

20

【0242】

日本語用語対「シルキーホワイト」「ミルキーホワイト」については、以下のように判断される。この日本語用語対における編集箇所は、「シ」と「ミ」でありカタカナを表している。シルキーホワイトとミルキーホワイトはJUMAN辞書において未定義であるため、読み方を調べることができない。よってルール4が適用され、日本語異表記対でないと判定される。

【0243】

次に、本実験で用いた機械学習部16の機械学習手法について、詳細に説明する。本機械学習手法は、サポートベクトルマシン法である。サポートベクトルマシン法は、上述したように、空間を超平面で分割することにより、2つの分類からなるデータを分類する手法である。このとき2つの分類が正例と負例からなるとすると、学習データにおいてこの2つの間隔が大きいものほど誤った分類をする可能性が低いと判断される。この間隔を最大にする超平面を求め、それを求めて分類を行うことが基本とされる。しかし、ここでは、学習データにおいて間隔の内部領域に少数の事例を含んでもよいとする手法や超平面の線形の部分を非線形にするなどの拡張がされたものを用いる。これらの拡張された方法は、識別関数を用いて分類することと等価となり、識別関数の出力値が正か負かによって2つに分類を判別することができる。また、3つ以上からなるデータを扱う場合にはペアワイズ手法というのを並行して用いる。ペアワイズ手法はN個の分類をもつデータの場合、異なる2つの分類先のあらゆるペアを作り、各ペアごとにどちらがよいかを2値分類器(サポートベクトルマシン法)で求め最終的に分類先の多数決により求める方法である。以降、サポートベクトルマシン法はSVMと、適宜、表記する。

30

40

【0244】

SVMによって編集距離の小さい(例えば、編集距離が1)日本語異表記対を抽出するために用いる素性は、上述したS1からS68の素性である。これらの素性は、大規模類似語リストからランダムで取り出した編集距離の小さい日本語用語対から取り出す。素性によってそれぞれの機械学習は、日本語用語対が日本語異表記対であるか日本語異表記対でないかを判定をする。日本語用語対からできるだけ多くの情報を得るために、種々の素

50

性を用いた。また、それぞれの素性について、上述したように、G 1 なら G 7 に分類できる。G 1、G 2、G 3 は、すべての編集距離の小さい日本語用語対に対応できる素性である。字種は対象の文字がひらがな、カタカナ、数字、アルファベット、その他のどの種類を表しているかの情報である。品詞は用語に J U M A N を用いて形態素解析をかけ、用語を単語に区切り、品詞情報を取得する。そして、対象の文字がどの品詞に属しているかの情報である。位置情報は、対象の文字が品詞に属している中でさらに、その品詞の先頭、最後尾、それ以外のどの位置を示しているかの情報である。類似度は、大規模類似語リストを生成する際に用いた類似度の情報である。

【0245】

スタッキングアルゴリズムとは、上述したように、実験データを本来の目的とは別の分類方法で分類させたデータを機械学習で学習させ、学習結果の分類情報を素性に加えることである。本実験において、スタッキングアルゴリズムに使用するデータは、実験で用いる 28370 組の日本語用語対以外の、大規模類似語リストから得られた J U M A N の代表表記が判別できる 904612 組の日本語用語対を用いる。904612 組の中で、正例は 25934 組、負例は 878678 組である。これにより、J U M A N 辞書において未定義の日本語用語対にも、近似的ではあるが S 6 4 の素性の情報を付与することができる。また、G 4、G 5 は、特徴がある編集距離の小さい日本語用語対に特化した素性である。G 4 は置換によって等しい文字列で、G 5 は削除によって等しい文字列になる日本語用語対が対象である。G 6、G 7 は、J U M A N 辞書、日本語ワードネット辞書、EDR 辞書を用いた素性である。J U M A N 辞書については、未定義とされている用語が出てくる日本語用語対に対して、素性の情報は付与しないこととする。

10

20

【0246】

上述した S 1 から S 6 8 の素性がどれくらい有効であるのかを有意差の分析により検討する。有意の検討はブートストラップ法を用いて求める。ブートストラップ法とは分類手法によって二つに分類されたデータを用いる。分類された二つのデータをそれぞれ、データ数（例えば、問いの数は 1400）は変えずに重複を許しランダムに取り出す（例えば、取り出したデータ数は 1400）。取り出したデータでそれぞれの F 値を求め、それぞれの F 値を比較する。取り出しと F 値の比較をする工程を 10000 回繰り返す。F 値とは以下の数式 2 で定義される。すなわち再現率と適合率の調和平均である。

30

【数 10】

$$F\text{値} = \frac{2 * \text{再現率} * \text{適合率}}{\text{再現率} + \text{適合率}}$$

【0247】

工程を 10000 回繰り返し比較した結果が、どちらかの手法の F 値よりも、もう一方の手法の F 値の方が高い回数が 9500 回（95%）以上の場合、有意水準 5% により F 値が高い方の手法は有意であるといえる。どちらも 9500 回（95%）以上ない場合、有意水準 5% により有意かどうかの判定はできない。なお、ここでは、有意水準 5% を適用するが、例えば、有意水準 10% を適用しても良い。

40

【0248】

また、本実験では全素性と全素性から 1 種類の素性だけを取り除いたデータを、S V M の学習結果により比較する。上述した 68 個すべて組み合わせた素性を全素性とし、取り除く素性は S 1 から S 6 8 におけるすべての素性でおこなう。この有意差の検討を、クローズドデータを 10 分割クロスバリデーション（10CV）で S V M による学習結果と、クローズドデータとオープンデータを使いオープンクローズ（OC）で S V M による学習結果で行う。なお、オープンクローズとは、クローズドデータを学習データとして、オープンデータをテストデータとして実験することをいう。以降は、10 分割クロスバリデーションによる S V M の実験は 10CV と表記し、オープンクローズによる S V M の実験は OC と表記する。

50

【0249】

次に、ベースライン手法と機械学習を利用した手法について実験を行った結果について報告する。

【0250】

本実験において、上述した大規模類似語リストに含まれる28370組の編集距離の小さい日本語用語対が、日本語異表記対であるか、日本語異表記対でないかについて判定を行った。図8に、用意したクローズドデータとオープンデータに対して、ベースラインの手法を適用した結果を示す。また、図8には、10CVとOCの結果も示す。実験で用いるSVMの実装としてTinySVMを採用し、1次の多項式カーネルでソフトマージンパラメータCを1に設定して利用した。それぞれの表での「全索性」はS1からS68のすべての索性を利用した実験を示し、「索性選択」は省いた索性以外の全索性を利用した実験を示す。

10

【0251】

図8における正解率は、それぞれの実験データに対して、編集距離が1の日本語異表記対であるのか、編集距離が1の日本語異表記対でないのかを、正しく判定した割合である。図8のF値は、それぞれの実験データに対して、編集距離が1の日本語異表記対を抽出する場合のF値である。10CV、OCに対して、全索性を利用したSVMの正解率、F値ともにベースラインの手法よりも高いことがわかる。編集距離1の日本語用語対から日本語異表記対を抽出する場合のF値は、ベースラインと比較して全索性を利用したSVMの方が、10CVでは0.433高く、オープンデータでは0.460高かった。ベースラインの結果より、編集距離が1の日本語用語対から日本語異表記対を抽出することは難しいといえるが、本報告で提案している種々の索性と機械学習を用いた手法は、ベースライン手法よりも多くの日本語異表記対が抽出できることがわかる。

20

【0252】

次に、上述したブートストラップ法を用いて索性が有効であるかどうかの検討をした結果を図9に示す。なお、S1からS68の索性を図10に示す。図9において、省いた索性は、ブートストラップ法により有意かどうかの判定が行われる索性である。全索性は本実験で扱った索性による手法であり、索性選択は省いた索性を全索性から省いた索性による手法である。それぞれの値は、全索性が索性選択よりF値が高かった回数あるいは、索性選択が全索性よりF値が高かった回数である。この実験では、全索性が索性選択よりF値が高い回数が9500回(95%)以上あれば、省いた索性は精度向上に役立っているということになり、省いた索性は有効であるといえる。全索性が索性選択よりF値が高い回数が9500回(95%)以上であった索性は、10CVの場合はS47、S55、S58、S67であり、OCの場合はS52、S54、S55、S58、S67であった。10CVとOCの両方で、全索性が索性選択よりF値が高い回数が9500回(95%)以上であった索性は、S55、S58、S67であった。この結果からS55、S58、S67の索性は、どのような編集距離の小さい日本語異表記対を抽出するデータにも、有効である索性といえる。S47とS52の索性はそれぞれの実験で使われたデータには有効である索性といえるが、編集距離の小さい日本語用語対のデータが変われば、有効でなくなる可能性がある索性といえる。そのためS47とS52は、どのような編集距離の小さい日本語異表記対を抽出するデータにも、有効であるとはいえない。

30

40

【0253】

次に、本異表記取得装置1の提案手法(以下、単に提案手法とも言う)が、編集距離が1の日本語異表記対を抽出できたのかを、種々の同義語辞書を用いて比較を行った結果について説明する。種々の同義語辞書は、EDR辞書、日本語ワードネット辞書、JUMAN辞書である。編集距離が1の日本語用語対は、EDR辞書には21224779組、日本語ワードネット辞書には890616組、JUMAN辞書には23348組あることがわかった。EDR辞書には人名に関する単語がある。本実験では人名は同義語でないと判断し、取り除いた。その結果、EDR辞書に含まれている編集距離が1の日本語用語対は933037組であった。JUMAN辞書は同じ代表表記をもつ単語対を日本語用語対として扱った。

50

【0254】

提案手法を用い、大規模類似語リストから編集距離が1の日本語異表記対と分類された用語対が、種々の辞書にどの程度の割合で含まれているかの検討結果を図11に示す。以降は大規模類似語リストにおける編集距離が1の日本語用語対すべてを、日本語用語対データベースとし、用語対DBと表記する。さらに、用語対DBにおいて、日本語異表記対であると提案手法が分類した日本語用語対すべてを、日本語異表記対データベースとし、異表記DBと表記する。用語対DBにおいて、日本語異表記対でない提案手法が分類した非日本語用語対すべてを、非日本語異表記対データベースとし、非異表記DBと表記する。EDR辞書は20.45%、日本語ワードネットは1.71%、JUMAN辞書は6.52%の割合で異表記DBの日本語異表記対が含まれていた。どの辞書においても、異表記DBの日本語用語対を含んでいる割合は高くない。これらの結果より、本明細書で記載した異表記取得装置により得られた異表記と既存辞書は重なりが小さいので、異表記取得装置により、既存辞書に対して多くの異表記を追加できることが分かる。また、例えば、EDR辞書では、約2割のカバー率であるが、相当な程度のカバー率である、と言える。

10

【0255】

また、種々の辞書に含まれる編集距離が1の日本語用語対を、提案手法により編集距離が1の日本語異表記対であるか、編集距離が1の日本語異表記対でないか分類した。SVMの分類における正解率を図12に示す。

【0256】

また、種々の辞書と用語対DBにおいて、編集距離が1の日本語異表記対であると分類された日本語用語対と、分類されなかった日本語用語対をそれぞれランダムに、5組ずつ取り出した結果を図13に示す。学習データはオープンデータとクローズデータを組み合わせたデータとし、テストデータ用語対DB、種々の辞書のそれぞれでOCにより、用語対DBと種々の辞書を分類した。

20

【0257】

図12において、日本語ワードネットにおいて分類の正解率が低かったのは、図13のように、日本語異表記対ではなく、日本語類義語対が多く含まれているからである。図12に示すように、JUMAN辞書の場合は8割という高い正解率で分類できている。また、JUMAN辞書には、日本語異表記対でないものが含まれるという問題が少なく、また、本提案手法により適切に異表記を抽出できるために、8割という高い正解率を達成できたものと考えられる。

30

【0258】

次に、編集距離が1の日本語異表記対抽出の評価について述べる。SVMは識別関数の出力値(機械学習部16が出力するスコア)が正か負かによって、データを分類することも可能であるが、ここでは、識別関数の出力値が正か負かによって、データを分類するのではなく、任意の値(閾値)によって正か負のデータを分類し、編集距離が1の日本語異表記対抽出の評価を行う。つまり、閾値判断手段172が、機械学習部16が取得したスコアが閾値格納手段171に格納されている閾値以上または閾値より大きいか否かを判断するものとする。正のデータを編集距離が1の日本語異表記対であると分類し、負のデータを編集距離が1の日本語異表記対でない日本語用語対と分類する手法では、精度が100%ではないため、誤って編集距離が1の日本語異表記対でない日本語用語対を、日本語異表記対であると判断し、抽出することがある。そのため、少量であっても確実に抽出を行いたい場合は、閾値を高く設定することで、日本語異表記対を確実に抽出できる。また、誤ったデータが含まれていても、網羅的に抽出を行いたい場合は、閾値を低く設定することで、可能となる。図14に、閾値の評価基準を示す。図14に示すように、閾値を-0.2に設定することで、F値0.9323と最も高い値を得られることがわかった。また、再現率と適合率の比率を図15に示す。図15によれば、再現率が高くしようとすると、カバー率を上げなくてはならなくなり、適合率は低くなり、再現率は低くなる。

40

50

(実験結果 2)

【 0 2 5 9 】

第 2 番目の実験において、正例 (例えば、「スパゲティ」と「スパゲッティ」との対) 7 4 5 個、負例 (正例に該当しない対) 1 3 , 4 4 0 個を持つ学習データから、正例 7 2 5 個、負例 1 3 , 4 6 0 個のテストデータの抽出が行なわれ、その F 値は、0 . 9 3 であった。なお、実験結果 2 において、実験結果 1 で利用した素性や学習データが完全に一致するものではないが、本提案手法の有効性を示すために足りる、素性や学習データの重複がある。

【 0 2 6 0 】

すべてを正例と判断する、即ちどんなものでも正例とするベースラインの方法であると、F 値は 0 . 0 9 7 2 程度であった。異表記かどうかを判定する対象の用語対において、編集箇所の文字また、編集箇所の文字の周辺の文字だけの素性を用いる従来の方法でも、F 値は 0 . 8 5 であった。つまり、提案手法のように、多数の素性 (ここでは、6 8) を用いた方法の効果は顕著であることが分かる。

【 0 2 6 1 】

また、既存の異表記辞書に基づく素性、また、スタッキング手法に基づく素性 (上記の辞書関連素性) を利用しなかった方法よりも、これらの方法を利用した方が有意に F 値が高いことも確かめており、これらの手法の有効性も確認している。

【 0 2 6 2 】

また、ルールベース的手法として、編集箇所の文字の字種が漢数字かアラビア数字であること、または、同じアルファベットであること、また、既存の異表記辞書を利用することで異表記と判定できるものを、異表記と決定的に推定する方法も試した。この場合の F 値は、0 . 4 2 0 2 であり、ルールベース的手法でなく教師あり機械学習を利用する方が良いことがわかる。

【 0 2 6 3 】

正しい異表記の対の差分データから、異表記になりやすい差分パターンを学習し、ある用語 A に対して、異表記の候補 B を上記差分パターンより生成し、用語 A と用語 B が異表記の対であるかを判定する操作を利用することにより、取得できる異表記が格段に増えるという効果がある。かかる操作については、実施の形態 2 で説明する。

(実験結果 3)

【 0 2 6 4 】

第 3 番目の実験において、1 0 万語の単語とそれの類似する 1 0 0 語の単語を用いた。1 0 万語の単語とそれの類似する 1 0 0 語の単語のすべての対のうち、1 文字のみ変化している用語対は 1 7 0 万個あった。なお、実験結果 3 において、実験結果 1 で利用した素性や学習データが完全に一致するものではないが、本提案手法の有効性を示すために足りる、素性や学習データの重複がある。そして、異表記取得装置 1 の技術を利用して、そこから 7 万対の異表記を取り出せる。以下に構築できる異表記の例を示す。

? B u s i n e s s W e e k B u s i n e s s W e e k

? J A V A S c r i p t J A V A S c r i p t

? 書いてた頃 書いていた頃

? アイリッシュトラッド アイリッシュ・トラッド

? 自サーバ 自サーバー

? でない場合 出ない場合

? WWWサーバ上 WWWサーバー上

? 日光彫 日光彫り

? 隣同士 隣り同士

【 0 2 6 5 】

なお、EDR (E l e c t r i c D i c t i o n a r y R e s e a r c h) 電子化辞書に含まれる差分が 1 文字の異表記は 2 4 , 1 8 5 語である。また、日本語 Word N e t に含まれる差分が 1 文字の異表記のようなものは 8 2 , 2 7 0 語ある。ただし、日本

10

20

30

40

50

語 WordNet には、異表記でないもの（類義語）も多く含まれており、適切に異表記を取り出すことが困難である。さらに、JUMAN の辞書に含まれる差分が 1 文字の異表記は 23, 348 語である。これらと比較しても本提案手法の技術の有効性がわかる。また、「JAVA」は登録商標です。

(実験結果 4)

【0266】

第 4 番目の実験において、上記したベースライン手法（上記のルール 1 からルール 4 を適用した方法）による精度を算出する。ベースライン手法では、有意差が高かった素性（S55、S58、S67）が yes と判定されたものを正例、すべて no と判定されたものを負例として F 値を求める。

10

【0267】

図 16 は、ベースライン手法で、オープンデータとクローズドデータの全部を用いて 10 分割クロスバリデーションによる実験をおこなった場合の結果である。図 16 において、「0」は負例、「1」は正例である。また、図 16 において、最も左側の列の「0」「1」は、正しい分類を示す。最も上の第一行の「0」「1」は、実験対象の手法（図 16 では、ベースライン手法）での出力結果を示す。つまり、正しい分類が「0」であり実験結果が「0」であったデータの数が 26892、正しい分類が「0」であり実験結果が「1」であったデータの数が 1018、正しい分類が「1」であり実験結果が「0」であったデータの数が 8、正しい分類が「1」であり実験結果が「1」であったデータの数が 452 であった。また、負例（「0」）の再現率は 99.97%、適合率は 96.35% であった。また、正例（「1」）の再現率は 30.75%、適合率は 98.26% であった。また、すべてのデータの再現率は 96.38%、適合率は 96.38% であった。さらに、「総数」は、実験データの数である。以上の再現率、適合率を、数式 10 に代入して、算出した負例の F 値は 0.9813、正例の F 値は 0.4684 であった。なお、図 17 から図 21 の各データの意味は、図 16 と同様であるので説明を省略する。

20

【0268】

図 17 は、ベースライン手法で、クローズドデータのみを用いて 10 分割クロスバリデーションによる実験をおこなった場合の結果である。図 17 において、負例の F 値は 0.9814、正例の F 値は 0.4833 であった。

30

【0269】

図 18 は、ベースライン手法で、オープンクローズを用いた場合の結果である。図 18 において、負例の F 値は 0.9812、正例の F 値は 0.4529 であった。

【0270】

実験結果 4 において、ベースライン手法は、正例の F 値が、提案手法における F 値（例えば、実験結果 1 の 0.912）と比較して極めて小さく、提案手法の有効性が極めて高い、と言える。

(実験結果 5)

【0271】

第 5 番目の実験において、すべてを正例としたベースライン手法の場合による精度を算出した。すべて正例としたベースライン手法の場合、再現率は「100%」、適合率は「0.0525%」であった。そして、かかる再現率および適合率を、数式 10 に代入し、算出された、正例（「1」）の F 値は「0.0998」であった。すべてを正例としたベースライン手法の正例の F 値は、提案手法における F 値と比較して極めて小さく、提案手法の有効性が極めて高い、と言える。なお、提案手法において、正解率「99.12%」、再現率「99.07%」、適合率「92.29%」、F 値「0.912」を得ている。なお、本実験で利用した素性や学習データは、提案手法の評価において利用した素性や学習データと完全に一致するものではないが、本提案手法の有効性を示すために足りる、素性や学習データの重複がある。

40

(応用例)

【0272】

50

以下、異表記取得装置 1 の応用例について説明する。応用例とは、異表記取得装置 1 を組み込んだ情報検索装置である。情報検索装置は、異表記取得装置 1 と検索部とを具備する。つまり、受付部 14 は、キーワード (KW1) を受け付ける。そして、異表記取得装置 1 は、受け付けた KW1 の異表記の用語 (KW2) を取得する。そして、検索部は、KW1 + KW2 (+ は OR) の検索式により、情報検索を行う。なお、情報検索の検索対象は問わないことは言うまでもない。また、検索部は、いわゆる Web の検索エンジンを起動するだけの処理でも良い。

【0273】

本情報検索装置を利用して、情報をキーワード検索する際に、ユーザが、「スパゲティ」と入力した場合に、情報検索装置は、「スパゲティ」の異表記である「スパゲッティ」を取得する。そして、情報検索装置は、これらの「スパゲティ」と「スパゲッティ」の双方をキーワードとして情報検索する。その結果、「スパゲティ」と「スパゲッティ」のいずれの表現が為されている情報もヒットするので、検索漏れの少ない情報検索が実現できる。

10

【0274】

特に情報検索装置は、検索漏れが許されない特許情報の検索に大きな効果をもたらす。例えば、情報検索装置が特許検索システムにおいて利用されることを考える。特許の明細書や特許請求の範囲や要約書等の特許の書類には、例えば、「コンピュータ」も「コンピューター」も存在するので、キーワードとして「コンピュータ+コンピューター」を入力しなければ、検索漏れが生じる。従って、検索者は検索時には細心の注意を払って検索しようとするキーワードの異表記を考える必要があった。「デジタル」と「ディジタル」となど、同義語であるにも拘わらず、異表記の文言は特に特許公報においては多い。しかしながら、本情報検索装置を採用することによってこのような配慮をすることなく、検索漏れのない特許情報の検索が可能となる。

20

【0275】

以上、本実施の形態によれば、用語対の分野を問わず、精度の高い異表記の用語対の抽出が可能となる。

【0276】

なお、本実施の形態によれば、主として、編集距離が 1 の用語対について、異表記の用語対であるか否かの判断手法について説明した。しかし、上述したとおり、異表記取得装置 1 は、編集距離が 2 の用語対についても、異表記の用語対であるか否かを判断できる。

30

【0277】

つまり、素性取得部 15 の差分文字取得手段 151 は、編集距離が 2 つの用語対について、2 つの差分文字の組を、それぞれ取得する。例えば、以下の 3 つの具体的な用語対を考える。(1)「できる」「出来る」(2)「理解できる」「できる」(3)「IX (ローマ数字の 9)」「9」を考える。かかる場合、差分文字取得手段 151 は、(1)の用語対について、「で」「出」と「き」「来」の 2 組の差分文字の組を取得する。また、差分文字取得手段 151 は、(2)の用語対について、「理」「」と「解」「」(「」は NULL である)の 2 組の差分文字の組を取得する。また、差分文字取得手段 151 は、(3)の用語対について、「I」「9」と「X」「」の 2 組の差分文字の組、または「I」「」と「X」「9」の 2 組の差分文字の組を取得する。

40

【0278】

そして、素性取得手段 152 は、差分文字取得手段 151 が取得した 2 つの差分文字を、独立に対象として、字種関連素性、辞書関連素性、類似度素性のうちの一以上を含む複数の素性を、2 組取得する。つまり、(1)の用語対について、素性取得手段 152 は、「で」「出」と「き」「来」の 2 組の差分文字の組のそれぞれを対象に素性の抽出を行い、それぞれ差分文字から抽出した素性は、別のもと考え、2 種類のテストデータを作成する。素性取得手段 152 は、例えば、用語対が有する 2 つの用語の編集箇所の字種が異なり、かつ、2 つの用語の編集箇所が同じ値の数字であるか否かを示す字種関連素性について、「で」「出」の編集箇所が同じ値の数字でないと判断し、当該字種関連素性「0」

50

を取得する。また、素性取得手段 152 は、例えば、「で」「出」について、2つの用語の読みが一致するか否かを示す辞書関連素性「1」を取得する。素性取得手段 152 は、用語辞書 13 から「出」の読み「で」を取得し、「で」と「出」の読みが一致すると判断する。また、素性取得手段 152 は、例えば、差分文字「で」「出」に対して、差分文字（編集箇所）の前後の文字という素性について、前の文字の素性「」（なし）、後の文字の素性「き」と「来」を取得する。また素性取得手段 152 は、例えば、差分文字「き」「来」に対して、差分文字の前後の文字という素性について、前の文字の素性「出」と「で」、後の文字の素性「る」を取得する。かかる処理により、別の差分文字も素性に含めることとなる。

【0279】

また、(2)の用語対について、素性取得手段 152 は、(1)と同様に、「理」「」と「解」「」の2組の差分文字の組のそれぞれを対象に素性の抽出を行い、それぞれ差分文字から抽出した素性は、別のものと考え、2種類のテストデータを作成する。さらに、(3)の用語対について、素性取得手段 152 は、(1)(2)と同様に、例えば、「I」「9」と「X」「」の2組の差分文字の組のそれぞれを対象に素性の抽出を行い、それぞれ差分文字から抽出した素性は、別のものと考え、2種類のテストデータを作成する。

【0280】

次に、機械学習部 16 は、(1)(2)(3)について、2種類のテストデータをそれぞれ、異表記の用語対であるか否かを判定する。そして、判定の結果、例えば、2種類のテストデータともに異表記の用語対であると判定された場合、元の用語対（例えば、「できる」「出来る」）は、異表記の用語対であるとして、出力部 17 は、判断結果を出力する。なお、出力部 17 は、上述したように、2種類のテストデータに対する2つのスコアのうちのスコアが0に近い方のスコアを採用して、採用したスコアが正の場合は正例（異表記の用語）、負の場合は負例（異表記の用語でない）と判断しても良いし、スコアの絶対値が大きい方のスコアを採用して、採用したスコアが正の場合は正例（異表記の用語）、負の場合は負例（異表記の用語でない）と判断しても良いし、2つのスコアのうちの、小さい方のスコアを取得し、当該小さい方のスコアが正の場合は正例（異表記の用語）、負の場合は負例（異表記の用語でない）と判断しても良い。

【0281】

また、(3)の2組の差分文字の組（例えば、「I」「9」と「X」「」、または「I」「」と「X」「9」）、つまり2つの問題（問題1、問題2）ができる場合、それぞれの差分文字を対象に素性の抽出を行い、それぞれ差分文字から抽出した素性は、別のものと考え、4種類のテストデータを作成する。そして、2つの問題ごとに、算出したスコアが0に近い方を取得し、問題ごとのスコアのうちの、絶対値が高いスコアを当該問題のスコアとし、スコアが正の場合は正例、負の場合は負例と判断しても良い。また、例えば、用語対が(3)「IX」「9」である場合、問題「I」「9」と「X」「」、および「I」「」と「X」ができる。そして、機械学習部 16 は、「I」「9」と「X」「」のスコアの小さい方を取得し、また、「I」「」と「X」「9」のスコアの小さい方を取得し、2つの取得されたスコアのうちの、値が大きい方を「IX」「9」の用語対におけるスコアとする。そして、機械学習部 16 は、当該スコアが正の場合は正例、負の場合は負例と判断しても良い。なお、例えば、機械学習部 16 は、「I」「9」と「X」「」のスコアが0に近い方を取得し、また、「I」「」と「X」「9」のスコアが0に近い方を取得し、2つの取得されたスコアのうちの、絶対値が大きい方を「IX」「9」の用語対におけるスコアとしても良い。そして、機械学習部 16 は、当該スコアが正の場合は正例、負の場合は負例と判断しても良い。

【0282】

また、本実施の形態において、編集距離が3以上の用語対についても、編集距離が2つの用語対と同様に、3以上のテストデータを作成し、3以上のテストデータの判断結果を用いて、元の用語対が異表記の用語対であるか否かを判定しても良い。かかる場合、例えば、3以上の差分文字のうちの1文字や2文字などを素性として用いるなど、新しい素性

10

20

30

40

50

を機械学習手法に導入しても良い。

【0283】

また、本実施の形態において、異表記取得装置1は、例えば、「あなた」「あんた」という日本語の用語対が異表記の用語対であると判断できたが、日本語以外の言語（例えば、英語）の用語対（例えば、「colour」「color」）も、異表記の用語対であると判断できる。

【0284】

また、本実施の形態において、異表記取得装置1は、用語対を構成する2つの用語の編集箇所の文字が2文字以上である場合、1文字ずつの対応とせずに、編集箇所をまとめて、当該まとめた文字列をそのまま機械学習しても良い。つまり、用語対「123組」「百二十三組」について、編集箇所を「123」「百二十三」とまとめて、処理しても良い。用語対「123組」「百二十三組」に対して、例えば、S1「一つ目の表記の編集箇所」「123」、S2「二つ目の表記の編集箇所」「百二十三」、S3「編集箇所の前方の1文字」"（なし）、S4「編集箇所の後方の1文字」「組」、S55「編集箇所が両方とも数字の場合であり、同じ値か違う値かどうか」「1」（同じ値）、S56「日本語用語対の編集箇所が両方ともひらがなの場合であり、同じ音声か違う音声かどうか」「0」等が得られる。そして、用語対「123組」「百二十三組」に対する学習データが構成され、学習データ格納部12に蓄積されて、利用されても良い。また、異表記取得装置1の素性取得部15は、編集箇所の文字が2文字以上の用語対のテストデータに対して、編集箇所をまとめて処理し、例えば、上述した68の素性を取得し、機械学習部16が、テストデータが異表記対か否かを判断しても良い。

【0285】

さらに、本実施の形態における処理は、ソフトウェアで実現しても良い。そして、このソフトウェアをソフトウェアダウンロード等により配布しても良い。また、このソフトウェアをCD-ROMなどの記録媒体に記録して流布しても良い。なお、このことは、本明細書における他の実施の形態においても該当する。なお、本実施の形態における情報処理装置を実現するソフトウェアは、以下のようなプログラムである。つまり、このプログラムは、記憶媒体に、編集距離が1以上の用語対、および、用語対の異なる文字である編集箇所の字種に関する素性である字種関連素性、用語辞書を用いて取得された素性である辞書関連素性、前記用語対を構成する2つの用語の類似度を示す素性である類似度素性のうちの1以上の素性を含む複数の素性と、前記用語対が異表記の用語対であることを示す情報である正負情報とを対応付けた学習データを2以上格納しており、コンピュータを、前記記憶媒体の用語対ごとに、字種関連素性、辞書関連素性、類似度素性のうちの1以上を含む複数の素性を取得する素性取得部と、前記各用語対に対して、前記記憶媒体の2以上の学習データと、前記素性取得部が取得した複数の素性とを用いて、教師あり機械学習法により、前記記憶媒体の各用語対が異表記の用語対であるか否かを判断する機械学習部と、前記機械学習部における判断結果を出力する出力部として機能させるためのプログラム、である。

（実施の形態2）

【0286】

本実施の形態において、置き換え文字対を保持し、置き換え文字対を用いて、用語から用語対を生成し、その用語対に対して、機械学習により異表記用語対を生成する異表記取得装置2について説明する。異表記取得装置2は、異表記取得装置1の機能に加えて、パターンを使った異表記用語対の生成機能を有する機能を有する。

【0287】

図19は、本実施の形態における異表記取得装置2のブロック図である。

異表記取得装置2は、用語対格納部11、異表記用語対格納部21、学習データ格納部12、用語辞書13、異表記パターン格納部22、受付部23、編集箇所取得部24、異表記パターン取得部25、異表記パターン蓄積部26、用語対生成部27、素性取得部15、機械学習部16、出力部17を備える。

【0288】

異表記用語対格納部21は、編集距離が1の異表記の用語対を1以上格納し得る。

【0289】

異表記パターン格納部22は、異表記のパターンを示す第一文字列と第二文字列とを対に有する異表記パターンを1以上格納し得る。

【0290】

受付部23は、ユーザからの入力を受け付ける。また、受付部23は、1以上の用語を受け付ける。この用語とは、用語対を生成する元となる用語である。受付部23が用語を受け付けるのは、ユーザからの入力でも良いし、記憶媒体からの読み込みや、通信手段を用いた受信でも良い。

【0291】

編集箇所取得部24は、異表記用語対格納部21に格納されている1以上の異表記の用語対の編集箇所を取得する。

【0292】

異表記パターン取得部25は、編集箇所取得部24が取得した編集箇所から、第一文字列と第二文字列とを対に有する異表記パターンを取得する。異表記パターン取得部25は、例えば、用語対「2番目」「二番目」から第一文字列「2」と第二文字列「二」とを対に有する異表記パターン「2」「二」を取得する。また、異表記パターン取得部25は、例えば、用語対「自サーバ」「自サーバー」から異表記パターン「del」「ー」を取得する。

【0293】

異表記パターン蓄積部26は、異表記パターン取得部25が取得した異表記パターンを、異表記パターン格納部22に蓄積する。

【0294】

用語対生成部27は、受付部23が受け付けた1以上の各用語に対して、異表記パターン格納部22の1以上の各異表記パターンを適用し、1以上の用語を生成し、1以上の各用語と生成した用語とを有する1以上の異表記の候補の用語対である異表記候補用語対を生成する。

【0295】

異表記用語対格納部21、異表記パターン格納部22は、不揮発性の記録媒体が好適であるが、揮発性の記録媒体でも実現可能である。

【0296】

異表記用語対格納部21に異表記用語対が記憶される過程は問わない。

【0297】

編集箇所取得部24、異表記パターン取得部25、異表記パターン蓄積部26、および用語対生成部27は、通常、MPUやメモリ等から実現され得る。編集箇所取得部24等の処理手順は、通常、ソフトウェアで実現され、当該ソフトウェアはROM等の記録媒体に記録されている。但し、ハードウェア(専用回路)で実現しても良い。

【0298】

次に、異表記取得装置2の動作について図20のフローチャートを用いて説明する。図20のフローチャートにおいて、異表記パターンを蓄積する処理、用語対を生成する処理について説明する。異表記取得装置2の動作について、異表記取得装置1の動作と同じである異表記用語対の判断処理、および判断結果の出力処理については、説明しない。

【0299】

(ステップS2001) 受付部23は、異表記パターンの生成指示を受け付けたか否かを判断する。異表記パターンの生成指示を受け付ければステップS2002に行き、受け付けなければステップS2009に行く。

【0300】

(ステップS2002) 編集箇所取得部24は、カウンタ*i*に1を代入する。

【0301】

10

20

30

40

50

(ステップ S 2 0 0 3) 編集箇所取得部 2 4 は、異表記パターン格納部 2 2 に i 番目の用語対が存在するか否かを判断する。i 番目の用語対が存在すればステップ S 2 0 0 4 に行き、i 番目の用語対が存在しなければステップ S 2 0 0 1 に戻る。

【0302】

(ステップ S 2 0 0 4) 編集箇所取得部 2 4 は、i 番目の用語対の差分文字(編集箇所)を取得する。

【0303】

(ステップ S 2 0 0 5) 異表記パターン取得部 2 5 は、ステップ S 2 0 0 4 で取得した差分文字(編集箇所)から、異表記パターンを構成する。

【0304】

(ステップ S 2 0 0 6) 異表記パターン蓄積部 2 6 は、ステップ S 2 0 0 5 で取得された異表記パターンが、異表記パターン格納部 2 2 に存在するか否かを判断する。存在すればステップ S 2 0 0 7 に行き、存在しなければステップ S 2 0 0 8 に行く。

【0305】

(ステップ S 2 0 0 7) 異表記パターン蓄積部 2 6 は、ステップ S 2 0 0 5 で取得された異表記パターンを、異表記パターン格納部 2 2 に蓄積する。

【0306】

(ステップ S 2 0 0 8) 編集箇所取得部 2 4 は、カウンタ i を 1、インクリメントする。ステップ S 2 0 0 3 に戻る。

【0307】

(ステップ S 2 0 0 9) 受付部 2 3 は、用語を受け付けたか否かを判断する。用語を受け付ければステップ S 2 0 1 0 に行き、受け付けなければステップ S 2 0 0 1 に戻る。

【0308】

(ステップ S 2 0 1 0) 用語対生成部 2 7 は、カウンタ i に 1 を代入する。

【0309】

(ステップ S 2 0 1 1) 用語対生成部 2 7 は、i 番目の異表記パターンが、異表記パターン格納部 2 2 に存在するか否かを判断する。存在すればステップ S 2 0 1 2 に行き、存在しなければ処理を終了する。

【0310】

(ステップ S 2 0 1 2) 用語対生成部 2 7 は、ステップ S 2 0 0 9 で受け付けられた用語が、i 番目の異表記パターンに合致するか否かを判断する。合致すればステップ S 2 0 1 3 に行き、合致しなければステップ S 2 0 1 6 に行く。なお、用語「WWWサーバ」に対して、異表記パターン「2」「二」は合致しない。異表記パターンは両方とも文字であり、当該いずれの文字も用語「WWWサーバ」が含まないからである。また、用語「WWWサーバ」に対して、異表記パターン「del」「ー」は合致する。異表記パターンに「del」が含まれる場合は、すべての用語が異表記パターンに合致することとなる。

【0311】

(ステップ S 2 0 1 3) 用語対生成部 2 7 は、ステップ S 2 0 0 9 で受け付けられた用語に対して、i 番目の異表記パターンを適用し、1 以上の異表記の用語を取得する。用語が「アイトラッキング」であり、i 番目の異表記パターンが「del」「・」である場合、用語対生成部 2 7 は、用語「アイトラッキング」に異表記パターン「del」「・」を適用し、「・」を各文字間に挿入し、7 つの異表記の用語「ア・イトラッキング」「アイ・トラッキング」「アイト・ラッキング」「アイトラ・ッキング」「アイトラッ・キング」「アイトラッキ・ング」「アイトラッキン・グ」を生成する。また、用語が「一番目」であり、i 番目の異表記パターンが「ー」「1」である場合、用語対生成部 2 7 は、用語「一番目」に異表記パターン「ー」「1」を適用し、「1 番目」を生成する。

【0312】

(ステップ S 2 0 1 4) 用語対生成部 2 7 は、ステップ S 2 0 0 9 で受け付けられた用語と、ステップ S 2 0 1 3 で生成した 1 以上の異表記の用語を用いて、1 以上の用語対を生成する。例えば、用語が「アイトラッキング」であり、i 番目の異表記パターンが「de

10

20

30

40

50

「・」である場合、用語対生成部 27 は、用語対「アイトラッキング」「ア・イトラッキング」、「アイトラッキング」「アイ・トラッキング」、「アイトラッキング」、「アイトラッキング」「イト・ラッキング」、「アイトラッキング」「イトラ・ッキング」、「アイトラッキング」「イトラッ・キング」、「アイトラッキング」「イトラッキ・ング」、「アイトラッキング」「イトラッキン・グ」の 7 つの用語対を生成する。また、用語が「一番目」であり、i 番目の異表記パターンが「一」「1」である場合、用語対生成部 27 は、用語対「一番目」「1 番目」を生成する。

【0313】

(ステップ S 2015) 用語対生成部 27 は、ステップ S 2013 で生成した 1 以上の用語対を、用語対格納部 11 に蓄積する。

10

【0314】

(ステップ S 2016) 用語対生成部 27 は、カウンタ i を 1、インクリメントする。ステップ S 2011 に戻る。

【0315】

以下、本実施の形態における異表記取得装置 2 の具体的な動作について説明する。

【0316】

異表記パターン取得部 25 が取得し、異表記パターン蓄積部 26 が異表記パターン格納部 22 に蓄積した異表記パターンの例を、図 21 に示す。図 21 において、「del」は、もう一方のパターン文字を削除することを示す。つまり、異表記パターン取得部 25 は、del のもう一方のパターン文字について、すべての大規模類似語リストに用いた用語を対象とし、用語対を生成する。

20

【0317】

かかる状況において、上述したように、用語対生成部 27 は、受け付けられた用語「アイトラッキング」が、1 番目の異表記パターン「del」「・」に合致する、と判断する。そして、用語が「アイトラッキング」が入力された場合、1 番目の異表記パターン「del」「・」が適用され、用語対生成部 27 は、用語「アイトラッキング」は、「・」を各文字間に挿入し、7 つの異表記の用語「ア・イトラッキング」「アイ・トラッキング」「イト・ラッキング」「イトラ・ッキング」「イトラッ・キング」「イトラッキ・ング」「イトラッキン・グ」を生成する。次に、用語対生成部 27 は、用語対「アイトラッキング」「ア・イトラッキング」、「アイトラッキング」「アイ・トラッキング」、「アイトラッキング」、「アイトラッキング」「イト・ラッキング」、「アイトラッキング」「イトラ・ッキング」、「アイトラッキング」「イトラッ・キング」、「アイトラッキング」「イトラッキ・ング」、「アイトラッキング」「イトラッキン・グ」の 7 つの用語対を生成する。そして、異表記パターン蓄積部 26 は、7 つの用語対を異表記パターン格納部 22 に蓄積する。

30

【0318】

次に、用語対生成部 27 は、受け付けられた用語「アイトラッキング」が、2 番目の異表記パターン「del」「-」に合致する、と判断する。次に、用語「アイトラッキング」に対して、2 番目の異表記パターン「del」「-」が適用され、用語対生成部 27 は、「ア・イトラッキング」「アイ・トラッキング」「イト・ラッキング」「イトラ・ッキング」「イトラッ・キング」「イトラッキ・ング」「イトラッキン・グ」を生成する。次に、用語対生成部 27 は、用語対「アイトラッキング」「ア・イトラッキング」、「アイトラッキング」「アイ・トラッキング」、「アイトラッキング」「イト・ラッキング」、「アイトラッキング」「イトラ・ッキング」、「アイトラッキング」「イトラッ・キング」、「アイトラッキング」「イトラッキ・ング」、「アイトラッキング」「イトラッキン・グ」の 7 つの用語対を生成する。そして、異表記パターン蓄積部 26 は、7 つの用語対を異表記パターン格納部 22 に蓄積する。

40

【0319】

次に、用語対生成部 27 は、受け付けられた用語「アイトラッキング」が、3 番目の異表記パターン「del」「い」に合致する、と判断する。そして、次に、用語「イトラッ

50

キング」に対して、3番目の異表記パターン「del」「い」が適用され、用語対生成部27は、「アイトラッキング」「アイイトラッキング」「アイトいラッキング」「アイトラッキング」「アイトラッキング」「アイトラッキング」「アイトラッキング」「アイトラッキング」「アイトラッキング」「アイトラッキング」「アイトラッキング」「アイトラッキング」「アイトラッキング」「アイトラッキング」「アイトラッキング」「アイトラッキング」「アイトラッキング」「アイトラッキング」「アイトラッキング」の7つの用語対を生成する。そして、異表記パターン蓄積部26は、7つの用語対を異表記パターン格納部22に蓄積する。

【0320】

次に、用語対生成部27は、受け付けられた用語「アイトラッキング」が、4番目の異表記パターン「-」「1」を構成する文字を含まないので、この異表記パターンに合致しない、と判断する。

【0321】

次に、用語対生成部27は、受け付けられた用語「アイトラッキング」が、5番目の異表記パターン「イ」「ィ」を構成する文字「イ」を含むので、この異表記パターンに合致する、と判断する。そして、次に、用語「アイトラッキング」に対して、5番目の異表記パターン「イ」「ィ」が適用され、用語対生成部27は、「アイトラッキング」を生成する。次に、用語対生成部27は、用語対「アイトラッキング」「アイトラッキング」の1つの用語対を生成する。そして、異表記パターン蓄積部26は、1つの用語対を異表記パターン格納部22に蓄積する。

【0322】

次に、同様に、用語対生成部27は、6番目以降の異表記パターンを適用して、処理していく。

【0323】

そして、用語対生成部27は、新たな用語対を用語対格納部11に蓄積する。

【0324】

以上、本実施の形態によれば、異表記の用語対の候補を自動生成できる。また、本実施の形態によれば、異表記の用語対の候補を自動生成するための異表記パターンを自動的に取得できる。

【0325】

なお、本実施の形態における情報処理装置を実現するソフトウェアは、以下のようなプログラムである。つまり、このプログラムは、記憶媒体に、用語対の異なる文字である編集箇所の字種に関する素性である字種関連素性、用語辞書を用いて取得された素性である辞書関連素性、前記用語対を構成する2つの用語の類似度を示す素性である類似度素性のうちの一以上の素性を含む複数の素性と、前記用語対が異表記の用語対であるかを示す情報である正負情報とを対応付けた学習データを2以上格納しており、記憶媒体に、異表記のパターンを示す第一文字列と第二文字列とを対に有する異表記パターンを1以上格納しており、コンピュータを、1以上の用語を受け付ける受付部と、前記受付部が受け付けた1以上の各用語に対して、前記記憶媒体の1以上の各異表記パターンを適用し、1以上の用語を生成し、前記1以上の各用語と前記生成した用語とを有する1以上の異表記の候補の用語対である異表記候補用語対を生成する用語対生成部と、前記用語対生成部が生成した1以上の異表記候補用語対ごとに、字種関連素性、辞書関連素性、類似度素性のうちの一以上の素性を含む複数の素性を取得する素性取得部と、前記用語対生成部が生成した各異表記候補用語対に対して、前記記憶媒体の2以上の学習データと、前記素性取得部が取得した複数の素性とを用いて、教師あり機械学習法により、前記用語対格納部の各異表記候補用語対が異表記の用語対であるか否かを判断する機械学習部と、前記機械学習部における判断結果を出力する出力部として機能させるためのプログラム、である。

また、図22は、本明細書で述べたプログラムを実行して、上述した実施の形態の異表記取得装置等を実現するコンピュータの外観を示す。上述の実施の形態は、コンピュータ

10

20

30

40

50

ハードウェア及びその上で実行されるコンピュータプログラムで実現され得る。図 22 は、このコンピュータシステム 340 の概観図であり、図 23 は、コンピュータシステム 340 のブロック図である。

【0326】

図 22 において、コンピュータシステム 340 は、FD ドライブ、CD-ROM ドライブを含むコンピュータ 341 と、キーボード 342 と、マウス 343 と、モニタ 344 とを含む。

【0327】

図 23 において、コンピュータ 341 は、FD ドライブ 3411、CD-ROM ドライブ 3412 に加えて、MPU 3413 と、CD-ROM ドライブ 3412 及び FD ドライブ 3411 に接続されたバス 3414 と、ブートアッププログラム等のプログラムを記憶するための ROM 3415 とに接続され、アプリケーションプログラムの命令を一時的に記憶するとともに一時記憶空間を提供するための RAM 3416 と、アプリケーションプログラム、システムプログラム、及びデータを記憶するためのハードディスク 3417 とを含む。ここでは、図示しないが、コンピュータ 341 は、さらに、LAN への接続を提供するネットワークカードを含んでも良い。

10

【0328】

コンピュータシステム 340 に、上述した実施の形態の異表記取得装置等の機能を実行させるプログラムは、CD-ROM 3501、または FD 3502 に記憶されて、CD-ROM ドライブ 3412 または FD ドライブ 3411 に挿入され、さらにハードディスク 3417 に転送されても良い。これに代えて、プログラムは、図示しないネットワークを介してコンピュータ 341 に送信され、ハードディスク 3417 に記憶されても良い。プログラムは実行の際に RAM 3416 にロードされる。プログラムは、CD-ROM 3501、FD 3502 またはネットワークから直接、ロードされても良い。

20

【0329】

プログラムは、コンピュータ 341 に、上述した実施の形態の異表記取得装置等の機能を実行させるオペレーティングシステム (OS)、またはサードパーティプログラム等は、必ずしも含まなくても良い。プログラムは、制御された態様で適切な機能 (モジュール) を呼び出し、所望の結果が得られるようにする命令の部分のみを含んでいれば良い。コンピュータシステム 340 がどのように動作するかは周知であり、詳細な説明は省略する。

30

【0330】

また、上記プログラムを実行するコンピュータは、単数であってもよく、複数であってもよい。すなわち、集中処理を行ってもよく、あるいは分散処理を行ってもよい。

【0331】

また、上記各実施の形態において、一の装置に存在する 2 以上の通信手段は、物理的に一の媒体で実現されても良いことは言うまでもない。

【0332】

また、上記各実施の形態において、各処理 (各機能) は、単一の装置 (システム) によって集中処理されることによって実現されてもよく、あるいは、複数の装置によって分散処理されることによって実現されてもよい。

40

【0333】

本発明は、以上の実施の形態に限定されることなく、種々の変更が可能であり、それらも本発明の範囲内に包含されるものであることは言うまでもない。

【産業上の利用可能性】

【0334】

以上のように、本発明にかかる異表記取得装置は、用語対の分野を問わず、精度の高い異表記の用語対の抽出が可能となる、という効果を有し、異表記取得装置等として有用である。

【符号の説明】

50

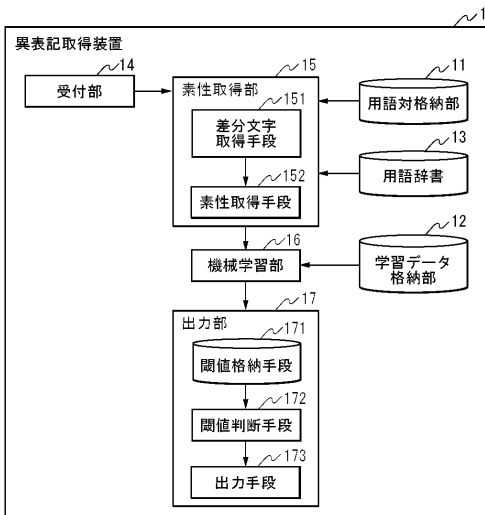
【 0 3 3 5 】

- 1、2 異表記取得装置
- 1 1 用語対格納部
- 1 2 学習データ格納部
- 1 3 用語辞書
- 1 4、2 3 受付部
- 1 5 素性取得部
- 1 6 機械学習部
- 1 7 出力部
- 2 1 異表記用語対格納部
- 2 2 異表記パターン格納部
- 2 4 編集箇所取得部
- 2 5 異表記パターン取得部
- 2 6 異表記パターン蓄積部
- 2 7 用語対生成部
- 1 5 1 差分文字取得手段
- 1 5 2 素性取得手段
- 1 7 1 閾値格納手段
- 1 7 2 閾値判断手段
- 1 7 3 出力手段

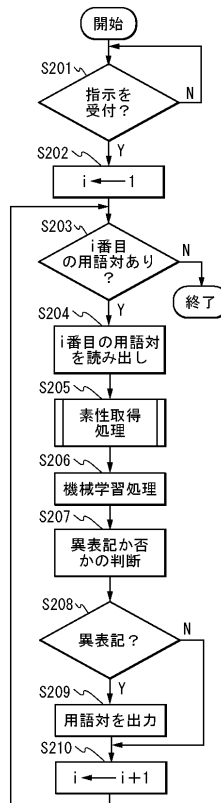
10

20

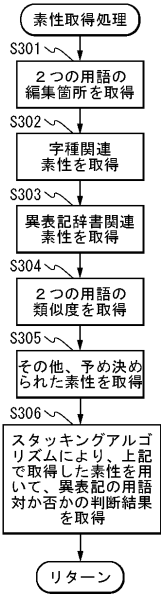
【 図 1 】



【 図 2 】



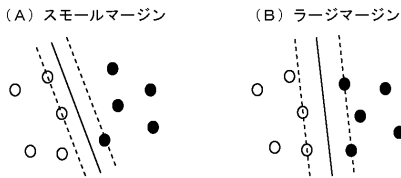
【 図 3 】



【 図 4 】

用語	読み	品詞	代表表記	カテゴリ	ドメイン
子ども	こども	名詞	子供/こども	人
リンゴ	りんご	名詞	林檎/りんご	植物; 人工物	料理・食事
おくれた	おくれた	動詞	送れる/おくれる

【 図 5 】



【 図 6 】

日本語用語対	クローズドデータにおける割合	オープンデータにおける割合
日本語異表記対である	5.25% (745/14185)	5.11% (725/14185)
日本語異表記対でない	94.74% (13440/14185)	94.88% (13460/14185)

【 図 7 】

一致度Kの値	一致度の評価
$0 \leq K < 0.2$	ごく軽度の一致
$0.2 \leq K < 0.4$	軽度の一致
$0.4 \leq K < 0.6$	中等度の一致
$0.6 \leq K < 0.8$	高度の一致
$0.8 \leq K$	ほぼ完全な一致

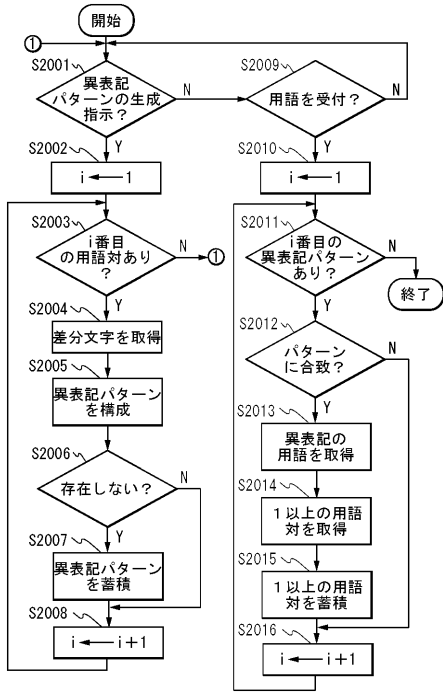
【 図 8 】

手法	10CV		OC	
	正解率	F値	正解率	F値
ベースライン	96.41%	0.483	96.36%	0.452
全素性	99.15%	0.916	99.12%	0.912

【 図 9 】

省いた素性	10CV		OC		10CV		OC	
	全素性	素性選択	全素性	素性選択	全素性	素性選択	全素性	素性選択
S1	1789	7809	7775	1456	6021	3474	S53	553
S2	4585	4418	8652	847	0	0	S54	9462
S3	3918	5810	9018	847	2805	6079	S55	9996
S4	3390	5869	7688	2302	5639	4167	S56	6061
S5	6081	2803	3388	5243	3936	5844	S57	0
S6	0	9933	4131	5694	32	2816	S58	10000
S7	6352	0	0	0	2220	7769	S59	0
S8	5203	4477	1810	6978	3375	6123	S60	0
S9	5150	4294	777	9196	3661	0	S61	5062
S10	285	9376	128	9761	5283	3910	S62	6086
S11	3504	6156	5348	3857	0	0	S63	8658
S12	4272	5631	7275	2414	0	0	S64	3463
S13	9246	469	8754	772	3488	6006	S65	0
S14	6080	2802	6320	0	3422	6079	S66	0
S15	3480	3408	6355	0	8696	0	S67	9999
S16	3505	3423	6292	0	6325	0	S68	8075
S17	3662	4009	8610	0	9021	1854		1859
S18	3468	3499	6319	0	3167	6402		9509
S19	0	6217	0	0	1230	8505		6302
S20	1780	6126	0	0	0	0		0
S21	0	0	3477	5158	0	0		6341
S22	6137	1818	6380	0	2114	7708		10000
S23	0	6324	6010	3497	9496	0		0
S24	0	0	3429	3490	0	0		0
S25	6101	1746	7686	954	3514	3394		0
S26	3983	3901	6304	0	6304	0		0
			9546	0	9546	0		0

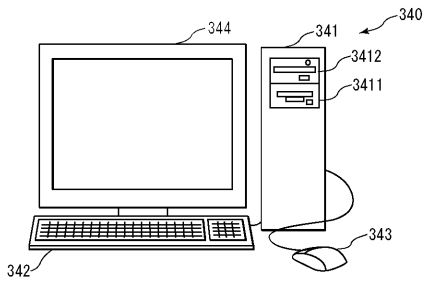
【図20】



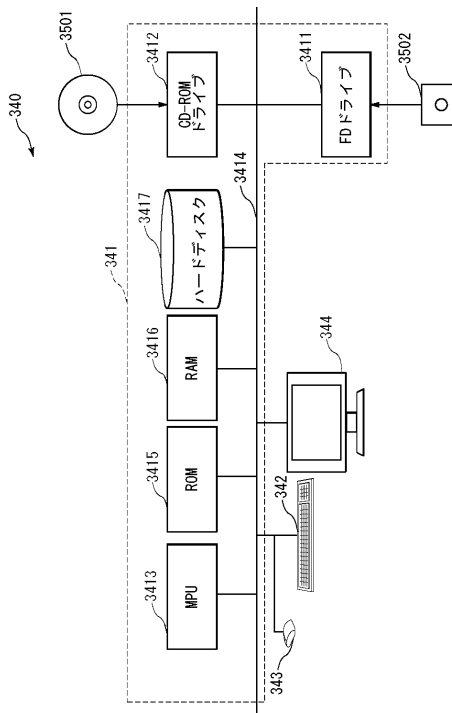
【図21】

異表記パターン	
del	・
del	—
del	い
—	1
イ	イ
取	と
当	あ
言	い
⋮	⋮

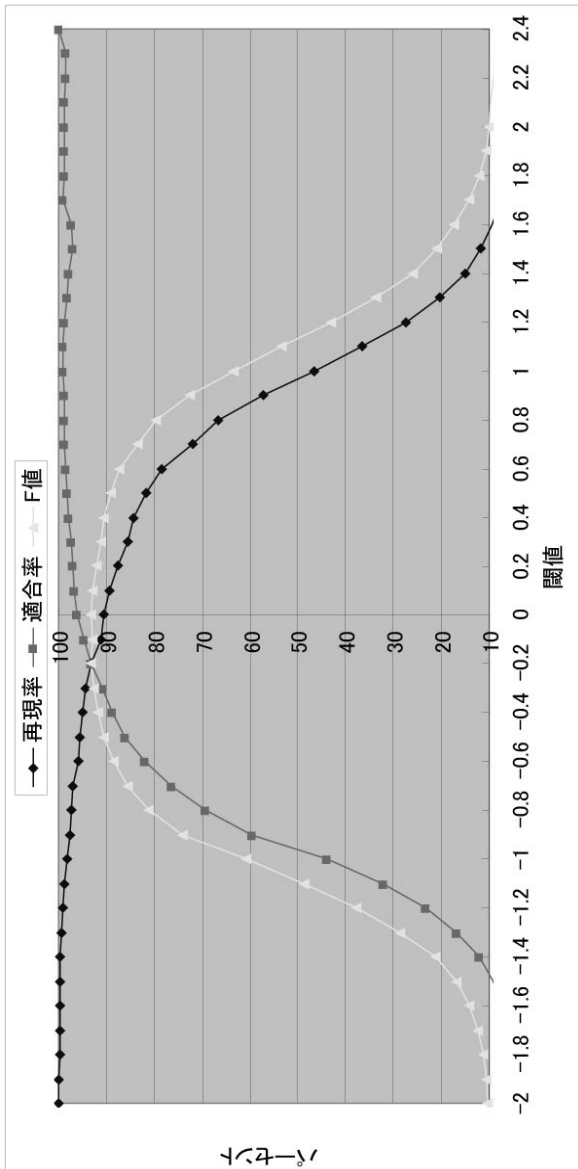
【図22】



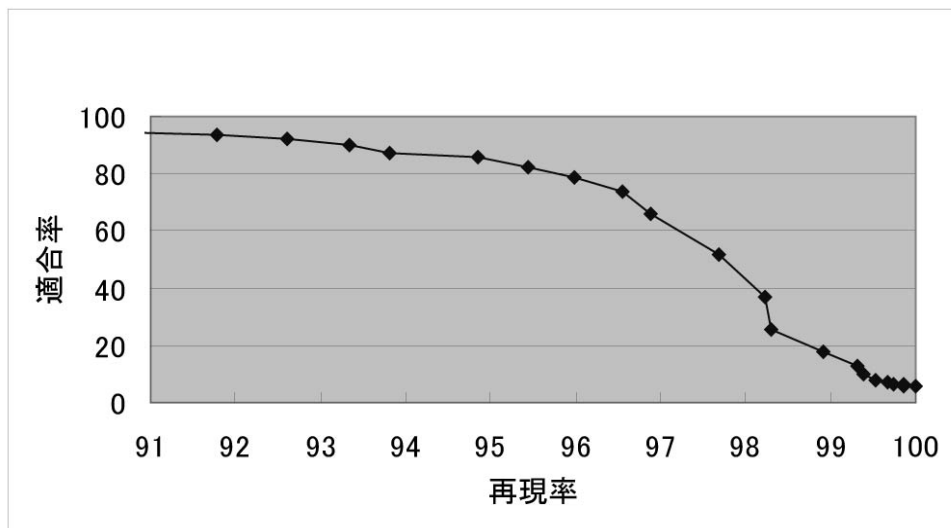
【図23】



【 図 1 4 】



【 図 1 5 】



フロントページの続き

- (72)発明者 風間 淳一
東京都小金井市貫井北町4 - 2 - 1 独立行政法人情報通信研究機構内
- (72)発明者 黒田 航
東京都小金井市貫井北町4 - 2 - 1 独立行政法人情報通信研究機構内
- (72)発明者 藤田 篤
東京都小金井市貫井北町4 - 2 - 1 独立行政法人情報通信研究機構内
- Fターム(参考) 5B075 ND03 QM08 UU01
5B091 AA15 AB17 CC01 CC16