

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号

特許第6849621号
(P6849621)

(45) 発行日 令和3年3月24日(2021.3.24)

(24) 登録日 令和3年3月8日(2021.3.8)

(51) Int.Cl. F I
G I O L 15/16 (2006.01) G I O L 15/16

請求項の数 7 (全 21 頁)

(21) 出願番号	特願2018-17224 (P2018-17224)	(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30年2月2日(2018.2.2)	(74) 代理人	110002147 特許業務法人酒井国際特許事務所
(65) 公開番号	特開2019-133084 (P2019-133084A)	(72) 発明者	小川 厚徳 東京都千代田区大手町一丁目5番1号 日 本電信電話株式会社内
(43) 公開日	令和1年8月8日(2019.8.8)	(72) 発明者	デルクロア マーク 東京都千代田区大手町一丁目5番1号 日 本電信電話株式会社内
審査請求日	令和1年12月18日(2019.12.18)	(72) 発明者	苅田 成樹 東京都千代田区大手町一丁目5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】 学習装置、学習方法及び学習プログラム

(57) 【特許請求の範囲】

【請求項1】

精度が既知である学習用の複数の系列の入力を受け付ける入力部と、
前記複数の系列のうち二つの系列の特徴量がそれぞれ与えられたとき、それら二つの系列の精度の高低を判定できるような、ニューラルネットワークで表されるモデルを学習する学習部と、

を有し、

前記モデルは、二つの系列を、再帰的ニューラルネットワークを用いて隠れ状態ベクトルに変換し、ニューラルネットワークを用いて、前記隠れ状態ベクトルを基に二つの系列の精度の高低の並びが正しいことを示す第1の事後確率及び二つの系列の精度の高低の並びが誤りであることを示す第2の事後確率を出力することを特徴とする学習装置。

10

【請求項2】

前記入力部は、音声認識精度が既知である学習用のNベスト仮説の入力を受け付け、
前記学習部は、前記Nベスト仮説の二つの仮説のうち音声認識精度がより高い仮説に他方の仮説よりも高い順位が付与されている場合に正解ラベルを付与して前記モデルに学習させ、前記二つの仮説のうち音声認識精度がより高い仮説に他方の仮説よりも低い順位が付与されている場合に誤りラベルを付与して前記モデルに学習させることを特徴とする請求項1に記載の学習装置。

【請求項3】

前記二つの仮説のうち一方の仮説は、最も音声認識精度が高いオラクル仮説であるこ

20

とを特徴とする請求項 2 に記載の学習装置。

【請求項 4】

前記二つの仮説のうちの他方の仮説は、前記オラクル仮説の次に高い音声認識精度を持つ第 1 の仮説、N ベスト仮説における音声認識スコアが最も高い第 2 の仮説、最も低い音声認識精度を持つ第 3 の仮説、及び、N ベスト仮説における音声認識スコアが最も低い第 4 の仮説の少なくともいずれかを含むことを特徴とする請求項 3 に記載の学習装置。

【請求項 5】

前記二つの仮説のうちの他方の仮説は、N ベスト仮説から、前記オラクル仮説、前記第 1 の仮説、前記第 2 の仮説、前記第 3 の仮説及び前記第 4 の仮説を除いた仮説から所定のルールにしたがって抽出した所定数の仮説及び前記第 1 から第 4 の仮説であることを特徴とする請求項 4 に記載の学習装置。

10

【請求項 6】

学習装置が実行する学習方法であって、
精度が既知である学習用の複数の系列の入力を受け付ける工程と、
前記複数の系列のうちの二つの系列の特徴量がそれぞれ与えられたとき、それら二つの系列の精度の高低を判定できるような、ニューラルネットワークで表されるモデルを学習する工程と、

を含み、

前記モデルは、二つの系列を、再帰的ニューラルネットワークを用いて隠れ状態ベクトルに変換し、ニューラルネットワークを用いて、前記隠れ状態ベクトルを基に二つの系列の精度の高低の並びが正しいことを示す第 1 の事後確率及び二つの系列の精度の高低の並びが誤りであることを示す第 2 の事後確率を出力することを特徴とする学習方法。

20

【請求項 7】

精度が既知である学習用の複数の系列の入力を受け付けるステップと、
前記複数の系列のうちの二つの系列の特徴量がそれぞれ与えられたとき、それら二つの系列の精度の高低を判定できるような、ニューラルネットワークで表されるモデルを学習するステップと、

をコンピュータに実行させ、

前記モデルは、二つの系列を、再帰的ニューラルネットワークを用いて隠れ状態ベクトルに変換し、ニューラルネットワークを用いて、前記隠れ状態ベクトルを基に二つの系列の精度の高低の並びが正しいことを示す第 1 の事後確率及び二つの系列の精度の高低の並びが誤りであることを示す第 2 の事後確率を出力する学習プログラム。

30

【発明の詳細な説明】

【技術分野】

【0001】

本発明は、学習装置、学習方法及び学習プログラムに関する。

【背景技術】

【0002】

音声認識は、人間が発した音声（発話）を計算機により単語列（テキスト）に変換する技術である。通常、音声認識システムは、入力された一つの発話に対して、音声認識スコアの最も高い仮説（音声認識結果）である一つの単語列（1 ベスト仮説）を出力する。ただし、音声認識装置による音声認識の精度は、100%ではない。このため、一つの入力発話に対して、1 ベスト仮説のみを出力するのではなく、N（2）個の仮説を出力して、N ベストリスコアリング装置を用いて、そのN 個仮説の中から音声認識精度が最も高いと推定される仮説を最終的な音声認識結果として出力する、N ベストリスコアリングと呼ばれる手法がある。なお、N ベストリスコアリング（モデル）とN ベストリランキング（モデル）とは同義として扱われている。

40

【先行技術文献】

【非特許文献】

【0003】

50

【非特許文献1】T. Oba, T. Hori, A. Nakamura, and A. Ito, "Round-Robin Duel Discriminative Language Models", IEEE Transactions on Audio, Speech, and Language Processing, vol. 20, no.4, pp.1244 - 1255, May 2012.

【非特許文献2】A. Ogawa and T. Hori, "Error detection and accuracy estimation in automatic speech recognition using deep bidirectional recurrent neural networks", Speech Communication, vol. 89, pp.70 - 83, May 2017.

【発明の概要】

【発明が解決しようとする課題】

【0004】

図8は、Nベストリスコアリングの処理手順を示す図である。このNベストリスコアリングでは、まず、1発話の入力を受け付けると(ステップS31)、音声認識を行い(ステップS32)、音声認識結果であるN個の仮説を、これらの各スコアを基に降順でソートして出力する(ステップS33)。出力した仮説は、Nベスト仮説である。このNベスト仮説を用いて、Nベストリスコアリング処理が実行される(ステップS34)。Nベストリスコアリング処理では、後処理として、モデルを用いて、そのNベスト仮説に対して再度スコア付けを行う。モデルは、例えば、Nベストリスコアリングモデルを用いる。

【0005】

そして、Nベストリスコアリング処理では、それらの再度付与されたスコアを基にNベスト仮説を降順にソートするリランキングを行う。Nベストリスコアリング処理では、ソートの結果、Nベスト仮説の最上位にランキングされた仮説を抽出し、抽出した仮説を、最終的な音声認識結果として出力する。

【0006】

Nの値としては、通常、100~1000程度が設定されることが多い。なお、Nを2以上に設定しても、一つの仮説しか得られない場合もある。その場合は、Nベストリスコアリングを行う意味はない。

【0007】

図9は、Nベスト仮説の具体例を示す図である。図9では、Nを5以上に設定して音声認識を行い、5位までの五つの仮説が得られている。図9において、「[]」は、本来そこには単語がないことを示す。仮説は、音声認識スコアを基準として降順にソートされている。図9の例では、3位仮説が最も音声認識精度が高い(最も誤りが少ない)オラクル仮説である。Nベストリスコアリングの処理によって、この3位仮説が1位にリランキングされることが期待される。

【0008】

ここで、Nベストリスコアリングモデルとして、音声認識仮説である単語列の言語としての正しさに着目して、単語のつながり易さを確率統計的に表現する(評価する)言語モデルが採用されることが多い。単語列の言語としての正しさは、単語列の自然さ、単語同士のつながりの正しさである。

【0009】

従来、例えば、Nベスト仮説中の各仮説に含まれる単語のn連鎖(nは通常1~3程度)を素性として、ログリニアモデルに基づき、認識精度がより高い仮説により高いスコアを与えるような識別的言語モデルが盛んに研究されていた。

【0010】

また、近年では、ニューラルネットワーク(Neural Network: NN)の発展に基づき、NNに基づくNN言語モデルのうち、再帰的ニューラルネットワーク(Recurrent Neural Network: RNN)に基づくRNN言語モデルがNベストリスコアリングモデルとして盛んに利用されている。

【0011】

ここで、識別的言語モデルは、Nベスト仮説を用いて学習されるため、音声認識誤りを考慮したNベストリスコアリングを行うためのモデルであるものの、最新のNNに基づく

10

20

30

40

50

モデルではない。

【0012】

一方、RNN言語モデルは、最新のNNに基づくモデルであるものの、その学習は誤りを含まない正しい単語列を用いて行われるため、音声認識誤りを考慮できない。また、RNN言語モデルは高いNベストリスコアリング精度を示すものの、本来は単語列が与えられたときに、その単語列の次にどの単語が生じやすいかを推定するモデルである。すなわち、RNN言語モデルは、厳密には、Nベストリスコアリングを行うためのモデルではない。言い換えると、RNN言語モデルの次単語を推定する機能は、Nベストリスコアリングを行う上で必要な機能以上の機能であると言える。

【0013】

このように、ある一つの入力に対する解の候補として挙げられた複数の系列に対して、最も精度が高い（最も誤りが少ない）候補を判定するために、複数の系列に対して、最も精度が高い候補を判定するうえで最適なモデルを、識別的言語モデルやRNN言語モデルではなく、最新のNNに基づき実現することが期待されている。

【0014】

本発明は、上記に鑑みてなされたものであって、ある一つの入力に対する解の候補として挙げられた複数の系列に対し、最も精度が高い候補を判定する上で最適なモデルを実現する学習装置、学習方法及び学習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】

【0015】

上述した課題を解決し、目的を達成するために、本発明に係る学習装置は、精度が既知である学習用の複数の系列の入力を受け付ける入力部と、複数の系列のうちの二つの系列の特徴量がそれぞれ与えられたとき、それら二つの系列の精度の高低を判定できるように、ニューラルネットワークで表されるモデルを学習する学習部と、を有することを特徴とする。

【発明の効果】

【0016】

本発明によれば、ある一つの入力に対する解の候補として挙げられた複数の系列に対し、最も精度が高い候補を判定する上で最適なモデルを実現する。

【図面の簡単な説明】

【0017】

【図1】図1は、実施の形態に係るリランキング装置の機能構成の一例を示す図である。

【図2】図2は、Nベストリスコアリングモデルの構築例を示す図である。

【図3】図3は、図1に示すリランキング装置が実行するリランキング処理の処理手順を示すフローチャートである。

【図4】図4は、実施の形態に係る学習装置の機能構成の一例を示す図である。

【図5】図5は、図4に示す学習装置が実行する学習処理の処理手順を示すフローチャートである。

【図6】図6は、Nベストリランキングの評価結果を示す図である。

【図7】図7は、プログラムが実行されることにより、リランキング装置及び学習装置が実現されるコンピュータの一例を示す図である。

【図8】図8は、Nベストリスコアリングの処理手順を示す図である。

【図9】図9は、Nベスト仮説の具体例を示す図である。

【発明を実施するための形態】

【0018】

以下、図面を参照して、本発明の一実施形態を詳細に説明する。なお、この実施の形態により本発明が限定されるものではない。また、図面の記載において、同一部分には同一の符号を付して示している。本実施の形態では、ある正解に対する候補として挙げられた複数の系列として、音声認識結果であるN(N-2)ベスト仮説を例として説明する。そして、本実施の形態では、Nベスト仮説のうち、最終的な音声認識結果である最も音声認

10

20

30

40

50

識精度が高い仮説（単語列）を得るためのNベストランキングモデルを用いたランキング装置、及び、Nベストランキングモデルを実現する学習装置について説明する。なお、本実施の形態については、Nベストリスコアリング（モデル）ではなく、Nベストランキング（モデル）と表現を統一して説明する。

【0019】

まず、本実施の形態に係るランキング装置がNベスト仮説のランキングを行う上で、Nベストランキングモデルが有すべき必要最低限な機能について述べる。従来の方法では、Nベスト仮説は、リスコアリングの結果、スコアが降順になるようにソートされる。しかしながら、Nベストリスコアリングの主な目的は、Nベスト仮説から最も音声認識精度が高い仮説（オラクル仮説）を、最終的な音声認識結果として見つけ出すことである。このため、リスコアリング後のNベスト仮説は、必ずしもソートされている必要はない。本実施の形態ではこの点に着目した。

10

【0020】

すなわち、本実施の形態では、Nベスト仮説の中からオラクル仮説をランキングにより見つけ出すためにNベストランキングモデルに必要な最低限な機能は、Nベスト仮説中の二つの仮説に着目したときに、どちらの仮説の方がより高い音声認識精度を有しているかを判定できることである点に着目した。言い換えると、Nベストランキングモデルに必要な最低限な機能は、Nベスト仮説中の二つの仮説を対象に、一対一の仮説比較を行うことができることである。

【0021】

20

そこで、本実施の形態に係るランキング装置は、NNで表され、一対一の二つの仮説の比較を行う機能を持つNベストランキングモデルを用いることによって、二つの仮説のうち音声認識精度がより高い仮説を判定する機能を持たせた。そして、本実施の形態に係るランキング装置は、音声認識精度がより高い仮説を次の判定対象の一方の仮説として残し、未判定の仮説から他方の仮説を選択して、Nベストランキングモデルを用いた比較を行う。本実施の形態に係るランキング装置は、前回の判定で音声認識精度がより高いと判定された仮説を判定対象の一方の仮説として選択し、未判定の仮説のいずれかを他方の仮説として選択し、Nベストランキングモデルによる二つの仮説に対する比較処理を繰り返す。これによって、本実施の形態では、Nベスト仮説の中からオラクル仮説を見つけ出すことを可能にした。

30

【0022】

[実施の形態]

[ランキング装置]

次に、実施の形態に係るランキング装置について説明する。このランキング装置は、音声認識結果であるNベスト仮説のうち二つの仮説に対する、NNで表されるNベストランキングモデルを用いた音声認識精度の高低の判定を繰り返し実行して、最も音声認識精度の高い仮説を最終的な音声認識結果として出力する。

【0023】

図1は、実施の形態に係るランキング装置の機能構成の一例を示す図である。実施の形態1に係るランキング装置10は、例えば、ROM（Read Only Memory）、RAM（Random Access Memory）、CPU（Central Processing Unit）等を含むコンピュータ等に所定のプログラムが読み込まれて、CPUが所定のプログラムを実行することで実現される。

40

【0024】

ランキング装置10は、音声認識装置2から出力されたNベスト仮説の入力を受け付ける。そして、ランキング装置10は、このNベスト仮説のうち、二つの仮説に対する音声認識精度の高低についての判定を、全Nベスト仮説について実行し、音声認識精度が高い仮説として残った仮説を、最終的な音声認識結果として出力する。なお、音声認識装置2は、1発話が入力されると、例えば、音声認識用のモデルを用いて音声認識を行い、音声認識結果としてNベスト仮説を出力する。音声認識用のモデルは、学習用の複数の発

50

話と、各発話に対応する書き起こし（正解単語列）を学習データとして用いて学習（モデルパラメータが最適化）されている。

【0025】

リランキング装置10は、Nベストリランキングモデル記憶部11、仮説入力部12、仮説選択部13、特徴量抽出部14、判定部15、実行制御部16及び出力部17を有する。

【0026】

Nベストリランキングモデル記憶部11は、Nベストリランキングモデルを記憶する。Nベストリランキングモデルは、NNで表されるモデルである。Nベストリランキングモデルは、音声認識精度が既知である学習用のNベスト仮説を用いて予め学習される。Nベストリランキングモデルは、学習用のNベスト仮説のうち二つの系列の複数の組み合わせについて、二つの系列の特徴量が与えられたときに、その二つの系列の音声認識精度の高低を判定できるように学習される。Nベストリランキングモデルは、二つの仮説を、RNNを用いて隠れ状態ベクトルに変換する。そして、Nベストリランキングモデルは、NNを用いて、隠れ状態ベクトルを基に二つの仮説の精度の高低の並びが正しいことを示す第1の事後確率及び二つの仮説の精度の高低の並びが誤りであることを示す第2の事後確率を出力する。言い換えると、Nベストリランキングモデルでは、RNNの後段に2クラス分類FFNNが接続される。この2クラス分類FFNNは、RNNが変換した隠れ状態ベクトルを基に、二つの仮説のNベスト仮説における順位の上下関係が正しいことを示す第1の事後確率及び二つの仮説のNベスト仮説における順位の上下関係が誤りであることを示す第2の事後確率を出力する。

【0027】

仮説入力部12は、Nベスト仮説の入力を受け付ける。Nベスト仮説は、音声認識装置2が出力する。或いは、他の装置が、ネットワーク等を介して、Nベスト仮説をリランキング装置10に入力してもよい。

【0028】

仮説選択部13は、入力を受け付けたNベスト仮説のうち、一対一の比較対象である二つの仮説を選択する。仮説選択部13は、一定のルールに従い、Nベスト仮説の中から、任意の二つの仮説を1組とし選択する。具体的には、仮説選択部13は、二つの仮説の一方の仮説として、比較対象時に最高の精度を持つと推定される仮説を選択する。仮説選択部13は、二つの仮説の他方の仮説として、前回比較対象となった仮説の順位の次の順位の仮説を選択する。このように、仮説選択部13は、全Nベスト仮説について一対一の比較が実行されるように、Nベスト仮説から、比較対象の二つの仮説を選択する。

【0029】

特徴量抽出部14は、一対一の比較対象である二つの仮説について、それぞれの特徴量を抽出する。特徴量抽出部14は、一対一の比較対象であるNベスト仮説中のu位の仮説（単語列）と、Nベスト仮説中のv位（ $u < v \leq N$ ）の仮説とについて、それぞれの特徴量を抽出する。特徴量抽出部14は、仮説中の各単語単位で特徴量ベクトルを抽出する。各単語の特徴量ベクトルは、例えば、離散値である単語IDをNNによる単語の埋め込み処理により連続値のベクトルとして表現した単語ベクトルに、音声認識処理により得られる単語単位の音響スコア（対数尤度）や言語スコア（対数確率）などを補助特徴量として、単語ベクトルに連結したものである。

【0030】

判定部15は、一対一の比較対象の二つの仮説に対し、Nベストリランキングモデルを用いて、いずれの仮説がより高い音声認識精度を有しているかを判定する。具体的には、一対一の比較対象であるu位の仮説と、v（ $u < v \leq N$ ）位の仮説との特徴量をNベストリランキングモデルに入力し、Nベストリランキングモデルによる出力結果を用いて、どちらの仮説が高い音声認識精度を有しているかを判定する。u位及びv位で表す仮説の順位は、Nベスト仮説において既に付与されているものである。リランキング装置10では、順位の設定を行わない。

10

20

30

40

50

【 0 0 3 1 】

ここで、Nベストランキングモデルは、u位の仮説の特徴量及びv位の仮説の特徴量が入力されると、u位の仮説がv位の仮説よりも音声認識精度が高いことを示す第1の事後確率と、v位の仮説がu位の仮説よりも音声認識精度が高いことを示す第2の事後確率とを出力する。判定部15は、第1の事後確率が第2の事後確率よりも高い場合には、u位の仮説がv位の仮説よりも音声認識精度が高いと判定する。また、判定部15は、第1の事後確率が第2の事後確率よりも低い場合には、v位の仮説よりもu位の仮説よりも音声認識精度が高いと判定する。

【 0 0 3 2 】

なお、リランキング装置10では、特徴量抽出部14の機能を、Nベストランキングモデルが有してもよい。この場合、判定部15は、比較対象である二つの仮説をNベストランキングモデルに入力する。

【 0 0 3 3 】

そして、判定部15は、比較対象の二つの系列のうち、より精度が高いと判定した仮説を次の判定時における比較対象として残し、他方の仮説を以降比較対象から外す。仮説選択部13は、判定部15によって精度が高いと判定された仮説を二つの系列の一方の仮説として選択し、Nベスト仮説のうち、判定部15による判定が行われていない仮説のいずれかを他方の仮説として選択する。具体的には、前述したように、仮説選択部13は、判定部15が残した仮説を二つの仮説の一方の仮説として選択し、Nベスト仮説のうち、前回比較対象となった仮説の順位の次の順位の仮説を二つの仮説の他方の仮説として選択する。

【 0 0 3 4 】

実行制御部16は、判定部15による判定処理と仮説選択部14による選択処理とを、所定条件に達するまで繰り返す制御を行う。この場合、実行制御部16は、全Nベスト仮説について一対一の比較が実行されるように、仮説選択部13における比較対象の二つの仮説の選択処理、特徴量抽出部14における特徴量抽出処理、及び、判定部15における判定処理を繰り返す制御を行う。具体的に、実行制御部16は、比較対象である仮説の順位がNになるまで、仮説の選択処理、特徴量抽出処理及び判定処理を繰り返す制御を行う。

【 0 0 3 5 】

出力部17は、仮説の選択処理、特徴量抽出処理、判定処理及び順位の設定処理が繰り返された結果、Nベスト仮説のうち、所定条件に達した場合、比較対象として残っている仮説を、最も音声認識精度が高い仮説、すなわち、最終的な音声認識結果として出力する。出力部17は、最後の判定処理で精度が高いと判定された仮説を最終的な音声認識結果として出力する。

【 0 0 3 6 】

次に、Nベストランキングモデルに必要な最低限な機能要件を数式で定義する。 $W^{(u)} = w_1^{(u)}, w_2^{(u)}, \dots, w_{L(W^{(u)})}^{(u)}$ を、Nベスト仮説中のu位の仮説(単語列)と定義する。また、 $L(W^{(u)})$ を、 $W^{(u)}$ の長さ(単語数)と定義する。

【 0 0 3 7 】

また、 $A^{(u)} = a_1^{(u)}, a_2^{(u)}, \dots, a_{L(W^{(u)})}^{(u)}$ を $W^{(u)}$ に対応する補助特徴量ベクトル列と定義する。 $W^{(u)}$ 中のi番目の単語 $w_i^{(u)}$ の補助特徴量ベクトル $a_i^{(u)}$ は、例えば、音声認識装置による音声認識処理の結果として得られる音響スコア(対数尤度)や言語スコア(対数確率)などである(詳細は、例えば、非特許文献2を参照)。

【 0 0 3 8 】

また、 $X^{(u)} = x_1^{(u)}, x_2^{(u)}, \dots, x_{L(W^{(u)})}^{(u)}$ を $W^{(u)}$ に対応する特徴量ベクトル列と定義する。 $W^{(u)}$ 中のi番目の単語 $w_i^{(u)}$ の特徴量ベクトル $x_i^{(u)}$ は、 $x_i^{(u)} = \text{concat}(\text{embed}(w_i^{(u)}), a_i$

10

20

30

40

50

(u))で得られる。ここで、 $\text{concat}(\cdot)$ は、ベクトルの連結処理を表す。また、 $\text{embed}(\cdot)$ は、NNによる単語の埋め込み処理（離散値の単語IDを連続値のベクトルで表現する処理）（詳細は、例えば、坪井祐太，海野裕也，鈴木潤，深層学習による自然言語処理，MLP機械学習プロフェッショナルシリーズ，講談社，2017．（以降、参考文献1とする。）を参照）を表す。なお、 $\text{embed}(\cdot)$ を行うNNもNベストリスコアリングモデルの一部であり、そのパラメータは、後述のエンコーダRNN及び2クラス分類FFNNのパラメータと同時に学習（最適化）される。

【0039】

Nベスト仮説中の u 番目の仮説 $W^{(u)}$ と v 番目の仮説 $W^{(v)}$ ($u < v \leq N$)の特徴量ベクトル列 $X^{(u)}$ ， $X^{(v)}$ が与えられたとき、リランキング装置10におけるNベストリランキングモデルは、2クラスの記号 $y = \{0, 1\}$ の事後確率 P を出力する。 $y = 0$ は、 $W^{(u)}$ 及び仮説 $W^{(v)}$ の順位の上下関係が正しいことを示す。また、 $y = 1$ は、 $W^{(u)}$ 及び仮説 $W^{(v)}$ の順位の上下関係が誤りであることを示す。 $P(0 | X^{(u)}, X^{(v)})$ は、 u 位の仮説と v 位の仮説との順位の上下関係が正しさを確率的に表現する第1の事後確率である。 $P(1 | X^{(u)}, X^{(v)})$ は、 u 位の仮説と v 位の仮説との順位の上下関係が誤りであることを確率的に表現する第2の事後確率である。

10

【0040】

判定部15は、Nベストリランキングモデルから出力された第1の事後確率 $P(0 | X^{(u)}, X^{(v)})$ 及び第2の事後確率 $P(1 | X^{(u)}, X^{(v)})$ を取得し、取得した二つの事後確率の大小を比較して、 u 位の仮説及び v 位の仮説のいずれがより音声認識精度が高いかを判定する。判定部15は、第1の事後確率 $P(0 | X^{(u)}, X^{(v)})$ が第2の事後確率 $P(1 | X^{(u)}, X^{(v)})$ よりも高い場合には、 u 位の仮説が v 位の仮説よりも音声認識精度が高いと判定する。また、判定部15は、第1の事後確率 $P(0 | X^{(u)}, X^{(v)})$ が第2の事後確率 $P(1 | X^{(u)}, X^{(v)})$ よりも低い場合には、 v 位の仮説が u 位の仮説よりも音声認識精度が高いと判定する。

20

【0041】

すなわち、判定部15は、以下の(1-1)式及び(1-2)式に示すように、 u 位の仮説及び v 位の仮説のいずれがより音声認識精度が高いかを判定する。

【0042】

$$P(0 | X^{(u)}, X^{(v)}) > P(1 | X^{(u)}, X^{(v)})$$

$$\text{if } \text{acc}(W^{(u)}) > \text{acc}(W^{(v)}) \quad \dots (1-1)$$

$$P(0 | X^{(u)}, X^{(v)}) < P(1 | X^{(u)}, X^{(v)})$$

$$\text{otherwise} \quad \dots (1-2)$$

30

【0043】

ここで、 $\text{acc}(\cdot)$ は、与えられた仮説（単語列）の音声認識精度を返す関数 $y = P(y | X^{(u)}, X^{(v)}) = 1$ である。(1-1)式の1段目に示す不等式が満足される場合、判定部15は、仮説 $W^{(u)}$ は仮説 $W^{(v)}$ 以上の音声認識精度を持つと判定する。また、(1-2)式の不等式が満足される場合、判定部15は、 $W^{(u)}$ は $W^{(v)}$ よりも低い音声認識精度を持つと判定する。

【0044】

したがって、(1-1)式の1段目に示す不等式が満足される場合、 $W^{(u)}$ 及び $W^{(v)}$ のランキングの上下関係 ($u < v$) が正しいと推定される。このため、判定部15は、 $W^{(u)}$ を、 $W^{(v)}$ との一対一の仮説比較において $W^{(v)}$ よりも音声認識精度が高い仮説として残し、次の一対一の仮説比較でも $W^{(u)}$ として引き続き使用する。なお、判定部15は、 $W^{(v)}$ を、 $W^{(u)}$ よりも音声認識精度が低い仮説として扱い、最も音声認識精度が高い仮説の候補、すなわち、最終的な音声認識結果の候補から除外する。

40

【0045】

そして、(1-2)式の1段目不等式が満足される場合は、 $W^{(u)}$ 及び $W^{(v)}$ のランキングの上下関係は、誤りであると推定される。すなわち、 $W^{(u)}$ 及び $W^{(v)}$ のランキングの上下関係は逆であると推定される。このため、判定部15は、 $W^{(v)}$ を、 W

50

(u) との一対一の仮説比較において $W^{(u)}$ よりも音声認識精度が高い仮説として残し、次の一対一の仮説比較では $W^{(u)}$ として使用する。なお、判定部 15 は、元の $W^{(u)}$ を、元の $W^{(v)}$ よりも音声認識精度が低い仮説として扱い、最も音声認識精度が高い仮説の候補、すなわち、最終的な音声認識結果の候補から除外する。なお、N ベストリランキングモデルは、第 1 の事後確率 $P(0 | X^{(u)}, X^{(v)})$ 及び第 2 の事後確率 $P(1 | X^{(u)}, X^{(v)})$ の事後確率の大小を比較して、 u 位の仮説及び v 位の仮説のいずれがより音声認識精度が高いかを判定し、仮説の残存の判定までを推定してもよい。

【0046】

[N ベストリランキングモデルの構築例]

図 2 は、N ベストリランキングモデルの構築例を示す図である。なお、図 2 では、簡単のため、単語の埋め込み処理 $embed(\cdot)$ を行う NN は省略されている。以下、その詳細について説明する。

10

【0047】

比較対象の仮説 $W^{(u)}$ の長さ (単語数) $L(W^{(u)})$ と仮説 $W^{(v)}$ ($u < v \leq N$) の長さ $L(W^{(v)})$ とが異なる可能性がある。この長さの違いを吸収するため、N ベストリランキングモデルは、二つの仮説を、RNN を用いて隠れ状態ベクトルに変換する。具体的には、N ベストリランキングモデルは、この処理を行うために、エンコーダ-デコーダモデル (詳細は、例えば、参考文献 1 参照) のエンコーダ RNN 111 を有する。

【0048】

N ベストリランキングモデルは、エンコーダ RNN 111 を用いて $W^{(u)}$ と $W^{(v)}$ を固定長の隠れ状態ベクトルで表現することができる。そして、N ベストリランキングモデルは、これらの隠れ状態ベクトルを用いることによって、 $W^{(u)}$ と $W^{(v)}$ とを公平に比較することが可能になる。

20

【0049】

エンコーダ RNN 111 の処理について説明する。エンコーダ RNN 111 は、RNN の一種である長短期記憶メモリ (long short-term memory: LSTM) ユニット (詳細は、例えば、参考文献 1 参照) を有する。LSTM ユニットは、 $W^{(u)}$ の i 番目の単語 $w_i^{(u)}$ の特徴量ベクトル $x_i^{(u)}$ と、 $i-1$ 番目の隠れ状態ベクトル $h_{\{i-1\}}^{(u)}$ が与えられたとき、 i 番目の隠れ状態ベクトル $h_i^{(u)}$ を以下の (2) 式のように与える。

30

【0050】

$$h_i^{(u)} = \text{lstm}(x_i^{(u)}, h_{\{i-1\}}^{(u)}) \quad \dots (2)$$

【0051】

ここで、 $\text{lstm}(\cdot)$ は、1 層単方向 (unidirectional) の LSTM ユニットの処理を示す。また、 $h_i^{(u)} = 0$ (ゼロベクトル) である。 $h_i^{(u)}$ は、単語列 $w_1^{(u)}, w_2^{(u)}, \dots, w_i^{(u)}$ の特徴量ベクトル列 $x_1^{(u)}, x_2^{(u)}, \dots, x_i^{(u)}$ をエンコード (符号化) したものである。エンコーダ RNN 111 は、この処理を、特徴量ベクトル列 $X^{(u)}$ 中の各特徴量ベクトル $x_i^{(u)}$ に対して繰り返すことで、 $X^{(u)}$ をエンコードした隠れ状態ベクトル $h_{L(W^{(u)})}^{(u)}$ を得ることができる。

40

【0052】

エンコーダ RNN 111 は、同様の処理を特徴量ベクトル列 $X^{(v)}$ に対しても行い、 $X^{(v)}$ をエンコードした隠れ状態ベクトル $h_{L(W^{(v)})}^{(v)}$ を得る。なお、 $X^{(u)}$ に対して処理を行う LSTM ユニットと、 $X^{(v)}$ に対して処理を行う LSTM ユニットは同じもの、すなわち、パラメータが共有されていてもよいし、別の LSTM ユニットであってもよい。また、図 2 では、 $x_{L(W^{(u)})}^{(u)}, x_{L(W^{(v)})}^{(v)}, h_{L(W^{(u)})}^{(u)}, h_{L(W^{(v)})}^{(v)}$ の下付き部分 $L(W^{(u)})$ は、 $L(W^{(u)})$ と示している。

【0053】

N ベストリランキングモデルは、以上で得た二つの隠れ状態ベクトル $h_{L(W^{(u)})}^{(u)}$

50

(u) , h L (w (v)) (v) を連結した隠れ状態ベクトル h { (u , v) } をエンコーダ RNN 1 1 の出力として以下の (3) 式のように得る。

【 0 0 5 4 】

$$h \{ (u , v) \} = \text{concat} (h_{L (w (u))} (u) , h_{L (w (v))} (v)) \dots (3)$$

【 0 0 5 5 】

そして、Nベストリランキングモデルは、エンコーダ RNN 1 1 1 の後段に、クラス分類 (y = 0 or 1) を行うための NN を連結する。例えば、Nベストリランキングモデルは、2クラス分類のための NN として、1層のフィードフォワード型 NN (FFNN) 1 1 2 (詳細は、例えば、参考文献 1 を参照) を用いる。エンコーダ RNN 1 1 1 の出力として得た隠れ状態ベクトル h { (u , v) } が、1層の2クラス分類 FFNN に入力され、最終的に、2クラスの記号 y = { 0 , 1 } の事後確率 P (y | X (u) , X (v)) を以下の (4) , (5) 式のように得ることができる。

【 0 0 5 6 】

$$z \{ (u , v) \} = \text{linear} (h \{ (u , v) \}) \dots (4)$$

$$P (y | X (u) , X (v)) = \text{softmax} (z \{ (u , v) \})_y \dots (5)$$

【 0 0 5 7 】

ここで、linear (·) は、線形変換処理 (詳細は、例えば、参考文献 1 を参照) を表す。softmax (·) は、ソフトマックス処理を表す。また、softmax (·)_y は、ソフトマックス処理の結果として得られる事後確率ベクトルの y 番目の要素 (確率値) を表す。

【 0 0 5 8 】

[Nベストリランキングモデルの他の構築例 1]

なお、図 2 に示すエンコーダ RNN 1 1 1 の LSTM ユニットは、1層単方向の LSTM ユニットとしたが、複数層または双方向 (bidirectional) の LSTM ユニットであってもよい。

【 0 0 5 9 】

[Nベストリランキングモデルの他の構築例 2]

また、LSTM ユニットの代わりに、単純な (下記の sigmoid 関数等を活性化関数として持つ。) RNN や、Gated Recurrent Unit (GRU) を用いてもよい。

【 0 0 6 0 】

[Nベストリランキングモデルの他の構築例 3]

さらに、Nベストリランキングモデルは、図 2 の構築例では、2クラス分類 NN として、1層のフィードフォワード型 NN を用いたが、複数層のフィードフォワード型 NN を用いてもよい。Nベストリランキングモデルは、複数層のフィードフォワード型 NN を用いる場合、活性化関数として、sigmoid 関数、tanh 関数、Rectified Linear Unit (ReLU) 関数、Parametric ReLU (PReLU) 関数などを用いることができる。なお、Nベストリランキングモデルの他の構築例 1 ~ 3 の用語の詳細については、例えば、参考文献 1 を参照いただきたい。

【 0 0 6 1 】

[Nベストリランキングモデルの他の構築例 4]

また、Nベストリランキングモデルは、従来の Nベストリスコアリングモデル (例えば RNN 言語モデル) により計算されたスコアを、特徴量ベクトルにおける新たな次元として追加して利用することも可能である。

【 0 0 6 2 】

[リランキング処理の処理手順]

次に、図 1 に示すリランキング装置 1 0 が実行するリランキング処理の処理手順について説明する。図 3 は、図 1 に示すリランキング装置 1 0 が実行するリランキング処理の処

10

20

30

40

50

理手順を示すフローチャートである。

【0063】

まず、仮説入力部12が、リランキング対象のNベスト仮説の入力を受け付けると(ステップS1)、仮説選択部13は、入力を受け付けたNベスト仮説のうち、順次、一対一の比較対象であるu位及びv位の二つの仮説を選択する($u < v \leq N$)。まず、仮説選択部13は、 $u = 1$ 、 $v = 2$ に設定する(ステップS2)。そして、仮説選択部13は、入力を受け付けたNベスト仮説から、u位及びv位の二つの仮説 $W^{(u)}$ 、 $W^{(v)}$ をNベスト仮説から選択する(ステップS3)。続いて、特徴量抽出部14は、仮説 $W^{(u)}$ 、 $W^{(v)}$ の特徴量を抽出する(ステップS4)。判定部15は、仮説 $W^{(u)}$ 、 $W^{(v)}$ の特徴量($X^{(u)}$ 、 $X^{(v)}$)をNベストリランキングモデルに入力する(ステップS

10

【0064】

判定部15は、Nベストリランキングモデルからの出力結果を取得する(ステップS6)。具体的には、判定部15は、第1の事後確率 $P(0 | X^{(u)}, X^{(v)})$ 及び第2の事後確率 $P(1 | X^{(u)}, X^{(v)})$ を取得する。

【0065】

そして、(1-1)式及び(1-2)式において説明したように、判定部15は、 $acc(W^{(u)}) > acc(W^{(v)})$ であるか否かを判定する(ステップS7)。判定部15は、 $P(0 | X^{(u)}, X^{(v)}) > P(1 | X^{(u)}, X^{(v)})$ の場合、 $acc(W^{(u)}) > acc(W^{(v)})$ であると判定する。一方、判定部15は、 $P(0 | X^{(u)}, X^{(v)}) < P(1 | X^{(u)}, X^{(v)})$ の場合、 $acc(W^{(u)}) < acc(W^{(v)})$ でないと判定する。

20

【0066】

判定部15が $acc(W^{(u)}) > acc(W^{(v)})$ であると判定した場合(ステップS7: Yes)、順位設定部16は、kについて $k = u$ と設定する(ステップS8)。kは、最も音声認識精度が高い仮説のNベスト仮説における順位(ランキング)である。一方、判定部15が $acc(W^{(u)}) < acc(W^{(v)})$ でないと判定した場合(ステップS7: No)、順位設定部16は、 $k = v$ と設定する(ステップS9)。

【0067】

続いて、実行制御部16は、 $v = N$ であるか否かを判定する(ステップS10)。実行制御部16は、 $v = N$ でないと判定した場合(ステップS10: No)、必要な一対一の仮説比較処理がまだ全ては終了していないため、仮説選択部13に対し、比較対象の次の仮説の選択を行わせる。具体的には、仮説選択部13は、 $u = k$ 、 $v = v + 1$ に設定し(ステップS11)、ステップS3に戻り、次の判定対象のNベスト仮説 $W^{(u)}$ 、 $W^{(v)}$ を選択する。そして、リランキング装置10は、このNベスト仮説 $W^{(u)}$ 、 $W^{(v)}$ に対して、ステップS4~ステップS10の処理を実行する。

30

【0068】

また、実行制御部16は、 $v = N$ であると判定した場合(ステップS10: Yes)、必要な一対一の比較処理が全て終了したため、k位の $W^{(k)}$ を最も音声認識精度が高いと推定される仮説、すなわち、最終的な音声認識結果として出力し(ステップS12)、処理を終了する。このように、リランキング装置10では、任意の二つの仮説を1組とし、複数の組についてそれぞれ音声認識精度の高低の判定を繰り返すことで、最も音声認識精度が高いと推定される仮説を、最終的な音声認識結果として出力することができる。

40

【0069】

[学習装置]

次に、リランキング装置10が用いるNベストリランキングモデルを学習する学習装置について説明する。図4は、実施の形態に係る学習装置の機能構成の一例を示す図である。実施の形態1に係る学習装置20は、例えば、ROM、RAM、CPU等を含むコンピュータ等に所定のプログラムが読み込まれて、CPUが所定のプログラムを実行することで実現される。図4に示すように、Nベストリランキングモデル記憶部21、学習装置2

50

0 は、仮説入力部 2 2 及び学習部 2 3 を有する。

【 0 0 7 0 】

N ベストリランキングモデル記憶部 2 1 は、学習対象の N ベストリランキングモデルを記憶する。N ベストリランキングモデルは、NN で表される。N ベストリランキングモデルは、N ベスト仮説のうち二つの仮説を、RNN を用いて隠れ状態ベクトルに変換する。そして、N ベストリランキングモデルは、NN を用いて、隠れ状態ベクトルを基に二つの仮説の精度の高低の並びが正しいことを示す第 1 の事後確率及び二つの仮説の精度の高低の並びが誤りであることを示す第 2 の事後確率を出力する。

【 0 0 7 1 】

仮説入力部 2 2 は、音声認識精度が既知である学習用の N ベスト仮説の入力を受け付ける。学習用の N ベスト仮説として、学習データ中の各発話に対して音声認識が行われ、各発話の N ベスト仮説が得られているものとする。また学習データであるので、全ての仮説の音声認識精度は、既知である。また、N ベスト仮説中の全ての仮説に対して、前述のように、特徴量ベクトル列が抽出されているものとする。

10

【 0 0 7 2 】

学習部 2 3 は、学習用の N ベスト仮説のうち二つの仮説の特徴量がそれぞれ与えられたときに、それら二つの仮説の精度の高低が判定できるような、N ベストリランキングモデルを学習する。学習部 2 3 では、学習用の N ベスト仮説のうち二つの仮説の特徴量ベクトル列と、これらに対応する教師ラベル（後述）とを、N ベストリランキングモデルに与える。これによって、学習部 2 3 は、N ベストリランキングモデルがこれら二つの仮説の音声認識精度の高低を正しく判定できるように、N ベストリランキングモデルの学習（パラメータの最適化）を行う。具体的には、学習部 2 3 は、特徴量ベクトル列と、対応する教師ラベルとを N ベストリランキングモデルに入力し、N ベストリランキングモデルがこれらの特徴量ベクトルを与えられたときに対応する教師ラベルを正しく出力できるように、N ベストリランキングモデルの学習を行う。学習部 2 3 は、教師ラベル付与部 2 3 1 及び入替部 2 3 2 を有する。

20

【 0 0 7 3 】

教師ラベル付与部 2 3 1 は、二つの仮説のうち音声認識精度がより高い仮説に他方の仮説よりも高い順位が付与されている場合に正解を表す教師ラベル ($y = 0$) を付与して、N ベストリランキングモデルに学習させる。また、教師ラベル付与部 2 3 1 は、二つの仮説のうち音声認識精度がより高い仮説に他方の仮説よりも低い順位が付与されている場合に誤りを表す教師ラベル ($y = 1$) を付与し、N ベストリランキングモデルに学習させる。

30

【 0 0 7 4 】

入替部 2 3 2 は、学習用の N ベスト仮説のうち二つの仮説の順位を入れ換え、対応する教師ラベルも入れ換えて、N ベストリランキングモデルの学習を行う。例えば、教師ラベルとして $y = 0$ が付与されている二つの仮説については、二つの仮説の順位を入れ換え、教師ラベル y を 1 に変える。一方、教師ラベルとして $y = 1$ が付与されている二つの仮説については、二つの仮説の順位を入れ換え、教師ラベル y を 0 に変える。

【 0 0 7 5 】

[学習処理の処理手順]

次に、図 4 に示す学習装置 2 0 が実行する学習処理の処理手順について説明する。図 5 は、図 4 に示す学習装置が実行する学習処理の処理手順を示すフローチャートである。図 5 では、N ベスト仮説から二つの仮説として $W^{(u)}$, $W^{(v)}$ ($u < v \leq N$) が与えられ、かつ、 $acc(W^{(u)}) > acc(W^{(v)})$ であるときの学習処理の処理手順を示す。

40

【 0 0 7 6 】

図 5 に示すように、教師ラベル付与部 2 3 1 が、教師ラベル $y = 0$ を付与し（ステップ S 2 1）、 $W^{(u)}$, $W^{(v)}$ の特徴量 $X^{(u)}$, $X^{(v)}$ を N ベストリランキングモデルに入力し（ステップ S 2 2）、N ベストリランキングモデルを学習させ、N ベストリラ

50

ンキングモデルのモデルパラメータを更新させる（ステップS23）。すなわち、（1-1）式に従うと、この二つの仮説の $W^{(u)}$ 、 $W^{(v)}$ の特徴量ベクトル $X^{(u)}$ 、 $X^{(v)}$ をNベストリランキングモデルに入力した場合、モデルは、理想的には、 $P(0 | X^{(u)}, X^{(v)}) = 1$ との事後確率を出力すべきである。このため、教師ラベル付与部231は、教師ラベルとして、 $y = 0$ を与える。以上の入力を基に、学習部23は、モデルパラメータ（エンコーダRNN（LSTMユニット）及び2クラス分類FFNN及び単語の埋め込み処理 $embed(\cdot)$ を行うNNのパラメータを同時に）を更新させる。

【0077】

そして、入替部232は、仮説 $W^{(u)}$ 、 $W^{(v)}$ の順位を入れ替える（ステップS24）。すなわち、元々、 $W^{(v)}$ であった仮説を $W^{(u)}$ とし、元々、 $W^{(u)}$ であった仮説を $W^{(v)}$ とする。この場合には、 $acc(W^{(u)}) > acc(W^{(v)})$ ではない。よって、（1-2）式に従えば、この二つの仮説 $W^{(u)}$ 、 $W^{(v)}$ の特徴量ベクトル $X^{(u)}$ 、 $X^{(v)}$ をNベストリランキングモデルに入力した場合、モデルは理想的には、 $P(1 | X^{(u)}, X^{(v)}) = 1$ との事後確率を出力すべきである。このため、教師ラベル付与部231は、教師ラベルとして、 $y = 1$ を付与し（ステップS25）、 $W^{(u)}$ 、 $W^{(v)}$ の特徴量 $X^{(u)}$ 、 $X^{(v)}$ をNベストリランキングモデルに入力する（ステップS26）。学習部23は、以上の入力を基に、Nベストリランキングモデルを学習させ、Nベストリランキングモデルのモデルパラメータを更新させて（ステップS27）、二つの仮説 $W^{(u)}$ 、 $W^{(v)}$ に対する学習処理を終了する。

【0078】

学習装置20は、上記の手順を、学習データ中の各発話のNベスト仮説について繰り返し、更にはその繰り返し自体を何度か（何エポックか）繰り返す。学習部23は、学習の更なる具体的な手順については、従来のNNの学習（詳細は、例えば、参考文献1参照）と同様に行うことができる。

【0079】

[学習処理の効率化例1]

図5に示す学習処理の処理手順は、計算コストが高い。例えば、Eをエポック数、Mを学習データ中の発話数とすると、上記の学習手順におけるモデルパラメータの更新回数は、最大で、 $E \times M \times N \times 2 \times N C_2$ になる。通常、Eは数十程度、Mは少なくとも数万、Nは上記の通り100～1000程度であるので、モデルパラメータの更新回数は、膨大な数に達する。このため、本実施の形態では、学習の効率化を図ることが好ましい。そこで、以下に、学習の効率化例1について述べる。

【0080】

上述したように、Nベストラスコアリングの主な目的は、Nベスト仮説からオラクル仮説を最終的な音声認識結果として見つけ出すことである。言い換えれば、オラクル仮説をその他のN-1個の仮説から精度よく区別できればよい。これを実現するために、学習の際に、Nベストリランキングモデルに入力する二つの仮説のうち的一方をオラクル仮説とする。これにより、モデルパラメータの更新回数を、 $E \times M \times N \times 2 \times (N - 1)$ に削減することができる。

【0081】

[学習処理の効率化例2]

次に、学習の効率化例2について説明する。学習の効率化例1では、Nベスト仮説が与えられたとき、その中に含まれるオラクル仮説とその他のN-1個の仮説とを比較していた。学習処理の効率化例2では、オラクル仮説と比較するその他の仮説の個数を絞り込む。

【0082】

例えば、まず、下の典型的な四つの仮説を選択する。

仮説1は、オラクル仮説の次に高い音声認識精度を持つ仮説である。

仮説2は、音声認識スコアが最も高い仮説である。

仮説3は、最も低い音声認識精度を持つ仮説である。

10

20

30

40

50

仮説 4 は、音声認識スコアが最も低い仮説である。

【 0 0 8 3 】

仮説 1 と仮説 2 とは、音声認識精度が高い（または高いと推定される）仮説で、オラクル仮説との区別が難しい仮説である。一方、仮説 3 と仮説 4 とは、音声認識精度が低い（または低いと推定される）仮説で、オラクル仮説との区別が容易な（確実に区別しないとはいけない）仮説である。その他の仮説をこの四つのみに絞り込む場合は、モデルパラメータの更新回数は、 $E \times M \times N \times 2 \times 4$ にまで削減することができる。

【 0 0 8 4 】

ただし、上記の四つの仮説のみではオラクル仮説の対立仮説としての多様性が十分に確保できないと考えられる場合、 N ベスト仮説から、オラクル仮説とこれらの四つの仮説を除いた、残りの $N - 5$ 個の仮説から、所定のルールにしたがって抽出した所定数の仮説を選択して前記四つの仮説と共に対立仮説として用いてもよい。例えば、二つの仮説のうち他方の仮説として、オラクル仮説とこれらの四つの仮説を除いた、残りの $N - 5$ 個の仮説から、等間隔に、或いは、はランダムに、 Q 個の仮説を選択して四つの仮説と共に他方の仮説として用いる。このとき、モデルパラメータの更新回数は、 $E \times M \times N \times 2 \times (4 + Q)$ となる。例えば、 Q は、 $5 \sim 50$ である。

【 0 0 8 5 】

[評価]

実際に、本実施の形態における N ベストリスコアリングと、従来の RNN 言語モデルとの比較評価を行った。 N ベストリランキングモデルを使用する際（評価時）、モデルは、 $(1 - 1)$ 式及び $(1 - 2)$ 式にしたがい、2 クラスの事後確率 $P(y | X^{(u)}, X^{(v)})$, $y = \{0, 1\}$ を推定する。リランキング装置 10 は、これらの事後確率をそのまま用いて N ベストリランキングを行ってもよい。また、リランキング装置 10 は、従来の N ベストリランキングモデルと同様に、 (6) 式を用いて、元々の音声認識スコアと N ベストリランキングモデルによるスコア（事後確率の対数値）とを重み付け加算し、その値を基に、 N ベストリランキングを行ってもよい。

【 0 0 8 6 】

スコア = $(1 - \alpha) \times$ 音声認識スコア + $\alpha \times N$ ベストリランキングモデルによるスコア
 $\dots (6)$

【 0 0 8 7 】

なお、 (6) 式において、 α は、 N ベストリランキングモデルの重みであり、 $0 < \alpha < 1$ である。リランキング装置 10 は、 $\alpha = 1$ に設定した場合は、音声認識スコアを用いず、 N ベストリランキングモデルによるスコアのみを用いて、 N ベストリランキングを行う。

【 0 0 8 8 】

[評価結果]

図 6 は、 N ベストリランキングの評価結果を示す図である。図 6 では、 N ベストリランキングの評価結果の例として、日本語話し言葉コーパスを用いて、従来の RNN 言語モデル、本実施の形態のリランキング装置 10 が用いる N ベストリランキングモデル、及び、 RNN 言語モデルのスコアを特徴量ベクトルの次元として加えた N ベストリランキングモデル（他の構築例 4）を比較評価した結果を示す。評価は、 (6) 式に従い、音声認識スコアと N ベストリランキングモデルによるスコアとを重み加算したスコアを用いて N ベストリランキングを行っている。

【 0 0 8 9 】

図 6 に示すように、従来の RNN 言語モデルよりも、本実施の形態の N ベストリランキングモデルの方が着実に音声認識精度を改善できることが分かる。また、構築例 4 のように、 RNN 言語モデルのスコアを特徴量ベクトルの次元として加えることによって、 N ベストリランキングモデルの音声認識精度をさらに改善できることが分かる。

【 0 0 9 0 】

さらに、図 6 より、従来の RNN 言語モデルでは、音声認識スコアを使用する必要があ

10

20

30

40

50

り、かつ、重み に比較的狭い最適値があることが分かる。本評価では、従来の RNN 言語モデルは、 $\alpha = 0.8$ 付近である。一方、本実施の形態の N ベストランキングモデルを用いた場合、従来の RNN 言語モデルの場合と比較して、 α の最適値の範囲が広いことが分かる。すなわち、本実施の形態の N ベストランキングモデルは、 α の値に頑健である。或いは、本実施の形態の N ベストランキングモデルを用いた場合、 $\alpha = 1$ において最高か最高に近い音声認識精度が得られているので、音声認識スコアを使用しなくてもよいことが分かる。

【0091】

[実施の形態の効果]

本実施の形態に係るランキング装置 10 では、音声認識結果である N ベスト仮説の入力を受け付け、N ベスト仮説中の二つの仮説に対し、NN で表される N ベストランキングモデルを用いて、いずれの仮説がより高い音声認識精度を有しているかを判定する。

10

【0092】

前述したように、ランキング装置 10 が N ベスト仮説のランキングを行う上で、N ベストランキングモデルが有すべき必要最低限の機能は、N ベスト仮説から最も高精度な仮説（オラクル仮説）を、最終的な音声認識結果として見つけ出すことである。このため、リスコアリング後の N ベスト仮説は、必ずしもソートされている必要はない。

【0093】

そこで、本実施の形態では、N ベスト仮説の中からオラクル仮説をランキングにより見つけ出すために、N ベストランキングモデルに、N ベスト仮説中の二つの仮説のうちどちらの仮説の方がより高い音声認識精度を有しているかを判定できる機能を持たせた。言い換えると、本実施の形態では、N ベストランキングモデルに、N ベスト仮説中の二つの仮説を対象に、一対一の仮説比較を実行できる機能を持たせた。

20

【0094】

具体的には、ランキング装置 10 は、NN で表され、一対一の二つの仮説の比較を行う機能を持つ N ベストランキングモデルを用い、N ベストランキングモデルを用いた一対一の二つの仮説に対する比較処理を繰り返すことによって、N ベスト仮説の中からオラクル仮説を見つけて出すことを可能にしている。

【0095】

さらに、学習装置 20 は、N ベストランキングモデルに、音声認識精度が既知である学習用の N ベスト仮説のうちの二つの仮説を 1 組として、複数の組についてそれぞれ音声認識精度の高低を判定できるように予め学習させている。したがって、学習装置 20 は、N ベストランキングを行う上で最適なモデルを、最新の NN に基づき実現することができる。そして、ランキング装置 10 は、学習装置 20 において学習された N ベストランキングモデルを使用することによって、一対一の二つの仮説の比較を精度よく行うことができ、このランキング装置 10 によるオラクル仮説の抽出を高精度にできる。

30

【0096】

このように、本実施の形態によれば、N ベスト仮説の中から最終的な音声認識結果を、精度よく得ることができる、NN で表された N ベストランキングモデルを実現することができる。そして、本実施によれば、N ベストランキングモデルを用いることによって、最終的な音声認識結果を精度よく得ることができる。

40

【0097】

なお、本実施の形態では、一対一の仮説比較をオラクル仮説（最も精度が高い仮説）と推定される仮説が見つかった時点で処理を終了していたが、オラクル仮説と推定される仮説を除いた $N - 1$ 個の仮説に対して、オラクル仮説を見つけるのと同様の処理を行うことで、二番目に精度が高いと推定される仮説を見つけて出すことができる。以降、この処理を繰り返すことによって、N ベスト仮説のソートも可能である。

【0098】

また、本実施の形態では、音声認識の N ベスト仮説をランキングするためのモデルとして、図 2 に例示する N ベストランキングモデルについて説明した。ただし、本実施の

50

形態のモデルは、音声認識のNベスト仮説への適用にとどまらず、Nベスト仮説を採用しているあらゆるタスクに適用可能である。例えば、機械翻訳や文章要約などにも本実施の形態を適用することが可能である。また、文字列に限らず、数字やアルファベットを含む複数の系列にも本実施の系列を適用することが可能である。

【0099】

このため、本実施の形態は、ある一つの入力に対する解の候補として挙げられた複数の系列であれば、このうちの二つの系列に対し、NNで表されるモデルを用いて、二つの系列のうちより精度が高い（誤りが少ない）系列を判定できる。そして、本実施の形態では、二つの系列のうち、より精度が高いと判定した系列を比較対象として残し、他方の系列を比較対象から外し、精度が高いと判定した系列を二つの系列の一方の仮説として選択し、複数の系列のうち、判定が行われていない系列のいずれかを他方の仮説として選択する。そして、本実施の形態では、判定処理と選択処理とを、所定条件に達するまで順次実行させせる。これによって、本実施の形態によれば、所定条件に達した場合に比較対象として残っている系列を、最も精度が高い系列、すなわち、最終的な出力として出力することができる。

10

【0100】

また、この場合には、本実施の形態では、精度が既知である学習用の複数の系列のうちの二つの系列の特徴量が与えられたとき、それら二つの系列の精度の高低が判定できるような、NNで表されるモデルを学習する。そして、本実施の形態では、二つの系列のうち精度がより高い（誤りがより少ない）系列に他方の系列よりも高い順位が付与されている場合に正解を示す教師ラベルを付与してモデルに学習させる。そして、本実施の形態では、二つの系列のうち精度がより高い（誤りがより少ない）系列に他方の系列よりも低い順位が付与されている場に誤りを示す教師ラベルを付与してモデルに学習させる。本実施の形態では、このモデルによって、一対一の二つの系列の比較が高精度で行うことができ、この結果、最も精度の高い系列を精度よく得ることができる。

20

【0101】

[システム構成等]

図示した各装置の各構成要素は機能概念的なものであり、必ずしも物理的に図示の如く構成されていることを要しない。すなわち、各装置の分散・統合の具体的形態は図示のものに限られず、その全部又は一部を、各種の負荷や使用状況等に応じて、任意の単位で機能的又は物理的に分散・統合して構成することができる。例えば、リランキング装置10及び学習装置20は、一体の装置であってもよい。さらに、各装置にて行なわれる各処理機能は、その全部又は任意の一部が、CPU及び当該CPUにて解析実行されるプログラムにて実現され、あるいは、ワイヤードロジックによるハードウェアとして実現され得る。

30

【0102】

また、本実施形態において説明した各処理のうち、自動的に行われるものとして説明した処理の全部又は一部を手動적으로おこなうこともでき、あるいは、手動적으로おこなわれるものとして説明した処理の全部又は一部を公知の方法で自動的におこなうこともできる。また、本実施形態において説明した各処理は、記載の順にしたがって時系列に実行されるのみならず、処理を実行する装置の処理能力あるいは必要に応じて並列的あるいは個別に実行されてもよい。この他、上記文書中や図面中で示した処理手順、制御手順、具体的な名称、各種のデータやパラメータを含む情報については、特記する場合を除いて任意に変更することができる。

40

【0103】

[プログラム]

図7は、プログラムが実行されることにより、リランキング装置10或いは学習装置20が実現されるコンピュータの一例を示す図である。コンピュータ1000は、例えば、メモリ1010、CPU1020を有する。また、コンピュータ1000は、ハードディスクドライブインタフェース1030、ディスクドライブインタフェース1040、シリ

50

アルポートインタフェース 1050、ビデオアダプタ 1060、ネットワークインタフェース 1070 を有する。これらの各部は、バス 1080 によって接続される。

【0104】

メモリ 1010 は、ROM 1011 及び RAM 1012 を含む。ROM 1011 は、例えば、BIOS (Basic Input Output System) 等のブートプログラムを記憶する。ハードディスクドライブインタフェース 1030 は、ハードディスクドライブ 1031 に接続される。ディスクドライブインタフェース 1040 は、ディスクドライブ 1041 に接続される。例えば磁気ディスクや光ディスク等の着脱可能な記憶媒体が、ディスクドライブ 1041 に挿入される。シリアルポートインタフェース 1050 は、例えばマウス 1110、キーボード 1120 に接続される。ビデオアダプタ 1060 は、例えばディスプレイ 1130 に接続される。

10

【0105】

ハードディスクドライブ 1031 は、例えば、OS 1091、アプリケーションプログラム 1092、プログラムモジュール 1093、プログラムデータ 1094 を記憶する。すなわち、リランキング装置 10 或いは学習装置 20 の各処理を規定するプログラムは、コンピュータ 1000 により実行可能なコードが記述されたプログラムモジュール 1093 として実装される。プログラムモジュール 1093 は、例えばハードディスクドライブ 1031 に記憶される。例えば、リランキング装置 10 或いは学習装置 20 における機能構成と同様の処理を実行するためのプログラムモジュール 1093 が、ハードディスクドライブ 1031 に記憶される。なお、ハードディスクドライブ 1031 は、SSD (Solid State Drive) により代替されてもよい。

20

【0106】

また、上述した実施形態の処理で用いられる設定データは、プログラムデータ 1094 として、例えばメモリ 1010 やハードディスクドライブ 1031 に記憶される。そして、CPU 1020 が、メモリ 1010 やハードディスクドライブ 1031 に記憶されたプログラムモジュール 1093 やプログラムデータ 1094 を必要に応じて RAM 1012 に読み出して実行する。

【0107】

なお、プログラムモジュール 1093 やプログラムデータ 1094 は、ハードディスクドライブ 1031 に記憶される場合に限らず、例えば着脱可能な記憶媒体に記憶され、ディスクドライブ 1041 等を介して CPU 1020 によって読み出されてもよい。あるいは、プログラムモジュール 1093 及びプログラムデータ 1094 は、ネットワーク (LAN (Local Area Network)、WAN (Wide Area Network) 等) を介して接続された他のコンピュータに記憶されてもよい。そして、プログラムモジュール 1093 及びプログラムデータ 1094 は、他のコンピュータから、ネットワークインタフェース 1070 を介して CPU 1020 によって読み出されてもよい。

30

【0108】

以上、本発明者によってなされた発明を適用した実施形態について説明したが、本実施形態による本発明の開示の一部をなす記述及び図面により本発明は限定されることはない。すなわち、本実施形態に基づいて当業者等によりなされる他の実施形態、実施例及び運用技術等は全て本発明の範疇に含まれる。

40

【符号の説明】

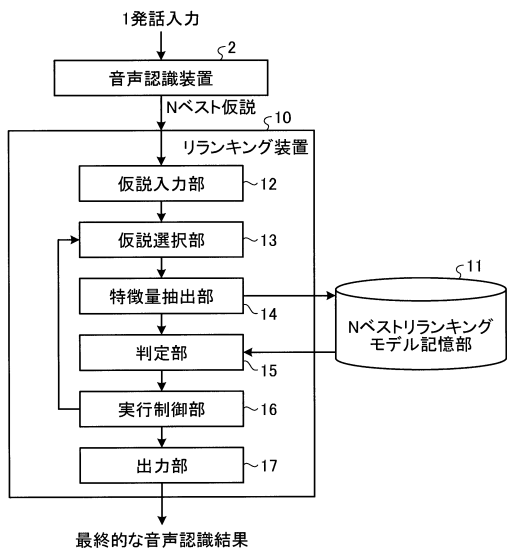
【0109】

- 2 音声認識装置
- 10 リランキング装置
- 11, 21 Nベストリランキングモデル記憶部
- 12 仮説入力部
- 13 仮説選択部
- 14 特徴量抽出部
- 15 判定部

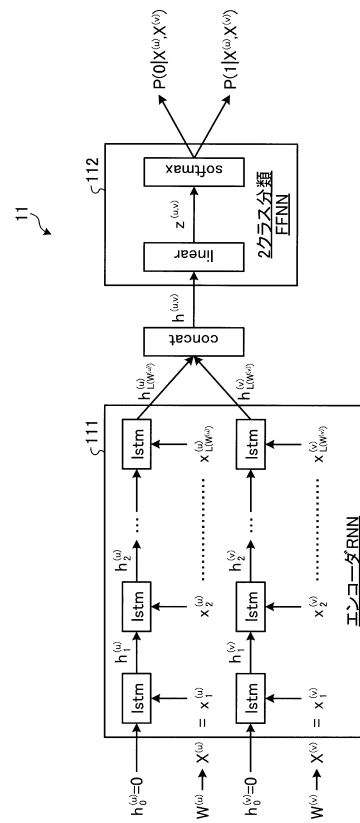
50

- 1 6 実行制御部
- 1 7 出力部
- 2 0 学習装置
- 2 2 仮説入力部
- 2 3 学習部
- 2 3 1 教師ラベル付与部
- 2 3 2 入替部

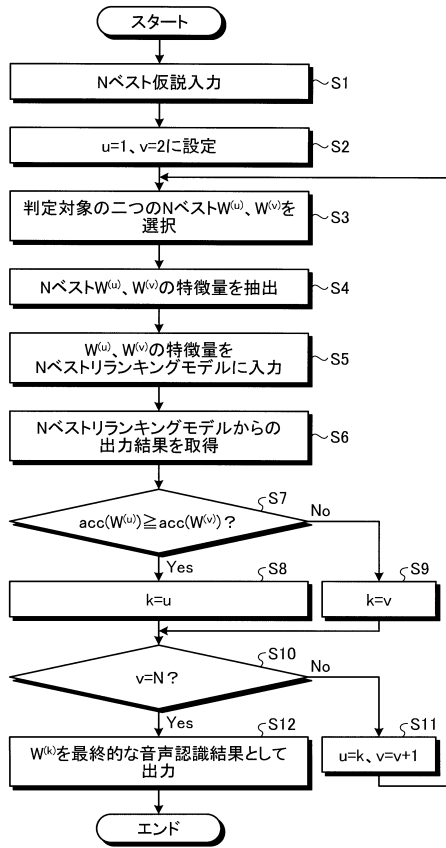
【図1】



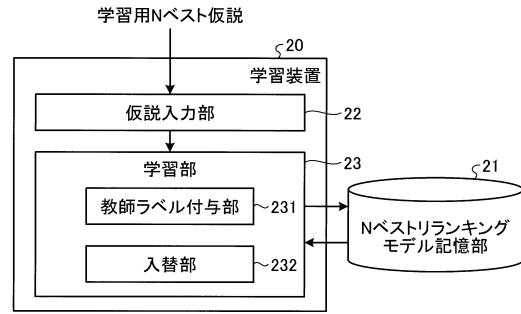
【図2】



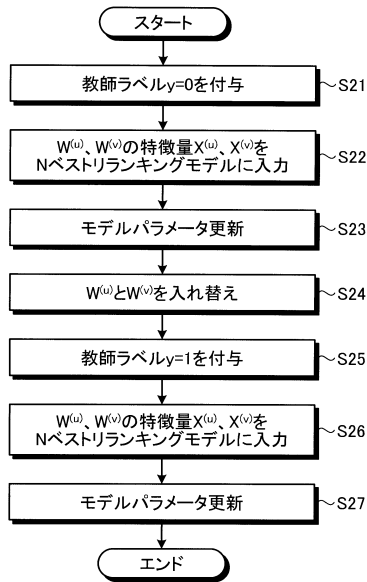
【図3】



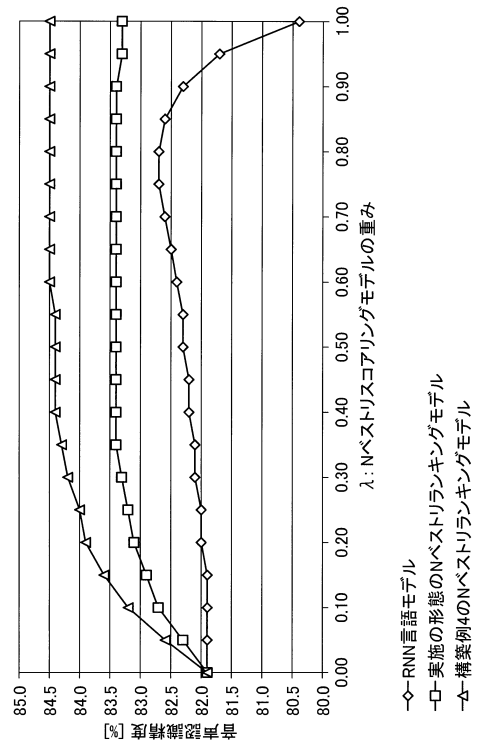
【図4】



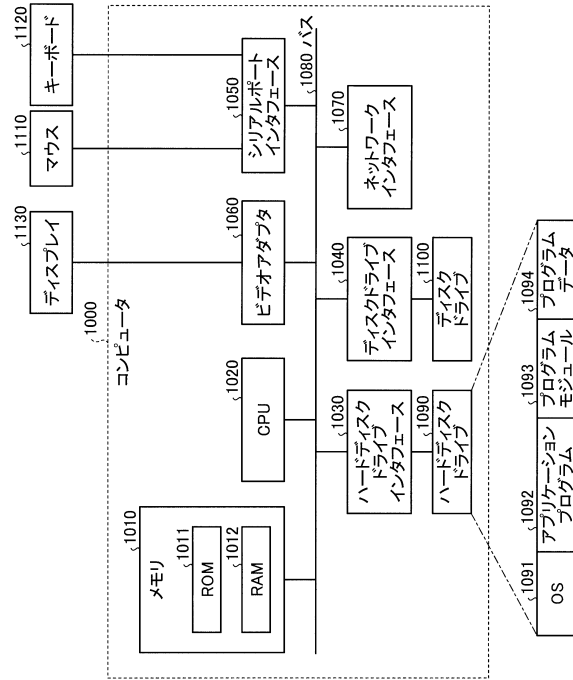
【図5】



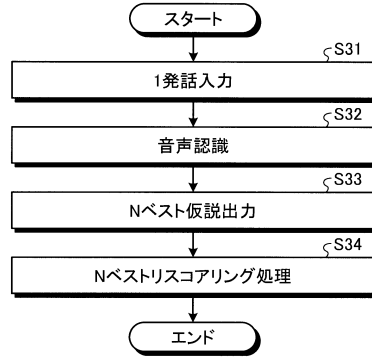
【図6】



【 図 7 】



【 図 8 】



【 図 9 】

発話(正解単語列): 私と友人は2018年4月に名古屋大学に入学します □

1位仮説(認識結果): 私と求人 2018年7月に名古屋大学に入学します □ 120.3

2位仮説(認識結果): 職と友人は2018年4月に名古屋大学に入学します □ 115.8

3位仮説(認識結果): 私と友人は2018年4月に名古屋大学に入学します □ 108.5

4位仮説(認識結果): 職と友人は2017年4月に名古屋大学に入学します □ 101.4

5位仮説(認識結果): 私と求人 2018年7月に名古屋大学に入学しました 98.2

↑
音声認識
スコア

フロントページの続き

(72)発明者 中谷 智広

東京都千代田区大手町一丁目5番1号 日本電信電話株式会社内

審査官 中村 天真

(56)参考文献 特開2011-243147(JP,A)

岩立将和,外2名,トーナメントモデルを用いた日本語係り受け解析,自然言語処理,2008年10月,第15巻,第5号,p.169-185

島岡聖世,外2名,オートエンコーダにおける単語ベクトルの学習,言語処理学会第19回年次大会発表論文集,2013年3月,p.612-615

小川厚徳,外3名,一対一の仮説比較を行うencoder-classifierモデルを用いたNベスト音声認識仮説のリスコアリング,日本音響学会講演論文集,2018年3月,p.23-24

(58)調査した分野(Int.Cl.,DB名)

G10L 15/00 - 15/34

G06F 40/20 - 40/58