

(19)日本国特許庁(JP)

(12)特許公報(B2)

(11)特許番号
特許第7192995号
(P7192995)

(45)発行日 令和4年12月20日(2022.12.20)

(24)登録日 令和4年12月12日(2022.12.12)

(51)国際特許分類 F I
G 1 0 L 15/08 (2006.01) G 1 0 L 15/08 2 0 0 Z
G 1 0 L 15/16 (2006.01) G 1 0 L 15/16

請求項の数 10 (全26頁)

(21)出願番号	特願2021-537548(P2021-537548)	(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
(86)(22)出願日	令和1年8月8日(2019.8.8)	(74)代理人	110002147 弁理士法人酒井国際特許事務所
(86)国際出願番号	PCT/JP2019/031517	(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
(87)国際公開番号	WO2021/024491	(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
(87)国際公開日	令和3年2月11日(2021.2.11)	(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
審査請求日	令和3年11月19日(2021.11.19)	(72)発明者	中谷 智広

最終頁に続く

(54)【発明の名称】 判定装置、学習装置、判定方法及び判定プログラム

(57)【特許請求の範囲】

【請求項1】

音声認識精度のスコアが対応付けられたNベスト仮説の入力を受け付ける入力部と、
入力を受け付けた前記Nベスト仮説のうち、判定対象である二つの仮説を選択する選択部と、

選択された二つの仮説が与えられたとき、前記二つの仮説を隠れ状態ベクトルに変換し、
前記二つの仮説の隠れ状態ベクトルを基に前記二つの仮説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複数の補助モデルと、前記複数の補助モデルでそれぞれ
変換された前記二つの仮説の隠れ状態ベクトルを基に、前記二つの仮説の精度の高低を
判定できるような、ニューラルネットワークで表されるメインモデルとを用いて、
前記二つの仮説の精度の高低を判定する判定部と、

を有することを特徴とする判定装置。

【請求項2】

前記選択部は、前記Nベスト仮説のスコアの昇順に前記二つの仮説を選択することを特徴とする請求項1に記載の判定装置。

【請求項3】

前記判定部は、前記メインモデルから出力された判定情報、または、各補助モデルから出力された判定情報と前記メインモデルから出力された判定情報とに対して計算した重み付け和の値、に基づいて前記二つの仮説の精度の高低を判定することを特徴とする請求項1または2に記載の判定装置。

【請求項 4】

各補助モデルは、前記二つの仮説を、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事後確率を出力し、

前記メインモデルは、ニューラルネットワークを用いて、前記複数の補助モデルでそれぞれ変換された前記二つの仮説の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事後確率を出力する

ことを特徴とする請求項 1～3 のいずれか一つに記載の判定装置。

【請求項 5】

音声認識精度が既知である学習用の二つの仮説の入力を受け付ける入力部と、

前記二つの仮説が与えられたとき、前記二つの仮説を隠れ状態ベクトルに変換し、前記二つの仮説の隠れ状態ベクトルを基に前記二つの仮説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複数の補助モデルと、前記複数の補助モデルでそれぞれ変換された前記二つの仮説の隠れ状態ベクトルを基に、前記二つの仮説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メインモデルとに対し、各ニューラルネットワークが前記二つの仮説の精度の高低を判定するタスクを個別に行うとみなしたマルチタスク学習を行わせる学習部と、

を有することを特徴とする学習装置。

【請求項 6】

前記学習部は、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複数の補助モデル及び前記メインモデルに学習させ、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複数の補助モデル及び前記メインモデルに学習させることを特徴とする請求項 5 に記載の学習装置。

【請求項 7】

前記学習部は、各ニューラルネットワークによって実行された各タスクについて所定の損失をそれぞれ計算し、各損失の重み付け和に基づいて、各ニューラルネットワークのパラメータの値を更新することを特徴とする請求項 5 または 6 に記載の学習装置。

【請求項 8】

各補助モデルは、前記二つの仮説を、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事後確率を出力し、

前記メインモデルは、ニューラルネットワークを用いて、前記複数の補助モデルでそれぞれ変換された前記二つの仮説の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事後確率を出力する

ことを特徴とする請求項 5～7 のいずれか一つに記載の学習装置。

【請求項 9】

判定装置が実行する判定方法であって、

音声認識精度のスコアが対応付けられた N ベスト仮説の入力を受け付ける工程と、

入力を受け付けた前記 N ベスト仮説のうち、判定対象である二つの仮説を選択する工程と、

選択された二つの仮説が与えられたとき、前記二つの仮説を隠れ状態ベクトルに変換し、前記二つの仮説の隠れ状態ベクトルを基に前記二つの仮説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複数の補助モデルと、前記複数の補助モデルでそれぞれ変換された前記二つの仮説の隠れ状態ベクトルを基に、前記二つの仮説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メインモデルとを用いて、前記二つの仮説の精度の高低を判定する工程と、

を含んだことを特徴とする判定方法。

【請求項 10】

音声認識精度のスコアが対応付けられた N ベスト仮説の入力を受け付けるステップと、

	10
	20
	30
	40
	50

入力を受け付けた前記Nベスト仮説のうち、判定対象である二つの仮説を選択するステップと、

選択された二つの仮説が与えられたとき、前記二つの仮説を隠れ状態ベクトルに変換し、前記二つの仮説の隠れ状態ベクトルを基に前記二つの仮説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複数の補助モデルと、前記複数の補助モデルでそれぞれ変換された前記二つの仮説の隠れ状態ベクトルを基に、前記二つの仮説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メインモデルとを用いて、前記二つの仮説の精度の高低を判定するステップと、

をコンピュータに実行させるための判定プログラム。

【発明の詳細な説明】

【技術分野】

【0001】

本発明は、判定装置、学習装置、判定方法及び判定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とNベストリランキングとは同義として扱われている。

【0003】

Nベストリスコアリング方法では、音声認識結果である仮説の中からスコアの高い所定数（N個）の仮説を出力する。そして、Nベストリスコアリング方法では、この中から尤もらしい仮説を音声認識結果として出力する。ここで、スコアが最大となる仮説が必ずしもベストな仮説とは限らない。このため、二つの仮説のうち尤もらしい仮説（正解に近い仮説）を選択する二択問題をトーナメント方式で繰り返し適用することで、尤もらしい仮説を選択するリランキング装置が提案されている（例えば、非特許文献1参照）。

【先行技術文献】

【非特許文献】

【0004】

【文献】Atsunori Ogawa, Marc Delcroix, Shigeki Karita, Tomohiro Nakatani, "RESCORING N-BEST SPEECH RECOGNITION LIST BASED ON ONE-ON-ONE HYPOTHESIS COMPARISON USING ENCODER-CLASSIFIER MODEL", IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP), pp. 6099-6103, 2018.

【発明の概要】

【発明が解決しようとする課題】

【0005】

非特許文献1に記載のリランキング方法では、N仮説をスコアの降順に並べ、先頭の仮説（スコアが最も高い仮説）から順に二つの仮説を選択し、学習済みの二択問題を解くニューラルネットワーク（NN）にこれらの仮説を入力することによって、いずれかの仮説を選択する処理を繰り返し行い、最終的に選択された仮説を音声認識結果として出力することが記載されている。非特許文献1に記載のリランキング方法では、一定の精度で音声認識結果を出力するが、さらに、近年では、音声認識結果の出力に対して、精度の安定化が要求されている。

【0006】

10

20

30

40

50

本発明は、上記に鑑みてなされたものであって、ある音声信号に対する解の候補として挙げられた複数の仮説に対し、最も精度が高い仮説を安定した精度で判定することができる判定装置、学習装置、判定方法及び判定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】

【0007】

上述した課題を解決し、目的を達成するために、本発明に係る判定装置は、音声認識精度のスコアが対応付けられたNベスト仮説の入力を受け付ける入力部と、入力を受け付けたNベスト仮説のうち、判定対象である二つの仮説を選択する選択部と、選択された二つの仮説が与えられたとき、二つの仮説を隠れ状態ベクトルに変換し、二つの仮説の隠れ状態ベクトルを基に二つの仮説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複数の補助モデルと、複数の補助モデルでそれぞれ変換された二つの仮説の隠れ状態ベクトルを基に、二つの仮説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メインモデルとを用いて、二つの仮説の精度の高低を判定する判定部と、を有することを特徴とする。

10

【0008】

また、本発明に係る学習装置は、音声認識精度が既知である学習用の二つの仮説の入力を受け付ける入力部と、二つの仮説が与えられたとき、二つの仮説を隠れ状態ベクトルに変換し、二つの仮説の隠れ状態ベクトルを基に二つの仮説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複数の補助モデルと、複数の補助モデルでそれぞれ変換された二つの仮説の隠れ状態ベクトルを基に、二つの仮説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メインモデルとに対し、各ニューラルネットワークが二つの仮説の精度の高低を判定するタスクを個別に行うとみなしたマルチタスク学習を行わせる学習部と、を有することを特徴とする。

20

【0009】

また、本発明に係る判定方法は、判定装置が実行する判定方法であって、音声認識精度のスコアが対応付けられたNベスト仮説の入力を受け付ける工程と、入力を受け付けたNベスト仮説のうち、判定対象である二つの仮説を選択する工程と、選択された二つの仮説が与えられたとき、二つの仮説を隠れ状態ベクトルに変換し、二つの仮説の隠れ状態ベクトルを基に二つの仮説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複数の補助モデルと、複数の補助モデルでそれぞれ変換された二つの仮説の隠れ状態ベクトルを基に、二つの仮説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メインモデルとを用いて、二つの仮説の精度の高低を判定する工程と、を含んだことを特徴とする。

30

【0010】

また、本発明に係る判定プログラムは、音声認識精度のスコアが対応付けられたNベスト仮説の入力を受け付けるステップと、入力を受け付けたNベスト仮説のうち、判定対象である二つの仮説を選択するステップと、選択された二つの仮説が与えられたとき、二つの仮説を隠れ状態ベクトルに変換し、二つの仮説の隠れ状態ベクトルを基に二つの仮説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複数の補助モデルと、複数の補助モデルでそれぞれ変換された二つの仮説の隠れ状態ベクトルを基に、二つの仮説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メインモデルとを用いて、二つの仮説の精度の高低を判定するステップと、をコンピュータに実行させる。

40

【発明の効果】

【0011】

本発明によれば、ある音声信号に対する解の候補として挙げられた複数の仮説に対し、最も精度が高い仮説を、安定した精度で判定することができる。

【図面の簡単な説明】

【0012】

【図1】図1は、実施の形態1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。

50

【図 2】図 2 は、第 1 補助モデル～第 M 補助モデル及びメインモデルの構成を説明する図である。

【図 3】図 3 は、第 1 補助モデルの構築例を示す図である。

【図 4】図 4 は、実施の形態 1 に係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 5】図 5 は、図 1 に示すリランキング装置が、N ベスト仮説に対して実行するリランキング処理を説明する図である。

【図 6】図 6 は、実施の形態 2 に係る学習装置の機能構成の一例を示す図である。

【図 7】図 7 は、図 6 に示す入替部の処理を説明する図である。

【図 8】図 8 は、実施の形態 2 に係る学習処理の処理手順を示すフローチャートである。

10

【図 9】図 9 は、実施の形態 3 に係るリランキング装置の要部構成を示す図である。

【図 10】図 10 は、実施の形態 3 に係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 11】図 11 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。

【発明を実施するための形態】

【0013】

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。

20

【0014】

本実施の形態では、音声認識結果である N (N ≥ 2) 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識精度が高い仮説 (単語列) を得るためのモデルを用いたリランキング装置、及び、N ベストのリランキング処理に用いる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ではなく、N ベストリランキングと表現を統一して説明する。

【0015】

まず、本実施の形態に係るリランキング装置が N ベスト仮説のリランキングを行う上で、本実施の形態におけるモデルが有すべき必要最低限な機能について述べる。本実施の形態では、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が目的である。

30

【0016】

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リランキングにより見つけ出すために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、本実施の形態における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。

【0017】

そこで、本実施の形態に係るリランキング装置は、一対一の二つの仮説の比較を行う機能を持つ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。さらに、本実施の形態では、モデルとして、ニューラルネットワーク (NN) で表されるメインモデルと、NN で表される複数の補助モデルとを用いる。各補助モデルは、二つの仮説が与えられたとき、二つの仮説を隠れ状態ベクトルに変換し、二つの仮説の隠れ状態ベクトルを基に二つの仮説の精度の高低を判定するモデルである。メインモデルは、複数の補助モデルでそれぞれ変換された二つの仮説の隠れ状態ベクトルを基に、二つの仮説の精度の高低を判定するモデルである。

40

【0018】

そして、本実施の形態に係るリランキング装置は、N ベスト仮説のスコアの昇順に二つの仮説を選択し、選択した二つの仮説のうち、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、未判定の仮説から昇順に他方の仮説を選択して、複数の補助

50

モデル及びメインモデルを用いた比較を行う。本実施の形態に係るリランキング装置は、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、未判定の仮説のうち最も順位の低い仮説を他方の仮説として選択し、複数の補助モデル及びメイン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、安定した精度で、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。

【0019】

[実施の形態1]

[リランキング装置]

まず、実施の形態1に係るリランキング装置について説明する。このリランキング装置は、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して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。

10

【0020】

図1は、実施の形態1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。実施の形態1に係るリランキング装置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。

20

【0021】

リランキング装置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、リランキング装置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、スコアと対応付けて、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。

30

【0022】

リランキング装置10は、モデル記憶部11、仮説入力部12、仮説選択部13(選択部)、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。

40

【0023】

モデル記憶部11は、補助モデル及びメインモデル110を記憶する。図1の例では、モデル記憶部11は、補助モデルとして、第1補助モデル111~第M補助モデル11Mを記憶する。第1補助モデル111~第M補助モデル11M及びメインモデル110は、NNで表されるモデルである。第1補助モデル111~第M補助モデル11M及びメインモデル110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。

【0024】

第1補助モデル111~第M補助モデル11Mは、選択された二つの仮説が与えられたとき、二つの仮説を隠れ状態ベクトルに変換し、二つの仮説の隠れ状態ベクトルを基に二つの仮説の精度の高低を判定できるような、NNで表される。第1補助モデル111~第M補助モデル11Mは、学習用のNベスト仮説のうち二つの仮説が与えられたときに、二つの仮説について、その二つの仮説の音声認識精度の高低を判定できるように学習される。第1補助モデル111~第M補助モデル11Mは、二つの仮説を、それぞれRNNを用いて隠れ状態ベクトルに変換する。そして、第1補助モデル111~第M補助モデル11M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事後確率をそれぞれ生成する。

50

【0025】

メインモデル110は、第1補助モデル111~第M補助モデル11Mにおいてそれぞ

50

れ変換された二つの仮説の隠れ状態ベクトルを基に、二つの仮説の精度の高低を判定できるような、NNで表される。メインモデル110は、第1補助モデル111～第M補助モデル11Mにおいてそれぞれ変換された学習用の二つの仮説の隠れ状態ベクトルを基に、二つの仮説の精度の高低を判定できるように学習される。メインモデル110は、NNを用いて、第1補助モデル111～第M補助モデル11Mにおいてそれぞれ変換された学習用の二つの仮説の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事後確率を生成する。

【0026】

第1補助モデル111～第M補助モデル11M及びメインモデル110は、各ニューラルネットワークが二つの仮説の精度の高低を判定するタスクを個別に行うとみなしたマルチタスク学習によって学習が実行される。

10

【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、Nベスト仮説のスコアの昇順に選択する。仮説選択部13は、最初の判定においては、Nベスト仮説のうち、スコアが最下位である仮説と、最下位の仮説より1つ順位が高い仮説とを判定対象として選択する。仮説選択部13は、以降の判定においては、二つの仮説の一方の仮説として、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選択する。そして、仮説選択部13は、二つの仮説の他方の仮説として、未判定の仮説のうち、最もスコアの順位が低い仮説を選択する。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、昇順に、比較対象の二つの仮説を選択する。

20

【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v位の仮説とNベスト仮説中のu ($u < v \leq N$)位の仮説(単語列)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。

30

【0030】

判定部15は、一対一の比較対象の二つの仮説に対し、第1補助モデル111～第M補助モデル11M及びメインモデル110を用いて、二つの仮説の精度の高低を判定する。判定部15は、一対一の比較対象であるv位の仮説とu位の仮説とを第1補助モデル111～第M補助モデル11Mにそれぞれ入力し、メインモデル110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。ランキング装置10では、順位の再設定を行わない。

40

【0031】

ここで、第1補助モデル111～第M補助モデル11M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事後確率を出力する。メインモデル110は、第1補助モデル111～第M補助モデル11Mにおいてそれぞれ変換された二つの仮説の隠れ状態ベクトルが入力されると、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事後確率を出力する。判定部15は、メインモデル110による事後確率が0.5以上である場合には、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、メインモデ

50

ル 1 1 0 による事後確率が 0 . 5 未満である場合には、 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。

【 0 0 3 2 】

なお、リランキング装置 1 0 では、特徴量抽出部 1 4 の機能を、第 1 補助モデル 1 1 1 ~ 第 M 補助モデル 1 1 M が有してもよい。この場合、判定部 1 5 は、比較対象である二つの仮説を第 1 補助モデル 1 1 1 ~ 第 M 補助モデル 1 1 M にそれぞれ入力する。

【 0 0 3 3 】

そして、判定部 1 5 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、他方の仮説を以降の比較対象から外す。仮説選択部 1 3 は、判定部 1 5 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、判定部 1 5 による判定が行われていない仮説のうち最もスコアの順位が低い仮説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 1 3 は、判定部 1 5 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 1 つ上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。

10

【 0 0 3 4 】

実行制御部 1 6 は、判定部 1 5 による判定処理と仮説選択部 1 3 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 1 6 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 1 3 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 1 4 における特徴量抽出処理、及び、判定部 1 5 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 1 6 は、1 位の仮説に対して比較処理が行われるまで、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。

20

【 0 0 3 5 】

出力部 1 7 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、所定条件に達した場合、 N ベスト仮説のうち、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 1 7 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 】

[定義]

まず、リランキング装置 1 0 に必要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。

30

【 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、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. (以降、参考文献 1 とする。) を参照) 。

40

【 0 0 3 8 】

補助特徴量ベクトル $a_i^{(u)}$ は、1 7 次元の基本の補助特徴量ベクトルを含む。さらに、補助特徴量ベクトル $a_i^{(u)}$ では、前向き (forword) L S T M L M の単語予測スコアを 1 8 次元の補助特徴量として用いてもよい。L S T M L M は、長短期記憶メモリ (long short-term memory : L S T M) ユニットを用いた再帰的ニューラルネットワーク (Recurrent Neural Network : R N N) 言語モデルであり、後述するように、本実施の形態 1 ~ 3 における各補助モデルを構成するモデルである。そして、補助特徴量ベクトル a

50

$i(u)$ では、後向き(backward) L S T M L Mの単語予測スコアを19次元の補助特徴量として用いてもよい。後向き L S T M L Mは、未来の単語列から現在の単語の生起確率を予測するものであるであり、前向き L S T M L Mと相補的な単語予測能力を持つことから各補助モデルが出力する判定情報の精度向上が期待できる。

【0039】

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、 $x_i(u) = \text{concat}(\text{embed}(w_i(u)), a_i(u))$ 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2とする。))を参照)を表す。なお、 $\text{embed}(\cdot)$ を行うNNも第1補助モデル111~第M補助モデル11M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。

10

【0040】

そして、 $P(0 | X(u), X(v))$ を、二つの仮説 $W(u), W(v)$ の精度の高低の並びが正しいことを示す事後確率と定義する。 $P(0 | X(u), X(v))$ は、第1補助モデル111~第M補助モデル11M及びメインモデル110のそれぞれにおいて、生成される。

20

【0041】

[補助モデル及びメインモデルの構成]

第1補助モデル111~第M補助モデル11M、メインモデル及び構成について説明する。図2は、第1補助モデル111~第M補助モデル11M及びメインモデル110の構成を説明する図である。図2では、処理の流れを説明するため、判定部15も記載される。

【0042】

図2に示すように、各補助モデルは、それぞれ、二つの第1変換部、第1結合部及び判定情報生成部を有する。具体的に、第1補助モデル111を例に説明する。第1補助モデル111は、二つの第1変換部111-1u, 111-1v、第1結合部111-2及び判定情報生成部111-3を有する。

30

【0043】

第1変換部111-1uは、比較対象の二つの仮説 $W(u), W(v)$ のうち、仮説 $W(u)$ の特徴量 $X(u)$ の入力を受け付け、隠れ状態ベクトルに変換する。第1変換部111-1vは、比較対象の二つの仮説 $W(u), W(v)$ のうち、仮説 $W(v)$ の特徴量 $X(v)$ の入力を受け付け、隠れ状態ベクトルに変換する。

【0044】

第1結合部111-2は、第1変換部111-1u, 第1変換部111-1vが変換した二つの隠れ状態ベクトルを結合する。判定情報生成部111-3は、二つの仮説 $W(u), W(v)$ の精度の高低の並びが正しいことを示す事後確率 $P(0 | X(u), X(v))$ を、判定情報として生成する。他の補助モデルも、第1補助モデル111と同じ構成であり、与えられた二つの仮説 $W(u), W(v)$ に対し、それぞれ、隠れ状態ベクトルの変換、隠れ状態ベクトルの結合、及び、判定情報の生成を含むタスクをそれぞれ実行できるようにしている。なお、各補助モデルは、学習時におけるランダム初期化時における初期値がそれぞれ異なる。

40

【0045】

メインモデル110は、メイン結合部110-1と、判定情報生成部110-2とを有する。メイン結合部110-1は、第1補助モデル111~第M補助モデル11Mでそれぞれ変換された二つの仮説の隠れ状態ベクトルを結合する。判定情報生成部111-3は、二つの仮説 $W(u), W(v)$ の精度の高低の並びが正しいことを示す事後確率 $P(0 | X(u), X(v))$ を、判定情報として生成する。

50

【 0 0 4 6 】

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 $X^{(u)}$, $X^{(v)}$ が各補助モデルに与えられたとき、メインモデル 110 は、記号 $y = \{ 0 \}$ の事後確率 $P(0 | X^{(u)}, X^{(v)})$ を出力する。

【 0 0 4 7 】

判定部 15 は、メインモデル 110 が出力した事後確率 $P(0 | X^{(u)}, X^{(v)})$ を受け取り、判定を行う。 $P(0 | X^{(u)}, X^{(v)})$ は、u 位の仮説と v 位の仮説との順位の間関係が正しさを確率的に表現する事後確率である。判定部 15 は、N ベストランキングモデルから出力された事後確率 $P(0 | X^{(u)}, X^{(v)})$ を取得し、取得した事後確率を所定の閾値と比較して、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。

10

【 0 0 4 8 】

具体的には、判定部 15 は、事後確率 $P(0 | X^{(u)}, X^{(v)})$ が 0.5 以上である場合には、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し、 $y = 0$ を出力する。また、判定部 15 は、事後確率 $P(0 | X^{(u)}, X^{(v)})$ が 0.5 未満である場合には、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し、 $y = 1$ を出力する。

【 0 0 4 9 】

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。

【 0 0 5 0 】

$$P(0 | X^{(u)}, X^{(v)}) = \begin{cases} 0.5 & \text{if WER (Word error rate) of } W^{(u)} \leq \text{WER of } W^{(v)} \cdots (1-1) \\ < 0.5 & \text{otherwise} \cdots (1-2) \end{cases}$$

20

【 0 0 5 1 】

ここで、与えられた仮説 (単語列) の音声認識精度を返す関数 $y = P(y | X^{(u)}, X^{(v)}) = 1$ であるため、(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。

30

【 0 0 5 2 】

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、 $W^{(v)}$ との一对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、次の一对一の仮説比較では、 $W^{(v)}$ として使用する。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。

【 0 0 5 3 】

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、 $W^{(u)}$ との一对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、次の一对一の仮説比較では、 $W^{(v)}$ として引き続き使用する。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。

40

【 0 0 5 4 】

[補助モデルの構築例]

第 1 補助モデル 111 ~ 第 M 補助モデル 11M の構築例について説明する。第 1 補助モデル 111 ~ 第 M 補助モデル 11M は、同じ構成であるため、図 3 を参照し、第 1 補助モデル 111 の構築例を説明する。図 3 は、第 1 補助モデル 111 の構築例を示す図である

50

。なお、図3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。

【0055】

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、第1補助モデル111は、二つの仮説の特徴量を、RNNを用いて隠れ状態ベクトルに変換する。具体的には、第1補助モデル111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2参照) のエンコーダ RNN_{111-1a} を第1変換部 $111-1u$, $111-1v$ として有する。

【0056】

第1補助モデル111は、エンコーダ RNN_{111-1a} 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。そして、第1補助モデル111~第M補助モデル11M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。

【0057】

エンコーダ RNN_{111-1a} の処理について説明する。エンコーダ RNN_{111-1a} 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2参照) 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。

【0058】

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

【0059】

ここで、 $lstm(\cdot)$ 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_{111-1a} 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。

【0060】

エンコーダ RNN_{111-1a} 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3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。

【0061】

第1補助モデル111は、以上で得た二つの隠れ状態ベクトル $h_{L(W^{(u)})}^{(u)}, h_{L(W^{(v)})}^{(v)}$ を、第1結合部111-2で連結した隠れ状態ベクトル $h_{\{(u,v)\}}$ をエンコーダ RNN_{111-a} の出力として以下の(3)式のように得る。

【0062】

$$h_{\{(u,v)\}} = concat(h_{L(W^{(u)})}^{(u)}, h_{L(W^{(v)})}^{(v)}) \dots (3)$$

【0063】

そして、第1補助モデル111は、エンコーダ RNN_{111-1a} 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、第1補助モデル111は、1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) $111-3a$ (詳細は、例えば、参考文献2を参照) を、判定情報生成部113として用いる。

10

20

30

40

50

エンコーダ RNN 1 1 1 - 1 a の出力として得た隠れ状態ベクトル $h\{(u, v)\}$ が、1 層の 1 クラス分類 FFNN 1 1 1 - 3 a に入力され、最終的に、1 クラスの $y = \{0\}$ 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。

【0064】

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

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

【0065】

ここで、 $\text{linear}(\cdot)$ は、線形変換処理（詳細は、例えば、参考文献 2 を参照）を表す。 $\text{sigmoid}(\cdot)$ は、シグモイド処理を表す。

10

【0066】

また、メインモデル 1 1 0 では、メイン結合部 1 1 0 - 1 は、第 1 結合部 1 1 1 - 2 と同様のベクトル連結処理を行う。また、メインモデル 1 1 0 では、判定情報生成部 1 1 0 - 3 は、判定情報生成部 1 1 1 - 3 の 1 層の 1 クラス分類 FFNN 1 1 1 - 3 a と同様の構成の 1 クラス分類 FFNN によって構成される。

【0067】

[補助モデル及びメインモデルの他の構築例 1]

なお、第 1 補助モデル 1 1 1 及びメインモデル 1 1 0 は、1 クラス分類 FFNN におけるシグモイド処理に代えて、ソフトマックス処理を行ってもよい。この場合、エンコーダ RNN 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 (6), (7) 式のように得ることができる。なお、 $y = 0$ 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。

20

【0068】

$$z\{(u, v)\} = \text{linear}(h\{(u, v)\}) \dots (6)$$

$$P(y | X^{(u)}, X^{(v)}) = \text{softmax}(z\{(u, v)\})_y \dots (7)$$

30

【0069】

ここで、 $\text{softmax}(\cdot)$ は、ソフトマックス処理を表す。また、 $\text{softmax}(\cdot)_y$ 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。

【0070】

この場合、判定部 1 5 は、メインモデル 1 1 0 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きさを比較して、 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 1 5 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 1 5 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。

40

【0071】

[補助モデルの他の構築例 2]

なお、図 3 に示すエンコーダ RNN 1 1 1 - 1 a 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。

【0072】

[補助モデルの他の構築例 3]

50

また、LSTMユニットの代わりに、単純な(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。

【0073】

[補助モデル及びメインモデルの他の構築例4]

さらに、補助モデル及びメインモデル110は、図3の構築例では、1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、補助モデル及びメインモデル110の他の構築例1~4の用語の詳細については、例えば、参考文献2を参照いただきたい。

10

【0074】

[補助モデルの他の構築例5]

また、補助モデルは、従来のNベストリスクアリングモデル(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。

【0075】

[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10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4は、実施の形態1に係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。

20

【0076】

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、スコアの昇順に、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($u < v < N$)。まず、仮説選択部13は、 $u = N - 1$ 、 $v = N$ 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($X^{(u)}$ 、 $X^{(v)}$)を各補助モデル(第1補助モデル111~第M補助モデル11M)に入力する(ステップS5)。

30

【0077】

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、事後確率 $P(0 | X^{(u)}, X^{(v)})$ を取得する。

【0078】

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、 $P(0 | X^{(u)}, X^{(v)}) > 0.5$ であるか否かを判定する(ステップS7)。 $P(0 | X^{(u)}, X^{(v)}) > 0.5$ である場合(ステップS7: Yes)、判定部15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し、実行制御部16は、kについて $k = u$ と設定する(ステップS8)。kは、比較処理後の仮説のうち、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、 $P(0 | X^{(u)}, X^{(v)}) < 0.5$ でない場合(ステップS7: No)、判定部15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し、実行制御部16は、 $k = v$ と設定する(ステップS9)。

40

【0079】

続いて、実行制御部16は、 $u = 1$ であるか否かを判定する(ステップS10)。 $u = 1$ でない場合(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、実行制御部16は、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、 $u = u - 1$ 、 $v = k$ に設定し(ステップS

50

11)、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。

【0080】

また、 $u = 1$ である場合(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、実行制御部16は、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。

10

【0081】

このように、実施の形態1に係るリランキング装置10は、一対一の二つの仮説の比較を行う機能を持つ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。さらに、リランキング装置10では、モデルとして、ニューラルネットワーク(NN)で表されるメインモデル110と、NNで表される複数の補助モデルとを用いる。

【0082】

すなわち、リランキング装置10では、複数の補助モデルを設け、入力された二つの仮説に対して、各補助モデルにタスクを実行させている。各補助モデルの構造は同じであっても、学習時においてパラメータのランダム初期化を行うので、同じ入力仮説に対しても異なる隠れ状態ベクトルを出力する。これにより、ある二つの入力仮説に対して、ある補助モデルが出力する隠れ状態ベクトルが適切なものでなかったとしても、別の補助モデルが適切な隠れ状態ベクトルを出力できる可能性が高まる。つまり、正確な仮説の判定結果を生成するのに適した隠れ状態ベクトルが、いずれかの補助ネットワークから出力される可能性が高くなる。この結果、リランキング装置10のメインモデル110には、適切な二つの仮説に対応する隠れ状態ベクトルが安定して入力されるため、メインモデル110の出力値の精度も安定する。このように、実施の形態1に係るリランキングモデルは、安定した精度で、Nベスト仮説の中からオラクル仮説を見つけ出すことができる。

20

【0083】

また、リランキング装置10は、Nベスト仮説のスコアの昇順に二つの仮説を選択する。言い換えると、リランキング装置10は、Nベスト仮説のうち、スコアが最も低い仮説から順に仮説ペアを選択する。図5は、図1に示すリランキング装置10が、Nベスト仮説に対して実行するリランキング処理を説明する図である。

30

【0084】

一般には、スコアが高い仮説の方が、尤もらしい仮説である可能性が高い。スコアの高い順に仮説を選択していくと、最もスコアの高い仮説は、N-1回の判定処理に勝ち抜かなければ、最終的な出力仮説として選択されず、尤もらしい仮説として選ばれにくくなってしまう。

【0085】

そこで、図5に示すように、リランキング装置10は、最終的に出力仮説として選ばれる可能性の高い仮説について、少ない判定回数で済むように、Nベスト仮説のうち、スコアが最も低い仮説から順に仮説ペアを選択する。言い換えると、リランキング装置10は、図5に示すように、スコアの最も高い仮説については、シード権を与え、Nベスト仮説全体に対する比較処理の後の方の処理で比較処理が行われるようにし、尤もらしい仮説として選ばれやすくしている。このように、リランキング装置10は、最終的に出力仮説として選ばれる可能性の高い仮説が、尤もらしい仮説として選ばれやすいため、安定した精度で、Nベスト仮説の中からオラクル仮説を見つけ出すことができる。

40

【0086】

[実施の形態2]

[学習装置]

50

次に、実施の形態 2 として、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 6 は、実施の形態 2 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 2 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 6 に示すように、モデル記憶部 21、学習装置 20 は、仮説入力部 22 及び学習部 23 を有する。

【0087】

モデル記憶部 21 は、学習対象の第 1 補助モデル 111 ~ 第 M 補助モデル 11M 及びメインモデル 110 を記憶する。第 1 補助モデル 111 ~ 第 M 補助モデル 11M 及びメインモデル 110 は、選択された二つの仮説が与えられたとき、二つの仮説を隠れ状態ベクトルに変換し、二つの仮説の隠れ状態ベクトルを基に二つの仮説の精度の高低を判定できるような、NN で表される。第 1 補助モデル 111 ~ 第 M 補助モデル 11M 及びは、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、第 1 補助モデル 111 ~ 第 M 補助モデル 11M 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事後確率を判定情報として生成する。

10

【0088】

メインモデル 110 は、第 1 補助モデル 111 ~ 第 M 補助モデル 11M においてそれぞれ変換された二つの仮説の隠れ状態ベクトルを基に、二つの仮説の精度の高低を判定できるような、NN で表される。メインモデル 110 は、NN を用いて、第 1 補助モデル 111 ~ 第 M 補助モデル 11M においてそれぞれ変換された学習用の二つの仮説の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事後確率を生成する。

20

【0089】

仮説入力部 22 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。

【0090】

学習部 23 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、第 1 補助モデル 111 ~ 第 M 補助モデル 11M 及びメインモデル 110 に対し、各 NN が二つの仮説の精度の高低を判定するタスクを個別に行うとみなしたマルチタスク学習を行わせる。学習部 23 は、各 NN によって実行された各タスクについて所定の損失をそれぞれ計算し、各損失の重み付け和を全体の損失関数とする。そして、学習部 23 は、この全体の損失関数に基づいて、各 NN のパラメータの値を更新する。

30

【0091】

なお、学習部 23 は、各損失に対し、等重みで重み付けをしてもよい。また、メインモデル 110 が出力する判定情報が判定部 15 における判定に使用されるため、学習部 23 は、メインモデル 110 に、他の補助モデルよりも多めの重みを付けてもよい。

【0092】

学習部 23 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、第 1 補助モデル 111 ~ 第 M 補助モデル 11M 及びメインモデル 110 に与える。これによって、学習部 23 は、第 1 補助モデル 111 ~ 第 M 補助モデル 11M 及びメインモデル 110 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、第 1 補助モデル 111 ~ 第 M 補助モデル 11M 及びメインモデル 110 の学習（パラメータの最適化）を行う。

40

【0093】

具体的には、学習部 23 は、特徴量ベクトル列と、対応する教師ラベルとを第 1 補助モデル 111 ~ 第 M 補助モデル 11M に入力し、第 1 補助モデル 111 ~ 第 M 補助モデル 11M 及びメインモデル 110 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、第 1 補助モデル 111 ~ 第 M 補助モデル 11M 及びメ

50

インモデル 1 1 0 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。

【 0 0 9 4 】

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、第 1 補助モデル 1 1 1 ~ 第 M 補助モデル 1 1 M 及びメインモデル 1 1 0 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、第 1 補助モデル 1 1 1 ~ 第 M 補助モデル 1 1 M 及びメインモデル 1 1 0 に学習させる。

【 0 0 9 5 】

入替部 2 3 2 は、学習用の N ベスト仮説のうちの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、第 1 補助モデル 1 1 1 ~ 第 M 補助モデル 1 1 M 及びメインモデル 1 1 0 の学習を行う。図 7 は、図 6 に示す入替部 2 3 2 の処理を説明する図である。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 (図 7 の (1) 参照)、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 (図 7 の (2) 参照)。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 9 6 】

[学習処理の処理手順]

次に、図 6 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 8 は、実施の形態 2 に係る学習処理の処理手順を示すフローチャートである。図 8 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、 $W^{(u)}$ の精度は、 $W^{(v)}$ の精度よりも高いときの学習処理の処理手順を示す。

【 0 0 9 7 】

図 8 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 $X^{(u)}$, $X^{(v)}$ を第 1 補助モデル 1 1 1 ~ 第 M 補助モデル 1 1 M に入力する (ステップ S 2 2)。そして、学習部 2 3 は、第 1 補助モデル 1 1 1 ~ 第 M 補助モデル 1 1 M 及びメインモデル 1 1 0 にマルチタスク学習を行わせて、第 1 補助モデル 1 1 1 ~ 第 M 補助モデル 1 1 M 及びメインモデル 1 1 0 のモデルパラメータを更新させる (ステップ S 2 3)。

【 0 0 9 8 】

すなわち、この二つの仮説の $W^{(u)}$, $W^{(v)}$ の特徴量ベクトル $X^{(u)}$, $X^{(v)}$ を第 1 補助モデル 1 1 1 ~ 第 M 補助モデル 1 1 M に入力した場合、第 1 補助モデル 1 1 1 ~ 第 M 補助モデル 1 1 M 及びメインモデル 1 1 0 は、理想的には、 $P(0 | X^{(u)}, X^{(v)}) = 1$ 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 2 3 1 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 2 3 は、第 1 補助モデル 1 1 1 ~ 第 M 補助モデル 1 1 M 及びメインモデル 1 1 0 のモデルパラメータ (エンコーダ RNN (LSTM ユニット)、1 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。

【 0 0 9 9 】

そして、入替部 2 3 2 は、仮説 $W^{(u)}$, $W^{(v)}$ の順位を入れ替える (ステップ S 2 4)。すなわち、入替部 2 3 2 は、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、 $W^{(u)}$ の精度は、 $W^{(v)}$ の精度よりも低い。よって、この二つの仮説 $W^{(u)}$, $W^{(v)}$ の特徴量ベクトル $X^{(u)}$, $X^{(v)}$ を第 1 補助モデル 1 1 1 ~ 第 M 補助モデル 1 1 M 及びメインモデル 1 1 0 に入力した場合、第 1 補助モデル 1 1 1 ~ 第 M 補助モデル 1 1 M 及びメインモデル 1 1 0 は、理想的には、 $P(0 | X^{(u)}, X^{(v)}) = 0$ の事後確率を出力すべきである。

【 0 1 0 0 】

このため、教師ラベル付与部 2 3 1 は、教師ラベルとして、 $y = 1$ を付与し (ステップ S 2 5)、 $W^{(u)}$, $W^{(v)}$ の特徴量 $X^{(u)}$, $X^{(v)}$ を第 1 補助モデル 1 1 1 ~ 第 M 補

10

20

30

40

50

助モデル 1 1 M 及びメインモデル 1 1 0 に入力する (ステップ S 2 6)。学習部 2 3 は、以上の入力を基に、第 1 補助モデル 1 1 1 ~ 第 M 補助モデル 1 1 M 及びメインモデル 1 1 0 にマルチタスク学習を行わせて、第 1 補助モデル 1 1 1 ~ 第 M 補助モデル 1 1 M 及びメインモデル 1 1 0 を更新させて (ステップ S 2 7)、二つの仮説 $W^{(u)}$, $W^{(v)}$ に対する学習処理を終了する。

【 0 1 0 1 】

学習装置 2 0 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 2 3 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 2 参照) と同様に行うことができる。

【 0 1 0 2 】

[実施の形態 2 の効果]

このように、実施の形態 2 に係る学習装置 2 0 は、第 1 補助モデル 1 1 1 ~ 第 M 補助モデル 1 1 M 及びメインモデル 1 1 0 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マルチタスク学習を行わせている。したがって、学習装置 2 0 は、N ベストランキングを行う上で最適な第 1 補助モデル 1 1 1 ~ 第 M 補助モデル 1 1 M 及びメインモデル 1 1 0 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 1 0 は、学習装置 2 0 において学習された第 1 補助モデル 1 1 1 ~ 第 M 補助モデル 1 1 M 及びメインモデル 1 1 0 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、安定した制度でオラクル仮説を抽出することができる。

【 0 1 0 3 】

[学習処理の効率化例 1]

図 8 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 1 0 0 ~ 1 0 0 0 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。

【 0 1 0 4 】

上述したように、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 ベストランキングモデルに入力する二つの仮説のうち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。

【 0 1 0 5 】

[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。

【 0 1 0 6 】

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。

仮説 4 は、音声認識スコアが最も低い仮説である。

【 0 1 0 7 】

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 (

10

20

30

40

50

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いけない) 仮説である。その他の仮説をこの四つの中に絞り込む場合は、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。

【0108】

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、二つの仮説のうちの他方の仮説として、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。

10

【0109】

[評価]

実際に、実施の形態1における N ベストリランキングと、非特許文献1記載の N ベストリランキングとの比較評価を行った。表1は、CSJ音声コーパスを用いて、非特許文献1記載の N ベストリランキングとの比較評価する100(= N) ベストリランキング評価を行った結果を示す表である。表の数値は、WER(Word error rate)率[%]であり、Dev(Development)、Eval(Evaluation)を示す。

【0110】

20

【表1】

(表1)

No.	Model	Dev	Eval
1	Single-encoder DDM	16.4	13.7
2	Eight-encoder DDM	16.1	13.4
3	2 with fwd & bwd LSTM scores	15.2	12.6
4	Oracle	11.6	9.7

30

【0111】

表1の通番「1」は、非特許文献1記載の N ベストリランキング結果である。表1の通番「2」は、実施の形態1に係るリランキング装置10であって8個の補助モデルを有する場合の結果である。通番「3」は、通番「2」の条件に加え、前向き及び後ろ向きLSTMの単語予測スコアを18時限目及び19次元目の補助特徴として用いている。通番「4」は、参考のために示されたオラクルである。

【0112】

40

表1に示すように、通番「1」の非特許文献1記載のリランキング方法でも、十分にWERを削減できるが、通番「2」の8個の補助モデルを有するリランキング装置10では、さらにWER削減が実現できる。また、通番「3」の評価結果に示すように、両方向のLSTMの単語予測スコアと8個の補助モデルを用いることで、相補的なWER削減効果が得られることが確認できた。また、通番「2」以外にも、補助モデルの個数を、2または4とした構成でも評価を行っており、補助モデルの個数を増やすほどWERを削減できる傾向が確認できた。

【0113】

この評価結果から、本実施の形態1に係るリランキング装置10は、非特許文献1記載のリランキング方法と比して、安定したWER削減が実現できる。

50

【 0 1 1 4 】

[実施の形態 3]

なお、実施の形態 1 のリランキング装置 1 0 は、メインモデル 1 1 0 の出力を用いて判定を行ったが、メインモデル 1 1 0 の出力の他に各補助モデルの出力を用いて判定を行ってもよい。図 9 は、実施の形態 3 に係るリランキング装置の要部構成を示す図である。

【 0 1 1 5 】

図 3 に示すように、リランキング装置は、判定部 1 5 の前段に、重み付け部 1 8 を有する。重み付け部 1 8 は、第 1 補助モデル 1 1 1 ~ 第 M 補助モデル 1 1 M 及びメインモデル 1 1 0 から出力された全ての判定情報を取得し、各判定情報に対して重み付け和を計算する。

10

【 0 1 1 6 】

なお、各判定情報に対応する重みは予め設定されている。重み付け部 1 8 は、各判定情報に対し、全補助モデル及びメインモデル 1 1 0 に対して等重みで重み付けをしてもよい。また、重み付け部 1 8 は、メインモデル 1 1 0 に、他の補助モデルよりも多めの重みを付けてもよい。また、重み付け部 1 8 は、予め各判定情報に対する重みを学習した 1 層の線形 NN を有し、各判定情報が入力されると各判定対象に対する重みを求めてもよい。

【 0 1 1 7 】

判定部 1 5 は、重み付け部 1 8 が計算した重み付け和の値に基づいて二つの仮説の精度の高低を判定する。例えば、判定部 1 5 は、判定情報のそれぞれが、仮説 $W^{(u)}$ が選択される確率を示すものとして、判定情報の重みづけ和を 0 ~ 1 の範囲に収まるように正規化した値が 0 . 5 以上であれば仮説 $W^{(u)}$ を選択し、そうでなければ仮説 $W^{(v)}$ を選択する。

20

【 0 1 1 8 】

[判定処理の処理手順]

図 1 0 は、実施の形態 3 に係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。

【 0 1 1 9 】

図 1 0 に示すステップ S 3 1 ~ ステップ S 3 6 は、図 4 に示すステップ S 1 ~ ステップ S 6 と同じ処理である。重み付け部 1 8 は、第 1 補助モデル 1 1 1 ~ 第 M 補助モデル 1 1 M 及びメインモデル 1 1 0 から出力された全ての判定情報を取得し、各判定情報に対して重み付け和を計算する重み付け処理を行う (ステップ S 3 7)。そして、重み付け部 1 8 が計算した重み付け和の値に基づいて二つの仮説の精度の高低を判定する。ステップ S 3 8 ~ ステップ S 4 3 は、図 4 に示すステップ S 7 ~ ステップ S 1 2 と同じ処理である。

30

【 0 1 2 0 】

[実施の形態 3 の効果]

この実施の形態 3 に示すように、メインモデル 1 1 0 による判定情報に加え、全補助モデルによる判定情報を用いて、判定を行うことも可能である。この際、実施の形態 3 では、各補助モデル或いはメインモデル 1 1 0 に応じて、各判定情報に対する重み付けを行い、重み付け和の値に基づいて二つの仮説の精度の高低を判定するため、オラクル仮説を抽出精度を保持することができる。

40

【 0 1 2 1 】

なお、本実施の形態では、全ての仮説に対して、比較処理を行うため、N ベスト仮説のソートも可能である。

【 0 1 2 2 】

また、本実施の形態 1 ~ 3 では、音声認識の N ベスト仮説をリランキングするためのモデルとして、図 1 に例示する第 1 補助モデル 1 1 1 ~ 第 M 補助モデル 1 1 M 及びメインモデル 1 1 0 について説明した。ただし、本実施の形態 1 ~ 3 の第 1 補助モデル 1 1 1 ~ 第 M 補助モデル 1 1 M 及びメインモデル 1 1 0 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

50

、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。

【0123】

このため、本実施の形態1～3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1～3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1～3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1～3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

【0124】

また、この場合には、本実施の形態1～3では、精度が既知である学習用の複数の系列のうち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で表される第1補助モデル111～第M補助モデル11M及びメインモデル110にマルチタスク学習を行わせる。そして、本実施の形態1～3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第1補助モデル111～第M補助モデル11M及びメインモデル110に学習させる。そして、本実施の形態1～3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第1補助モデル111～第M補助モデル11M及びメインモデル110に学習させる。本実施の形態1～3では、この第1補助モデル111～第M補助モデル11M及びメインモデル110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

【0125】

[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

【0126】

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

【0127】

[プログラム]

図11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

50

リ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。

【0128】

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

【0129】

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

【0130】

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。

【0131】

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

【0132】

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】

【0133】

- 2 音声認識装置
- 10 リランキング装置
- 11, 21 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

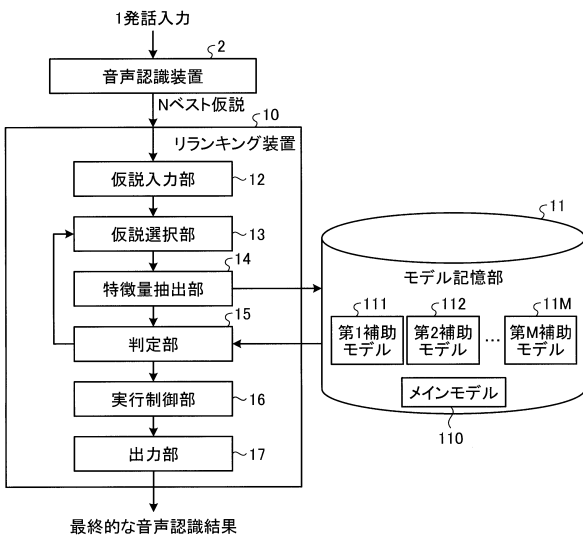
50

- 1 6 実行制御部
- 1 7 出力部
- 1 8 重み付け部
- 2 0 学習装置
- 2 2 仮説入力部
- 2 3 学習部
- 1 1 0 メインモデル
- 1 1 1 ~ 1 1 M 第 1 補助モデル ~ 第 M 補助モデル
- 2 3 1 教師ラベル付与部
- 2 3 2 入替部

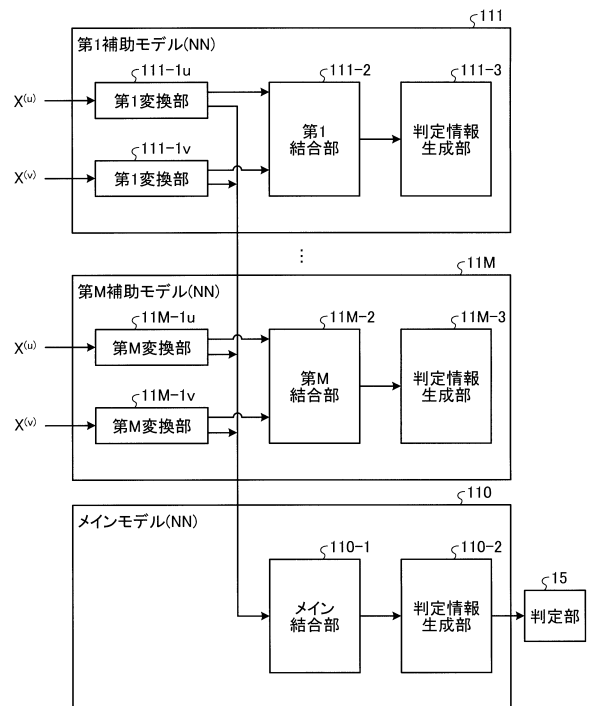
10

【図面】

【図 1】



【図 2】



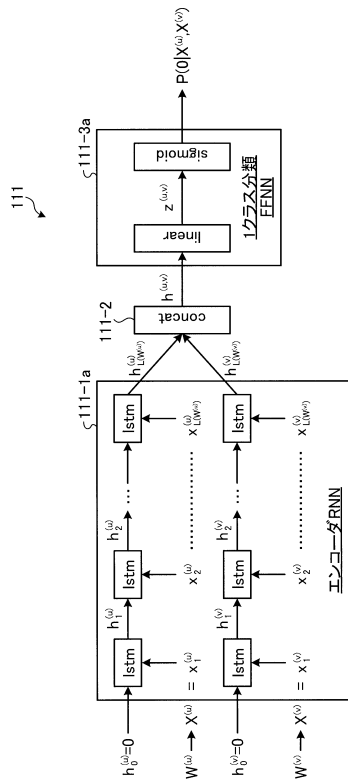
20

30

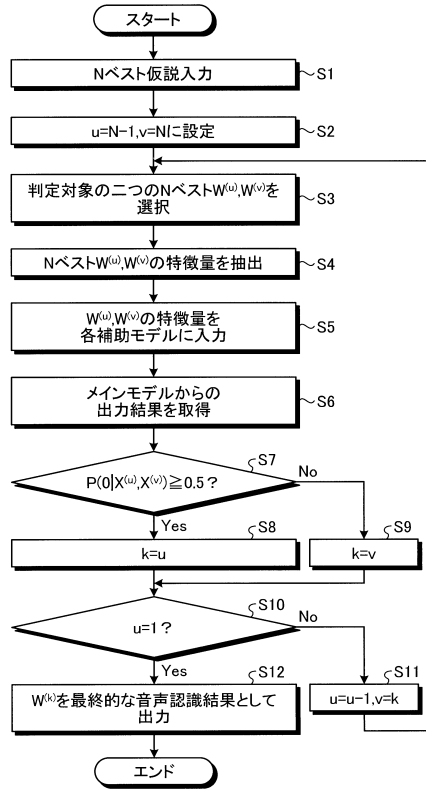
40

50

【図3】



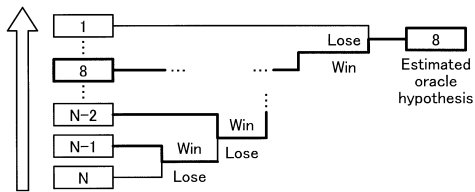
【図4】



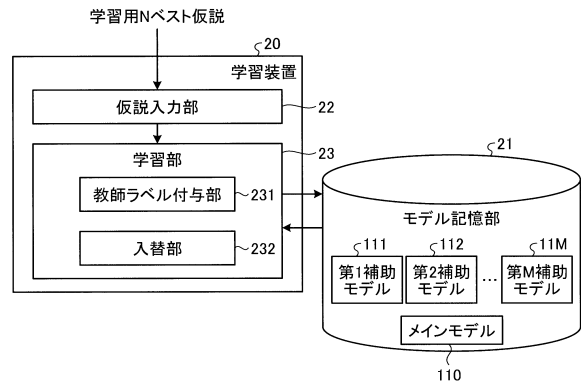
10

20

【図5】



【図6】

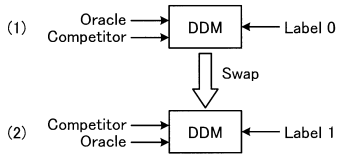


30

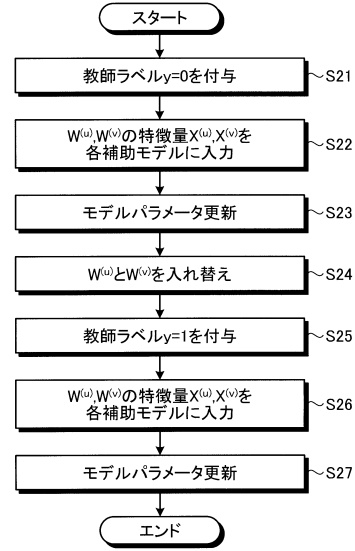
40

50

【 図 7 】



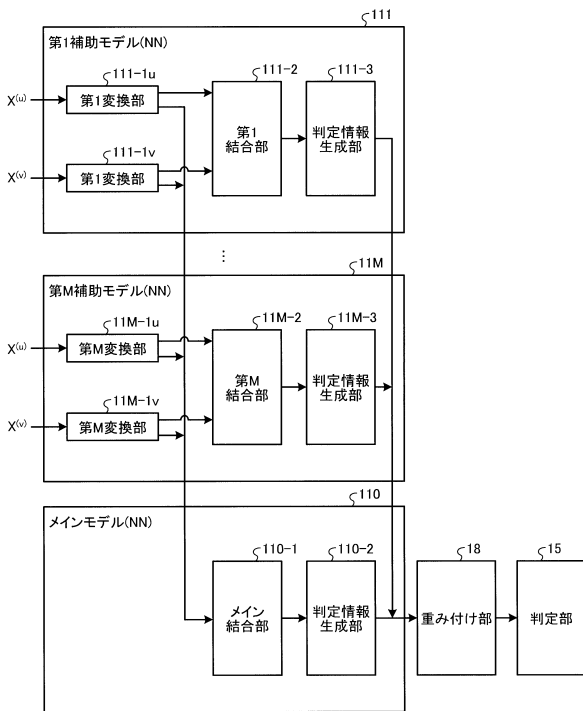
【 図 8 】



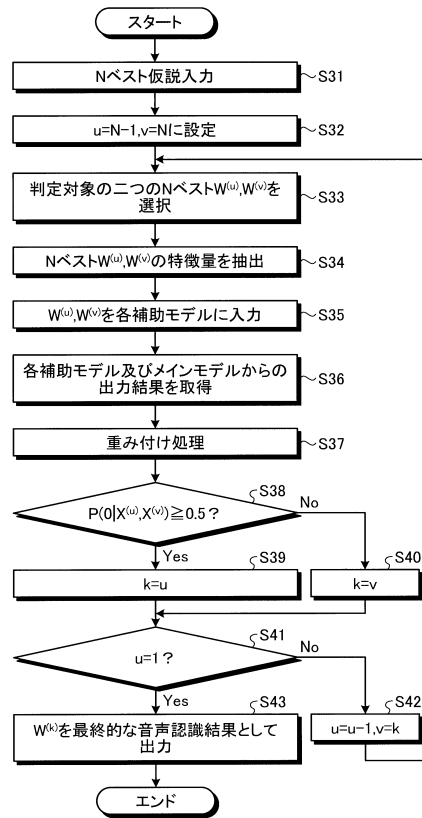
10

20

【 図 9 】



【 図 10 】

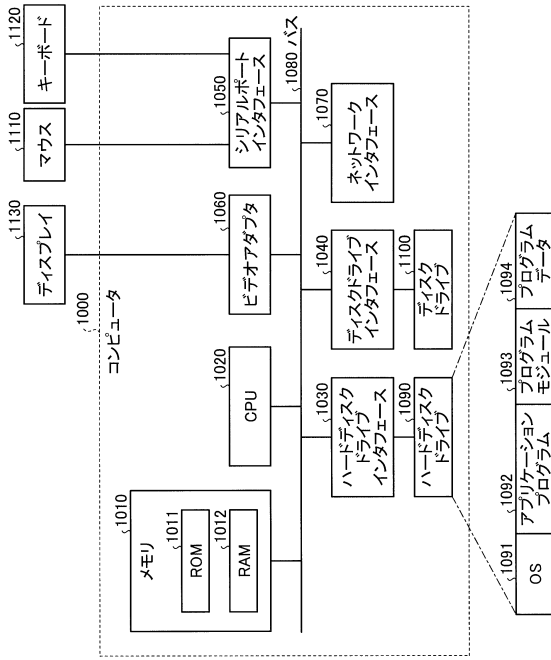


30

40

50

【 図 1 1 】



10

20

30

40

50

フロントページの続き

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 山下 剛史

- (56)参考文献 特開2011-243147(JP,A)
特開2018-60047(JP,A)
米国特許出願公開第2017/0221474(US,A1)
米国特許第10032463(US,B1)
小川厚徳他, 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, 日本音響学会2018年春季研究発表会講演論文集[CD-ROM], 2018年03月, pp.23-24
田中智大他, 複数仮説を考慮したニューラル誤り訂正言語モデルの検討, 電子情報通信学会技術研究報告, 2018年08月, Vol.118, No.198, pp.31-36
- (58)調査した分野 (Int.Cl., DB名)
G10L 15/00 - 15/34