

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号

特許第5327054号
(P5327054)

(45) 発行日 平成25年10月30日(2013.10.30)

(24) 登録日 平成25年8月2日(2013.8.2)

| (51) Int.Cl. | | F I | |
|----------------|---------------|------------------|-----------------------|
| G 1 0 L | 15/18 | (2013.01) | G 1 0 L 15/18 3 0 0 H |
| G 1 0 L | 15/10 | (2006.01) | G 1 0 L 15/10 3 0 0 G |
| G 1 0 L | 15/065 | (2013.01) | G 1 0 L 15/06 3 1 0 Z |

請求項の数 12 (全 16 頁)

| | | | |
|---------------|------------------------------|-----------|-------------------------|
| (21) 出願番号 | 特願2009-546202 (P2009-546202) | (73) 特許権者 | 000004237 |
| (86) (22) 出願日 | 平成20年11月27日(2008.11.27) | | 日本電気株式会社 |
| (86) 国際出願番号 | PCT/JP2008/071500 | | 東京都港区芝五丁目7番1号 |
| (87) 国際公開番号 | W02009/078256 | (74) 代理人 | 100102864 |
| (87) 国際公開日 | 平成21年6月25日(2009.6.25) | | 弁理士 工藤 実 |
| 審査請求日 | 平成23年9月8日(2011.9.8) | (72) 発明者 | 越仲 孝文 |
| (31) 優先権主張番号 | 特願2007-326132 (P2007-326132) | | 東京都港区芝五丁目7番1号 日本電気株式会社内 |
| (32) 優先日 | 平成19年12月18日(2007.12.18) | | |
| (33) 優先権主張国 | 日本国(JP) | | |
| | | 審査官 | 菊池 智紀 |

最終頁に続く

(54) 【発明の名称】 発音変動規則抽出装置、発音変動規則抽出方法、および発音変動規則抽出用プログラム

(57) 【特許請求の範囲】

【請求項1】

音声データを記憶する音声データ記憶手段と、
前記音声データの標準形発音を表す標準形発音データを記憶する標準形発音記憶手段と、
前記標準形発音データからサブワード言語モデルを生成するサブワード言語モデル生成手段と、
前記サブワード言語モデルを用いて前記音声データを認識する音声認識手段と、
前記音声認識手段が出力する認識結果と、前記標準形発音データとを比較して、これらの差分を抽出する差分抽出手段と、
前記サブワード言語モデルの重み値を制御する言語モデル重み制御手段を備え、
前記言語モデル重み制御手段は、
複数の重み値を出力し、
前記音声認識手段は、
前記複数の重み値のそれぞれについて、前記音声データを認識し、
前記言語モデル重み制御手段は、
前記重み値を制御する際、あらかじめ定められた値の組に基づいて、所定の回数だけ前記重み値を更新し、
前記言語モデル重み制御手段は、

前記差分に応じて、前記重み値を更新することの有無を随時決定する
発音変動規則抽出装置。

【請求項 2】

前記言語モデル重み制御手段は、
前記差分が所定のしきい値よりも小さい場合に、前記重み値を減少させるように更新す
る

請求項 1 記載の発音変動規則抽出装置。

【請求項 3】

前記言語モデル重み制御手段は、
前記差分が所定のしきい値よりも大きい場合に、前記重み値を増加させるように更新す
る

請求項 2 記載の発音変動規則抽出装置。

【請求項 4】

前記差分抽出手段は、
前記差分を、前記認識結果と前記標準形発音データとの間の編集距離として計算する
請求項 1 乃至 3 のいずれかに記載の発音変動規則抽出装置。

【請求項 5】

前記差分抽出手段は、
前記差分として、前記認識結果と前記標準形発音データとの相違箇所の文字列対と、前
記認識結果が得られたときに前記音声認識手段が前記言語モデル重み制御手段から受け取
ったサブワード言語モデルの重み値とを含む発音変動事例を抽出する

請求項 1 乃至 4 のいずれかに記載の発音変動規則抽出装置。

【請求項 6】

前記発音変動事例から発音変動の確率的規則を生成する発音変動確率推定手段を更に備
える

請求項 5 記載の発音変動規則抽出装置。

【請求項 7】

前記発音変動確率推定手段は、
ある発音変動事例が観測されたときのサブワード言語モデルの重み値の大きさに応じて
、前記ある発音変動事例の発現確率が高くなるように前記発音変動の確率的規則を生成す
る

請求項 6 記載の発音変動規則抽出装置。

【請求項 8】

音声データの標準形発音を表す標準形発音データを記憶することと、
前記標準形発音データからサブワード言語モデルを生成することと、
前記サブワード言語モデルを用いて前記音声データを認識することと、
前記認識することによる認識結果と、前記標準形発音データとを比較して、これらの差
分を抽出することと、

前記サブワード言語モデルの重み値を制御することと
を具備し、

前記制御することは、

複数の重み値を出力することを含み、

前記認識することは、

前記複数の重み値のそれぞれについて、前記音声データを認識することを含み

前記制御することは、

前記重み値を制御する際、あらかじめ定められた値の組に基づいて、所定の回数だけ前
記重み値を更新することと、

前記差分に応じて、前記重み値を更新することの有無を随時決定することと、

前記差分が所定のしきい値よりも小さい場合に、前記重み値を減少させるように更新す
ることと、

10

20

30

40

50

前記差分が所定のしきい値よりも大きい場合に、前記重み値を増加させるように更新することを含む

発音変動規則抽出方法。

【請求項 9】

前記抽出することは、

前記差分を、前記認識結果と前記標準形発音データとの間の編集距離として計算することと、

前記差分として、前記認識結果と前記標準形発音データとの相違箇所の文字列対と、前記認識結果が得られたときに受け取ったサブワード言語モデルの重み値とを含む発音変動事例を抽出することを含む

10

請求項 8 記載の発音変動規則抽出方法。

【請求項 10】

前記発音変動事例から発音変動の確率的規則を生成することを更に具備し、

前記確率的規則を生成することは、

ある発音変動事例が観測されたときのサブワード言語モデルの重み値の大きさに応じて、前記ある発音変動事例の発現確率が高くなるように前記発音変動の確率的規則を生成することを含む

請求項 9 記載の発音変動規則抽出方法。

【請求項 11】

音声データを記憶する音声データ記憶手段と、

20

前記音声データの標準形発音を表す標準形発音データを記憶する標準形発音記憶手段と、

前記標準形発音データからサブワード言語モデルを生成するサブワード言語モデル生成手段と、

前記サブワード言語モデルを用いて前記音声データを認識する音声認識手段と、

前記音声認識手段が出力する認識結果と、前記標準形発音データとを比較して、これらの差分を抽出する差分抽出手段と、

前記サブワード言語モデルの重み値を制御する言語モデル重み制御手段と、

としてコンピュータを機能させるための

発音変動規則抽出用プログラムであって、

30

前記言語モデル重み制御手段は、

複数の重み値を出力し、

前記音声認識手段は、

前記複数の重み値のそれぞれについて、前記音声データを認識し、

前記言語モデル重み制御手段は、

前記重み値を制御する際、あらかじめ定められた値の組に基づいて、所定の回数だけ前記重み値を更新し、

前記差分に応じて、前記重み値を更新することの有無を随時決定し、

前記差分が所定のしきい値よりも小さい場合に、前記重み値を減少させるように更新し

40

、前記差分が所定のしきい値よりも大きい場合に、前記重み値を増加させるように更新する

発音変動規則抽出用プログラム。

【請求項 12】

前記差分抽出手段は、

前記差分を、前記認識結果と前記標準形発音データとの間の編集距離として計算し、

前記差分として、前記認識結果と前記標準形発音データとの相違箇所の文字列対と、前記認識結果が得られたときに前記音声認識手段が前記言語モデル重み制御手段から受け取ったサブワード言語モデルの重み値とを含む発音変動事例を抽出し、

前記発音変動事例から発音変動の確率的規則を生成する発音変動確率推定手段を更に備

50

え、

前記発音変動確率推定手段は、

ある発音変動事例が観測されたときのサブワード言語モデルの重み値の大きさに応じて、前記ある発音変動事例の発現確率が高くなるように前記発音変動の確率的規則を生成する

請求項 1 1 記載の発音変動規則抽出用プログラム。

【発明の詳細な説明】

【技術分野】

【0001】

本発明は、発音変動規則抽出装置、発音変動規則抽出方法、および発音変動規則抽出用プログラムに関し、特に、対応する書き起こしテキストが付随した音声データ等から、自由な話し言葉によく現れる発音変動の規則を抽出することができる発音変動規則抽出装置、発音変動規則抽出方法、および発音変動規則抽出用プログラムに関する。

10

【背景技術】

【0002】

『堤、加藤、小坂、好田著「発音変形依存モデルを用いた講演音声認識」電子情報通信学会論文誌、第 J 89 - D 巻、2 号、305 ~ 313 頁、2006 年』、『秋田、河原著「話し言葉音声認識のための汎用的な統計的発音変動モデル」電子情報通信学会論文誌、第 J 88 - D 2 巻、9 号、1780 ~ 1789 頁、2005 年』に、発音変動規則抽出装置の一例が記載されている。図 1 に示すように、この発音変動規則抽出装置 200 は、標準形発音記憶手段 201 と、変形発音記憶手段 202 と、差分抽出手段 203 と、発音変動計数手段 204 とから構成されている。

20

【0003】

このような構成を有する発音変動規則抽出装置 200 は次のように動作する。すなわち、差分抽出部 203 は、標準形発音記憶部 201 および変形発音記憶部 202 からそれぞれ書き起こしテキストを抽出し、差分、すなわち相違箇所を抽出する。

【0004】

ここで標準形発音記憶部 201 および変形発音記憶部 202 には、長時間の音声データの発音内容を書き起こした結果である書き起こしテキストが記憶されている。より具体的には、標準形発音記憶部 201 には、例えば以下のような書き起こしテキストが格納されている。

30

その ような しゅじゅつ を ほぼ まいにち おこない ました

【0005】

また、変形発音記憶部 202 には、標準形発音記憶部 201 に記憶された書き起こしテキストと対応する形で、例えば以下のような書き起こしテキストが格納されている。

その ような しじつ を ほぼ まいんち おこない ました

【0006】

標準形発音記憶部 201 には、元となった音声データの標準的な発音、つまり正しく発音された場合に観測されるべき本来の発音を書き起こしテキストとして記憶されている。一方、変形発音記憶部 202 には、音声データを実際に人が聞いて、聞こえるままの発音を忠実に書き起こした書き起こしテキストが記憶されている。上の例では、標準形発音「しゅじゅつ(手術)」、「まいにち(毎日)」に対して、それぞれ「しじつ」、「まいんち」という変形発音が記憶されている。

40

【0007】

差分抽出部 203 は、標準形の書き起こしテキストと変形の書き起こしテキストとを比較し、相違する箇所の文字列対を抽出する。上の例では、「しゅじゅつ」と「しじつ」、「まいにち」と「まいんち」が抽出される。以下、これらの対を発音変動事例と呼ぶ。また、標準形発音と変形発音が等しい、すなわち変形がない場合の発音変動事例を特に、恒等発音変動と呼ぶことにする。

【0008】

50

発音変動計数部 204 は、差分抽出部 203 から発音変動事例を受け取り、同じ標準形、同じ変形ごとに分類し、恒等発音変動も含めて観測回数を計数する。さらに、計数結果を正規化して確率値に変換する。例えば、上の例で標準形発音「まいにち」に対する変形発音として「まいにち(恒等変形)」、「まいんち」、「まいち」、「まんいち」があり、それぞれ 966 回、112 回、13 回、2 回観測されたとする。標準形発音「まいにち」の観測回数は $966 + 112 + 13 + 2 = 1093$ であるから、確率値に変換すると、

まいにち まいにち 0.884 (966 / 1093)

まいにち まいんち 0.102 (112 / 1093)

まいにち まいち 0.012 (13 / 1093)

まいにち まんいち 0.002 (2 / 1093)

となる。この結果は、標準形発音「まいにち」に対する変形発音の出現傾向に関する確率的な規則と解釈できる。発音変動計数部 204 は、上記結果を発音変動規則として出力する。

【0009】

なお、上の例では、標準形発音や変形発音を単語単位で扱っているが、他の単位、例えば所定の長さの音素(母音、子音等、音声を構成する最小単位)の系列として扱うことも可能である。また、上記確率値を計算する際に、適当な平滑化操作、例えば観測回数が所定値に満たない特殊な発音変動規則を無視する等を行ってもよい。

【0010】

『緒方、有木著「発音変形と音響的誤り傾向を考慮した話し言葉音声認識の検討」日本音響学会 2003 年春季研究発表会講演論文集、9~10 頁、2003 年 3 月』、『緒方、後藤、浅野著「話し言葉音声認識のための動的発音モデリング法の検討」日本音響学会 2004 年春季研究発表会講演論文集、203~204 頁、2004 年 3 月』に、発音変動規則抽出装置の別の一例が記載されている。図 2 に示すように、この発音変動規則抽出装置 300 は、音声データ記憶部 301 と、標準形発音記憶部 302 と、音節辞書記憶部 303 と、音響モデル記憶部 304 と、音声認識部 305 と、差分抽出部 306 と、発音変動計数部 307 とから構成されている。

【0011】

このような構成を有する発音変動規則抽出装置 300 は次のように動作する。すなわち、音声認識部 305 は、音節辞書記憶部 303 に記憶された辞書、および音響モデル記憶部 304 に記憶された音響モデルを用いて、音声データ記憶部 301 に記憶された音声データに対して、公知の連続音節認識処理を行い、認識結果の音節系列を出力する。

【0012】

ここで、音節辞書記憶部 303 に記憶された辞書は、日本語の場合、あ、い、う、え、お、か、き、く、け、こ、... のようにあらゆる音節を記録したリストであり、各音節について、その音響的特徴が参照できるよう、音響モデルへのポインタが付与されている。他の言語の場合でも、その言語に即して適当な単位を定義し、辞書を構成することが可能である。また、音響モデル記憶部 304 に記憶された音響モデルは、所定の認識単位、すなわち音節、音素などに関する音響的特徴が、公知の隠れマルコフモデル等の手法に基づいて記述されたモデルである。

【0013】

差分抽出部 306 は、音声認識部 305 から認識結果を、標準形発音記憶部 302 から書き起こしテキストをそれぞれ受け取り、両者の差分、すなわち相違箇所を抽出する。ここで、標準形発音記憶部 302 に記憶された書き起こしテキストは、図 1 の標準形発音記憶部 201 に記憶された書き起こしテキストと同様であるが、音声データ記憶部 301 に記憶された音声データと対応付いている、すなわち、音声データ記憶部 301 の音声データの内容が正しく発音された場合に観測されるべき本来の発音を書き起こしテキストとして記憶されている。発音変動計数部 307 は、図 1 の発音変動計数部 204 と同様の動作により、差分抽出部 306 から発音変動事例を受取り、発音変動規則を出力する。

【0014】

10

20

30

40

50

『大西著「認識誤りの話者性を考慮した発声変形抽出と認識辞書拡張」日本音響学会2007年春季研究発表会講演論文集、65～66頁、2007年3月』に、発音変動規則抽出装置のさらに別の一例が記載されている。図3に示すように、この発音変動規則抽出装置400は、音声データ記憶部401と、標準形発音記憶部402と、単語言語モデル・辞書記憶部403と、音響モデル記憶部404と、音声認識部405と、差分抽出部406と、発音変動計数部407とから構成されている。

【0015】

このような構成を有する発音変動規則抽出装置400は次のように動作する。すなわち、音声認識部405は、単語言語モデル・辞書記憶部403に記憶された言語モデルと辞書、および音響モデル記憶部404に記憶された音響モデルを用いて、音声データ記憶部401に記憶された音声データに対して公知の連続単語認識処理を行い、認識結果の単語系列を出力する。

10

【0016】

ここで、単語言語モデル・辞書記憶部403に記憶された辞書および言語モデルは、一般的な大語彙連続音声認識システムが備える辞書および言語モデルと同様のものでよい。辞書は数万語の単語を含み、各単語について、その発音と、音響的特徴を参照するに足る音響モデルへのポイントが付与されている。言語モデルは、公知のn-gramモデルに基づき、n-1個の単語並びを仮定した場合に、次にどのような単語が現れるかを確率の形で規定したモデルとなる。

【0017】

20

また、音響モデル記憶部404に記憶された音響モデルは、図2の音響モデル記憶部304に記憶された音響モデルと同様、所定の認識単位、すなわち音節、音素などに関する音響的特徴が、公知の隠れマルコフモデル等の手法に基づいて記述されたモデルである。

【0018】

差分抽出部406は、図2の差分抽出部306と同様の動作により、音声認識部405から認識結果を、標準形発音記憶部402から書き起こしテキストをそれぞれ受け取り、両者の差分、すなわち相違箇所を抽出する。ここで、標準形発音記憶部402に記憶された書き起こしテキストは、図2の標準形発音記憶部302と同様であり、音声データ記憶部401に記憶された音声データと対応付いていることが必要である。発音変動計数部407は、図1の発音変動計数部204や、図2の発音変動計数部307と同様の動作により、差分抽出部406から発音変動事例を受取り、発音変動規則を出力する。

30

【0019】

これらの5つの文献に記載された発音変動規則抽出装置100, 200, 300における第1の問題点は、発音変動規則やその元となる発音変動事例を得るために多大な労力を要するという点である。その理由は、標準形発音とそれに対応する変形発音を大量に用意する必要があるためである。妥当性の高い発音変動規則を獲得するために、図1の発音変動規則抽出装置100では、大量の音声データの書き起こしにより、標準形発音記憶部201に記憶される標準形発音、および変形発音記憶部202に記憶される変形発音をあらかじめ作成しておく必要がある。しかしながら、標準形発音と変形発音、特に後者の作成は、音声の聞き取りに習熟した作業者が注意深く音声聞き、曖昧で判断が付きにくい変形発音を文字列として書き起こす作業となるため、とりわけ時間と労力がかかる。

40

【0020】

第2の問題点は、汎化性の高い発音変動規則を得ることが難しいということである。その理由は、自由な話し言葉の音声データから正確な発音変動事例を得ることが難しいためである。例えば、図1の発音変動規則抽出装置100では、人手により変形発音を書き起こすが、大量の書き起こしを得るためには、多数の作業者が分担して作業を行うのが普通である。しかしながら、話し言葉の発音は本質的に曖昧であるため、書き起こしには作業者の主観が多分に入り、作業結果にばらつきが生じる。また、図2の発音変動規則抽出装置200では、音声認識部により統一的な基準で変形発音を自動的に取得することが可能である。しかしながら、現在の音声認識の技術水準では、言語的な事前知識のない状況で

50

音節の並びを求める連続音節認識処理を正確に行うことは極めて難しい。例えば、「ひろしま」という発声を連続音節認識すると、「けるせま」、「かるりか」というような、実際の発音の変動とは程遠い結果がしばしば得られる。すなわち、連続音節認識を適用しても、ランダムで有用性の乏しい文字列が得られるのみである。

【0021】

図3の発音変動規則抽出装置300でも、単語辞書と言語モデルという事前知識が利用可能とはいえ、図2の発音変動規則抽出装置200と同様、音声認識の不正確さの問題がなお残る。さらに図3の発音変動規則抽出装置300では、単語辞書と言語モデルが音声認識処理における言語的な制約として働くことから、得られる発音変動事例は単語辞書と言語モデルの影響を受ける。よって、実際に起こっている発音変動現象とは一般に異なる発音変動事例が得られる。例えば、「せんたくき（洗濯機）」が「せんたっき」に変わったり、「しょくぱん（食ぱん）」が「しょっぱん」に変わったりするような現象は一般的にみられるが、図3の発音変動規則抽出装置300では、単語辞書に含まれる単語の組合せとしてしか音声認識結果が得られないため、「せんたっき」という発音と一致する認識結果が得られる保証はどこにもない。

【発明の開示】

【0022】

本発明の目的は、少ない労力で、発音変動事例を頑健に検出し、汎化性の高い発音変動規則を獲得することにある。

【0023】

本発明の一つ目のアスペクトによる発音変動規則抽出装置は、音声データ記憶部と、標準形発音記憶部と、サブワード言語モデル生成部と、音声認識部と、差分抽出部とを備える。音声データ記憶部は、音声データを記憶する。標準形発音記憶部は、音声データの標準形発音を表す標準形発音データを記憶する。サブワード言語モデル生成部は、標準形発音データからサブワード言語モデルを生成する。音声認識部は、サブワード言語モデルを用いて音声データを認識する。差分抽出部は、音声認識部が出力する認識結果と、標準形発音データとを比較して、これらの差分を抽出する。

【0024】

本発明の二つ目のアスペクトによる発音変動規則抽出方法は、記憶することと、生成することと、認識することと、抽出することとを具備する。記憶することは、音声データの標準形発音を表す標準形発音データを記憶する。生成することは、標準形発音データからサブワード言語モデルを生成する。認識することは、サブワード言語モデルを用いて音声データを認識する。抽出することは、認識することによる認識結果と、標準形発音データとを比較して、これらの差分を抽出する。

【0025】

本発明の三つ目のアスペクトによる発音変動規則抽出プログラムは、コンピュータを、音声データ記憶部と、標準形発音記憶部と、サブワード言語モデル生成部と、音声認識部と、差分抽出部ととして機能させる。音声データ記憶部は、音声データを記憶する。標準形発音記憶部は、音声データの標準形発音を表す標準形発音データを記憶する。サブワード言語モデル生成部は、標準形発音データからサブワード言語モデルを生成する。音声認識部は、サブワード言語モデルを用いて音声データを認識する。差分抽出部は、音声認識部が出力する認識結果と、標準形発音データとを比較して、これらの差分を抽出する。このプログラムは、コンピュータ読み取り可能な記録媒体に格納でき、その記録媒体からコンピュータに読み込ませることができる。

【0026】

本発明による効果は、正確で汎化性の高い発音変動規則を獲得できることにある。その理由は、制約のないサブワードを単位とした音声認識を基本として、音声データに対応する標準形発音という言語制約を任意の強さでかけながら音声認識を行うことにより、個々の音声データの違いに依存せず、多くの音声データに共通して現れる発音変動を抽出できるからである。また、人手作業で発生する、主観判断によるばらつきもないからである。

【図面の簡単な説明】

【0027】

【図1】従来技術の一例を示すブロック図である。

【図2】従来技術の一例を示すブロック図である。

【図3】従来技術の一例を示すブロック図である。

【図4】本発明による第1の発明を実施するための最良の形態の構成を示すブロック図である。

【図5】第1の発明を実施するための最良の形態の動作の具体例を示す図である。

【図6】第1の発明を実施するための最良の形態の動作の具体例を示す図である。

【図7】第1の発明を実施するための最良の形態の動作の具体例を示す図である。

10

【図8】第1の発明を実施するための最良の形態の動作を示す流れ図である。

【図9】本発明による第2の発明を実施するための最良の形態の構成を説明するブロック図である。

【発明を実施するための最良の形態】

【0028】

本発明を実施するための最良の形態の一つについて図面を参照して詳細に説明する。図4を参照すると、本発明の第1の実施の形態における発音変動規則抽出装置100は、音声データ記憶部101と、標準形発音記憶部102と、サブワード言語モデル・辞書生成部103と、音響モデル記憶部104と、音声認識部105と、差分抽出部106と、発音変動確率推定部107と、言語モデル重み制御部108とを含む。

20

【0029】

音声データ記憶部101は、発音変動事例が含まれると思われる多数の音声データを記憶する。標準形発音記憶部102は、音声データ記憶部101に記憶された音声データの書き起こしテキストを記憶する。ここに書き起こしテキストは、音声データの発音内容が、標準形で書き起こされたテキストデータであり、ひらがな、カタカナ、あるいは任意の発音記号の列で表される。ひらがなで表した書き起こしテキストの例を以下に示す。

みなさん こんにちは

発音を表す書き起こしであることから、「こんにちは」は「こんにちわ」と記述される。

【0030】

音響モデル記憶部104は、後述する音声認識部105が音声認識処理を行う際に必要となる音響モデルを記憶する。音響モデルは、隠れマルコフモデルに基づいて個々の音素（日本語の場合は母音 a , i , u , e , o , 子音 k , s , t , n , ...）をモデル化したもの等を用いることができる。

30

【0031】

サブワード言語モデル・辞書生成部103は、標準形発音記憶部102に記憶された書き起こしテキストを用いて、後述する音声認識部105が音声認識処理を行う際に必要となるサブワード言語モデル・辞書を生成する。ここに辞書は、例えばサブワードを音節とした場合、「あ、い、う、え、お、か、き、く、け、こ、...」の各音節を1単語として構成された辞書である。各単語、すなわち各音節についてその音響的特徴がわかるように、例えば「あ a」、「か ka」、「さ sa」、...のように、単語から音響モデルへの

40

【0032】

また、サブワード言語モデルは、サブワードを単語として、音声認識で広く用いられる n - g r a m モデルの考え方に基づき、各単語について、履歴 h に続いて単語 w が出現する確率 $P(w | h)$ を規定したモデルである。具体的には、例えば n = 3 のモデル (t r i g r a m モデル) の場合、音節 s_{i-2} , s_{i-1} がこの順に出現したとき、次に音節 s_i が出現する確率 $P(s_i | s_{i-2}, s_{i-1})$ が種々の s_{i-2} , s_{i-1} , s_i について規定されている。さらに、ここで生成されるサブワード言語モデルは、標準形発

50

音記憶部 102 に記憶された標準形の書き起こしテキストを学習データとして生成される。

【0033】

例えば、上述の例の みなさん こんにちは という1発話を学習データに用いて生成されるサブワード言語モデルは、図5のように表される。なお、図5に示されていない履歴 h を含む確率については等確率を与えることができる。また、図5に示された履歴 h を含むが図5に示されていない確率については0とすることができる。図5の " h " 欄において、 \square は空文字列であり、ここでは特に文頭を意味する。また、 $\#$ は単語間のポーズ（無音）を意味し、単語間にポーズが入る場合と入らない場合とで確率を二分している。上述のように、1発話のような短い単位の書き起こしテキストから学習されたサブ

10

【0034】

なお、ここでは1発話を単位としてサブワード言語モデルを生成しているが、第1の実施の形態はこの単位の取り方を制限するものではなく、数個の発話を1単位とする、あるいは1個ないし数個の単語を単位とすることも可能である。また、サブワード言語モデル・辞書を構成する単語の単位を、ここでは音節としているが、一般にサブワードと呼ばれる単位、すなわち、音節、半音節、モーラ、音素等を単位としてサブワード言語モデル・辞書を生成することが可能である。

【0035】

言語モデル重み制御部 108 は、サブワード言語モデルの重み値を少なくとも1回決定し、音声認識部 105 に送る。1回だけ決定する場合は、例えばあらかじめ実験的に定めた定数を使えばよい。また、複数回決定する場合は、同様にあらかじめ実験的に定めた複数個の定数を順に選択したり、あらかじめ定めた初期値から、あらかじめ定めた値を順次加算あるいは減算すればよい。ここにサブワード言語モデルの重み値とは、一般に正の値をとり、後述する音声認識部 105 が音響モデルやサブワード言語モデル・辞書を参照して音声認識処理を行う際に、サブワード言語モデルから計算されるスコアをどの程度重視するかを規定するパラメータである。

20

【0036】

音声認識部 105 は、音響モデル記憶部 104 から音響モデルを、サブワード言語モデル・辞書生成部 103 から言語モデル・辞書をそれぞれ受け取り、また、言語モデル重み制御部 108 からサブワード言語モデルの重み値を少なくとも1回受け取る。そして、サブワード言語モデルの重み値ごとに、音声データ記憶部 101 に記憶された音声データに対して音声認識処理を行い、認識結果の音節列を求める。なお、音声認識処理は次の数式1により表すことができ、認識対象の音声データ O に対して、認識結果 W が得られる。

30

【0037】

$$W = \arg \max_{W'} \left[\log P(O|W', \theta_{AM}) + \lambda_{LM} \log P(W'|\theta_{LM}) \right] \quad \dots(1)$$

40

ここに、右辺 $\arg \max$ 関数内の第1項、第2項は、それぞれ音響スコア、言語スコアと呼ばれる。 AM は音響モデルであり、音響モデル記憶部 104 に記憶されている。

LM はサブワード言語モデル・辞書であり、サブワード言語モデル・辞書生成部 103 により生成される。 LM はサブワード言語モデルの重み値であり、言語モデル重み制御部 108 により決定される。 W' は、認識結果 W になる候補であり、いずれかの W' が、 W として算出される。 $\arg \max$ は、変数 W' を動かしたときに、最大値を与える W' を求める関数である。

【0038】

サブワード言語モデルの重み値 LM が十分大きい場合、認識結果は、サブワード言語

50

モデルの学習データとなった書き起こしテキストと極めて高い確率で一致する。逆に、サブワード言語モデルの重み値 LM が小さい場合は、認識結果は先述の図 2 に示したような、連続音節認識の結果に近づく。なお、サブワード言語モデルの重み値を設定する代わりに、音響モデルの重み値を設定してもよい。すなわち、言語スコアの項に係数 LM をかける代わりに、音響スコアの項に同様の係数をかけても同じことである。サブワード言語モデルの重み値を大きくすることは、音響モデルの重み値を小さくすることと同値である。

【 0 0 3 9 】

差分抽出部 106 は、音声認識部 105 から少なくとも 1 つの認識結果を、また標準形発音記憶部 102 から標準形の書き起こしテキストをそれぞれ受け取り、両者の差分、すなわち相違箇所を抽出する。図 6 は、差分抽出部 106 が音声認識部 105 から受け取る認識結果の一例である。この例では、複数のサブワード言語モデルの重み値 ($10.0 \sim 0.5$) について、それぞれ得られた認識結果が示されている。差分抽出部 106 は、図 6 の認識結果を、標準形の書き起こしテキストと比較して、図 7 に示すように相違箇所をサブワード言語モデルの重み値とともに抽出する。図 7 の各行を、ここでは発音変動事例と呼ぶ。

10

【 0 0 4 0 】

なお、ここでは単語単位で差分すなわち発音変動事例を抽出しているが、第 1 の実施の形態は単位の取り方を単語に限定するものではなく、他の任意の単位でも実施可能である。例えば、前出の 2 つ目の文献では、所定の長さの音素系列を単位として差分抽出を行っているが、第 1 の実施の形態においても、このような形式による差分抽出を容易に適用することが可能である。

20

【 0 0 4 1 】

発音変動確率推定部 107 は、差分抽出部 106 から発音変動事例を受け取り、標準形発音、変形発音ごとに分類し、発音変動規則を得る。図 7 に示したように、発音変動事例を標準形発音、変形発音、サブワード言語モデルの重み値の組として、音声データ記憶部 101 に記憶された音声データから、数式 2 のような N 個の発音変動事例が得られたとする。

【 0 0 4 2 】

$$\{w_i, \tilde{w}_i, \lambda_i \mid i = 1, 2, \dots, N\} \quad \dots(2)$$

30

サブワード言語モデルの重み値が大きく、言語的制約が強い場合でも観測される変形発音は、一般的に発現しやすいであろうことを考慮すると、標準形発音 w を所与とした発音変動規則が数式 3 のように確率論的に定義される。

【 0 0 4 3 】

$$P(\tilde{w} | w) = \frac{\sum_i \lambda_i \delta_{w, w_i} \delta_{\tilde{w}, \tilde{w}_i}}{\sum_i \lambda_i \delta_{w, w_i}} \quad \dots(3)$$

40

ただし、 i, j はクロネッカのデルタ ($i = j$ なら 1、そうでなければ 0) である。なお、数式 3 の変形例として、サブワード言語モデルの重み値 i を考慮せず、 i を 1 に置き換えて計算してもよい。また、数式 3 の i を、 i を変数とする関数、例えば i の多項式関数などに置き換えてもよい。さらに、数式 3 の確率値を計算する際に、適当な平滑化操作を行ってもよい。適当な平滑化操作とは、例えば、サブワード言語モデルの重み値が小さい発音変動事例を無視する、観測回数が所定値に満たない発音変動事例を無視する等の操作に相当する。

【 0 0 4 4 】

50

次に、図4のブロック図および図8のフローチャートを参照して、第1の実施の形態における動作について詳細に説明する。まず、音声認識部105は、音響モデル記憶部104から音響モデルを読み込む(図8のステップA1)。次に、サブワード言語モデル・辞書生成部103は、標準形発音記憶部102に記憶された1発話分の書き起こしテキストを選択し(ステップA2)、読み込み(ステップA3)、これを学習データとしてサブワード言語モデルを生成し、および、必要に応じて辞書を生成する(ステップA4)。音声認識部105は、サブワード言語モデル・辞書生成部103が生成したサブワード言語モデル・辞書を読み込む(ステップA5)。次に、音声認識部105は、ステップA2で選択された書き起こしテキストに対応する音声データを、音声データ記憶部101から読み込む(ステップA6)。

10

【0045】

言語モデル重み制御部108は、サブワード言語モデルの重み値として所定の値、例えば十分大きな値をセットし、音声認識部105に送る(ステップA7)。音声認識部105は、言語モデル重み制御部108がセットしたサブワード言語モデルの重み値に基づいて、音声認識処理を行い、音声認識結果すなわち音節列と、サブワード言語モデルの重み値を記憶する(ステップA8)。言語モデル重み制御部108は、サブワード言語モデルの重み値を一定量だけ増加又は減少させるなどして、サブワード言語モデルの重み値を更新する(ステップA9)。サブワード言語モデルの重み値の更新が所定回数Iを超えていれば次のステップに進み、そうでなければ、前述のステップA8、A9を繰り返す(ステップA10)。差分抽出部106は、音声認識部105が行った音声認識処理の結果を、図6にすでに示したような形式で受け取り、またステップA3でサブワード言語モデル・辞書生成部103が選択した標準形書き起こしテキストを受け取る。そして、図7や数式2ですでに示したような形式で、両者の相違箇所すなわち発音変動事例を抽出する(ステップA11)。以上示したステップA2からA11までの処理を、未処理の発話がなくなるまで繰り返す(ステップA12)。最後に、発音変動確率推定部107は、差分抽出部106が求めたすべての発音変動事例を、数式3に従ってまとめ上げ、発音変動規則として出力する(ステップA13)。

20

【0046】

なお、第1の実施の形態における音声認識部105とわずかに異なる別の音声認識部を適用することが可能である。この別の音声認識部は、図示しない記憶部に無情報なサブワード言語モデルを記憶している。ここで無情報とは、任意の履歴において各単語が等確率で出現し得ることを意味する。無情報なサブワード言語モデルとは、サブワードを音節とした場合、任意の音節の並びである s_{i-2} , s_{i-1} , s_i について $P(s_i | s_{i-2}, s_{i-1}) = \text{const}$ と表されるモデルのことである。無情報なサブワード言語モデルを用いた場合、音声認識処理は数式1に代わり、数式4となる。

30

【0047】

$$W = \arg \max_{W'} [\log P(O|W', \theta_{AM}) + K \log \{(1 - \lambda_{LM}) P(W' | \theta_{LM0}) + \lambda_{LM} P(W' | \theta_{LM})\}] \quad \dots(4)$$

40

ここにLM0は、無情報なサブワード言語モデルを表す。また、Kはあらかじめ定める定数である(なくてもよい)。数式4を用いた場合でも、言語モデル重み制御部108がサブワード言語モデルの重み値LMを大小させることによって、第1の実施の形態における音声認識部105の場合と同様の結果が別の音声認識部から得られる。ただしこの場合、サブワード言語モデルの重み値には0 <math>LM < 1</math>なる制約が生ずる。よって、言語モデル重み制御部108は、この制約の中でサブワード言語モデルの重み値を決定するよう動作する。

【0048】

また、音声認識部105、差分抽出部106、および言語モデル重み制御部108の動

50

作についても、第 1 の実施の形態とわずかに異なる変形例を適用することが可能である。すなわち、第 1 の実施の形態では、図 8 のステップ A 8、A 9 および A 10 を所定の回数だけ反復するとしているが、以下に述べるように、変形例における差分抽出部の抽出結果に応じて適応的に反復回数を決定することも可能である。

【 0 0 4 9 】

例えば、ステップ A 7 にて、十分大きな値をサブワード言語モデルの重み値の初期値とし、ステップ A 9 で順次サブワード言語モデルの重み値が減少するように動作させる場合は、差分抽出部にて標準形発音と認識結果の差分が所定のしきい値よりも大きくなった時点で反復を止めればよい。ここで、標準形発音と認識結果の差分を定量的に測るには、例えば文字列間の相違度合いの一般的尺度として知られている編集距離などが利用できよう。

10

【 0 0 5 0 】

あるいは逆に、ステップ A 7 にて、十分小さな値をサブワード言語モデルの重み値の初期値とし、ステップ A 9 で順次サブワード言語モデルの重み値が増加するように動作させる場合は、差分抽出部にて標準形発音と認識結果の差分が所定のしきい値よりも小さくなった時点、または標準形発音と認識結果が完全に一致した時点で反復を止めればよい。

【 0 0 5 1 】

次に、第 1 の実施の形態の効果について説明する。第 1 の実施の形態では、標準形発音のみ受理可能なサブワード言語モデルを生成するサブワード言語モデル・辞書生成部 103 と、サブワード言語モデルの重み、すなわちサブワード言語モデルの重み値を決定する言語モデル重み制御部 108 と、サブワード言語モデルおよびその重み値を用いて標準形発音に対応する音声データを認識する音声認識部 105 と、音声認識部 105 が出力する認識結果を標準形発音と比較して相違箇所を発音変動事例として抽出する差分抽出部 106 と、発音変動事例をまとめ上げて発音変動規則を出力する発音変動確率推定部 107 とを備える。そして、いくつかのサブワード言語モデルの重み値で音声認識処理を実行した結果をそれぞれ標準形発音と比較し、抽出される差分を発音変動事例とし、この発音変動事例をサブワード言語モデルの重み値を考慮してまとめ上げるようにしているため、正確で汎化性が高く、発現のしやすさに応じて確率値が付与された発音変動規則を獲得できる。

20

【 0 0 5 2 】

次に、本発明による第 2 の発明を実施するための最良の形態について図面を参照して詳細に説明する。第 2 の実施の形態は、第 1 の実施の形態を、プログラムを用いて実現するものである。このプログラムは、コンピュータを、第 1 の実施の形態における部 101 ~ 108 が結合されたものとして機能させる。図 9 を参照すると、発音変動規則抽出用プログラム 92 は、コンピュータ読み取り可能な記録媒体 90 に格納されていて、コンピュータ 91 に読み込まれ、コンピュータ 91 の動作を制御する。

30

【 0 0 5 3 】

発音変動規則抽出用プログラム 92 は、コンピュータ 91 に読み込まれた後、起動すると、記憶装置 94 内の音声データ記憶部 941 を音声記憶部 101 として機能させ、標準形発音記憶部 942 を標準型発音記憶部 102 として機能させ、および音響モデル記憶部 943 を音響モデル記憶部 104 として機能させる。また、データ処理装置 93 は発音変動規則抽出用プログラム 92 の制御により、第 1 の実施の形態におけるサブワード言語モデル・辞書生成部 103、音声認識部 105、差分抽出部 106、発音変動確率推定部 107、および言語モデル重み制御部 108 として機能し、記憶装置 94 内の音声データ記憶部 941、標準形発音記憶部 942、および音響モデル記憶部 943 に記憶されたデータを処理し、発音変動規則を出力する。

40

【 0 0 5 4 】

本発明によれば、大規模な音声データから発音変動規則を抽出する発音変動抽出装置や、発音変動規則抽出装置をコンピュータに実現するためのプログラムといった用途に適用できる。また、情報入力、情報検索、書き起こし支援、映像インデクシング等に広く用い

50

られる音声認識装置が知られるが、このような音声認識装置が使用する音響モデルや言語モデルを発音変動に対して頑健に作成するための音声認識用モデル作成装置、あるいは発音練習装置、語学学習装置、といった用途にも適用可能である。

【0055】

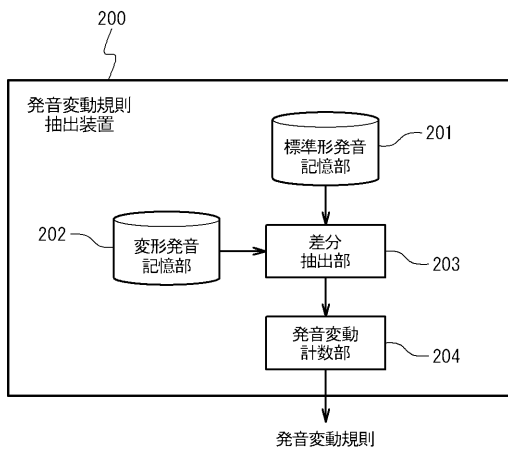
以上、実施の形態を参照して本願発明を説明したが、本願発明は上記実施の形態に限定されるものではない。本願発明の構成や詳細には、請求の範囲に記載された本願発明の技術的思想の範囲内において、当業者が適宜、様々な変形又は変更を加えることが可能である。

【0056】

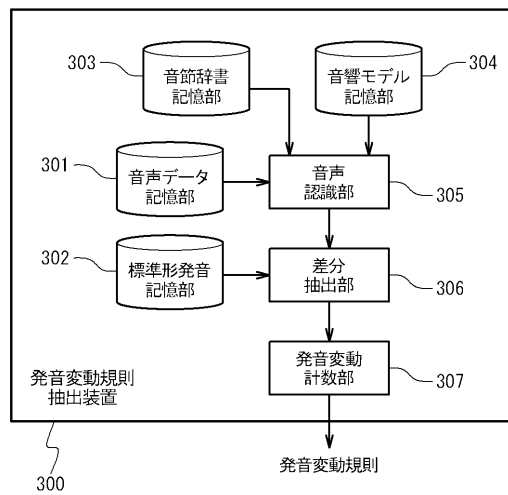
この出願は、2007年12月18日に出願された特許出願番号2007-326132号の日本特許出願に基づいている。本願は、この基礎出願により生じた優先権の利益を享受しており、この基礎出願における開示の内容の全てを、引用により、そっくりそのままここに取り込んでいる。

10

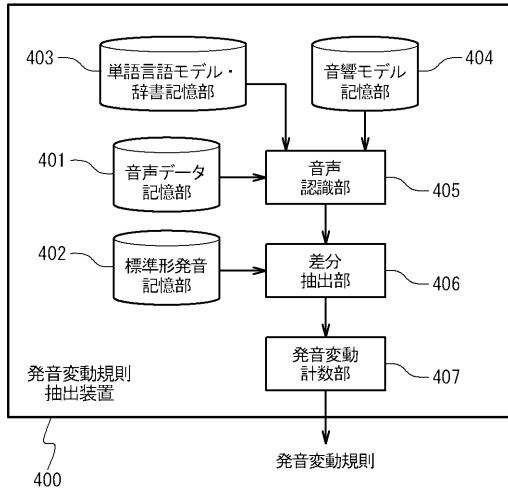
【図1】



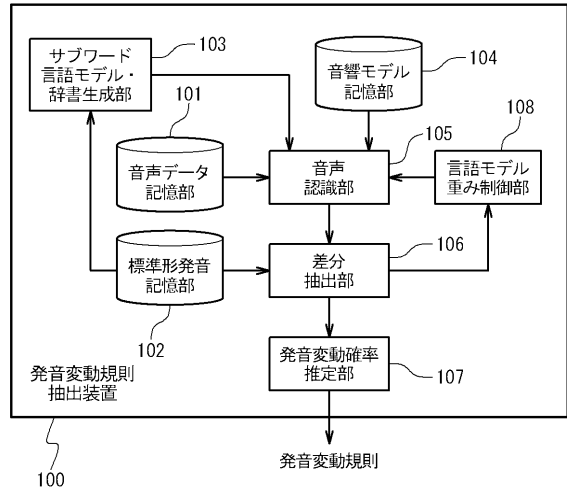
【図2】



【図3】



【図4】



【図5】

| h | w | P(w h) | 備考 |
|----|----|--------|--------------|
| φ | み | 1.0 | 文頭unigram |
| φ | み | 1.0 | 文頭bigram |
| み | な | 1.0 | |
| みな | さ | 1.0 | |
| みな | さん | 1.0 | |
| な | さん | 1.0 | |
| な | # | 0.5 | 単語間ポーズの有無で分岐 |
| ん | ん | 1.0 | |
| ん | # | 0.5 | 単語間ポーズの有無で分岐 |
| ん | ん | 1.0 | |
| ん | に | 1.0 | |
| ん | に | 1.0 | |
| ん | ち | 1.0 | |
| ん | ち | 1.0 | |

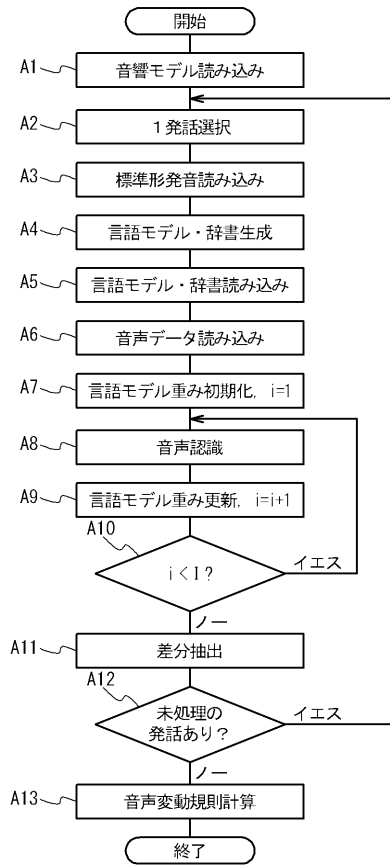
【図6】

| 言語モデル重み | 認識結果 |
|---------|--------------|
| 10.0 | #みなさん#こんにちわ# |
| ⋮ | ⋮ |
| 2.0 | #みなさん#こんにちわ# |
| 1.0 | #みんさん#こにつわ# |
| 0.5 | #みなさ#ごぬずば# |

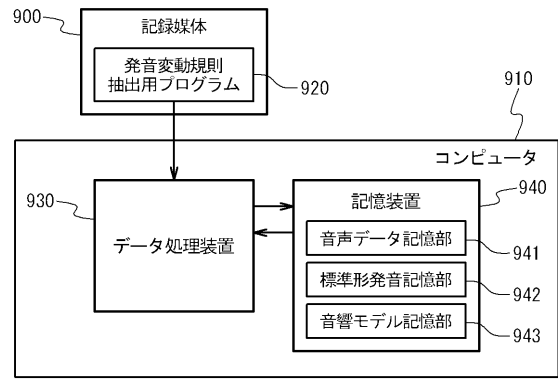
【図7】

| 標準形 | 変形 | 言語モデル重み | 備考 |
|-------|-------|---------|---------|
| みなさん | みなさん | 10.0 | 恒等変形 |
| こんにちわ | こんにちわ | 10.0 | 恒等変形 |
| みなさん | みなさん | 2.0 | 恒等変形 |
| こんにちわ | こにちわ | 2.0 | 脱落”ん” |
| みなさん | みんさん | 1.0 | 置換”な→ん” |
| こんにちわ | こにつわ | 1.0 | 置換”ち→つ” |
| ⋮ | ⋮ | ⋮ | ⋮ |

【図8】



【図9】



フロントページの続き

(56)参考文献 特開平08-123470(JP,A)

特開平10-308887(JP,A)

秋田祐哉 他, "話し言葉音声認識のための汎用的な統計的発音変動モデル", 電子情報通信学会論文誌D-II, 2005年9月1日, Vol.J88-D-II, No.9, p.1780-1789

塚田元, "重み付き有限状態トランスデューサの自動学習と発音変形のモデル化", 電子情報通信学会論文誌D-II, 2000年11月25日, Vol.J83-D-II, No.11, p.2457-2464

深田俊明 他, "発音ネットワークに基づく発音辞書の自動生成", 電子情報通信学会論文誌D-II, 1997年10月25日, Vol.J80-D-II, No.10, p.2626-2635

(58)調査した分野(Int.Cl., DB名)

G10L 11/00-25/93

JSTPlus(JDreamIII)