

(19) 日本国特許庁 (JP)

## (12) 公 開 特 許 公 報 (A)

(11) 特許出願公開番号

特開2015-179242

(P2015-179242A)

(43) 公開日 平成27年10月8日 (2015. 10. 8)

(51) Int.Cl.		F I		テーマコード (参考)
<b>G 1 0 L 15/07 (2013.01)</b>		G 1 0 L 15/07		
<b>G 1 0 L 15/16 (2006.01)</b>		G 1 0 L 15/16		
<b>G 1 0 L 15/06 (2013.01)</b>		G 1 0 L 15/06	4 0 0 U	

審査請求 未請求 請求項の数 9 O L (全 20 頁)

(21) 出願番号	特願2014-187022 (P2014-187022)	(71) 出願人	899000068
(22) 出願日	平成26年9月12日 (2014. 9. 12)		学校法人早稲田大学
(31) 優先権主張番号	特願2014-39639 (P2014-39639)		東京都新宿区戸塚町 1 丁目 1 〇 4 番地
(32) 優先日	平成26年2月28日 (2014. 2. 28)	(74) 代理人	100080089
(33) 優先権主張国	日本国 (JP)		弁理士 牛木 護
		(74) 代理人	100121153
			弁理士 守屋 嘉高
		(74) 代理人	100161665
			弁理士 高橋 知之
		(74) 代理人	100133639
			弁理士 矢野 卓哉
		(72) 発明者	新田 恒雄
			東京都新宿区戸塚町 1 丁目 1 〇 4 番地 学 校法人早稲田大学内

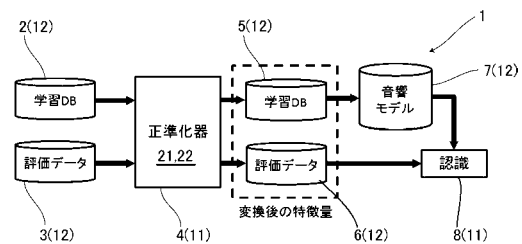
(54) 【発明の名称】 音声認識装置、音声認識方法及び音声認識プログラム

## (57) 【要約】

【課題】一度の計算で話者変動を効果的に抑圧して、全ての入力音声に対して高い認識性能を実現できる音声認識装置を提供する。

【解決手段】音声認識装置 1 は、任意話者の対数パワースペクトル  $X_1, X_2, \dots, X_{24}$  を、標準話者に正準化した対数パワースペクトル  $y_1, y_2, \dots, y_{24}$  に変換する話者正準化手段として正準化器 4 を備える。正準化器 4 は、標準話者の音響特徴を教師データとして、任意話者の対数パワースペクトル  $X_1, X_2, \dots, X_{24}$  を、標準話者に正準化した対数パワースペクトル  $y_1, y_2, \dots, y_{24}$  に周波数軸上で非線形に変換する M L P 3 1 を含む写像関数学習手段 2 2 を備えている。

【選択図】図 1



**【特許請求の範囲】****【請求項 1】**

任意話者の音声スペクトルを、標準話者の音声スペクトルに変換する話者正準化手段を備えた音声認識装置において、

前記標準話者の音響特徴を教師データとして、前記任意話者の音声スペクトルを、前記標準話者の音声スペクトルに周波数軸上で、スペクトル形状の違いに応じて非線形に変換するニューラルネットワークを含む写像関数学習手段を、前記話者正準化手段に備えたことを特徴とする音声認識装置。

**【請求項 2】**

前記写像関数学習手段は、前記ニューラルネットワークからの出力を、前記標準話者の音声スペクトルとしてそのまま音声認識に用いる周波数領域を制限する構成としたことを特徴とする請求項 1 記載の音声認識装置。

10

**【請求項 3】**

前記写像関数学習手段は、制限された前記周波数領域以外の周波数領域で、前記ニューラルネットワークからの出力と、前記任意話者の音声スペクトルのそれぞれを重み付けして合成し、前記標準話者の音声スペクトルに変換出力する構成としたことを特徴とする請求項 2 記載の音声認識装置。

**【請求項 4】**

前記多数話者の音声スペクトルをクラスター分析することにより、当該多数話者の中から前記標準話者を特定する標準話者確定手段をさらに備えたことを特徴とする請求項 1 ~ 3 の何れか一つに記載の音声認識装置。

20

**【請求項 5】**

任意話者の音声スペクトルを、標準話者の音声スペクトルに変換して話者正準化を行なう音声認識方法において、

前記話者正準化では、ニューラルネットワークを用い、前記標準話者の音響特徴を教師データとして、前記任意話者の音声スペクトルを、前記標準話者の音声スペクトルに周波数軸上で、スペクトル形状の違いに応じて非線形に変換することを特徴とする音声認識方法。

**【請求項 6】**

前記ニューラルネットワークからの出力は、前記標準話者の音声スペクトルとしてそのまま音声認識に用いる周波数領域が制限されていることを特徴とする請求項 5 記載の音声認識方法。

30

**【請求項 7】**

制限された前記周波数領域以外の周波数領域で、前記ニューラルネットワークからの出力と、前記任意話者の音声スペクトルのそれぞれを重み付けして合成し、前記標準話者の音声スペクトルに変換出力することを特徴とする請求項 6 記載の音声認識方法。

**【請求項 8】**

前記多数話者の音声スペクトルをクラスター分析することにより、当該多数話者の中から前記標準話者を特定することを特徴とする請求項 5 ~ 7 の何れか一つに記載の音声認識方法。

40

**【請求項 9】**

請求項 5 ~ 8 の何れか一つに記載の音声認識方法を、コンピュータに実行させるための音声認識プログラム。

**【発明の詳細な説明】****【技術分野】****【0001】**

本発明は、話者正準化の手法を用いて、話者の音声信号を高い精度で識別し得る音声認識装置、音声認識方法及び音声認識プログラムに関する。

**【背景技術】****【0002】**

50

一般に、この種の音声認識装置は、例えばコンピュータなどを利用して、発話を記録した学習用データから音声の特徴を記憶部に蓄積し、入力された話者の音声信号と、記憶部に蓄積された音声の特徴とを比較しながら、最も音声の特徴に近い言語系列を認識結果として出力するものである。ここで比較対象となる音声の特徴は、音声の音響的な特徴と言語的な特徴とに分離され、前者の音響的な特徴は、認識対象の音素がそれぞれどのような周波数特性を持っているかを表わす音声認識用の音響モデルとして、また後者の言語的な特徴は、音素の並び方に関する制約を表わす次単語予測用の言語モデルとして、記憶装置にそれぞれ蓄積される。

【 0 0 0 3 】

音響モデルと言語モデルとを用いたベース判定に基づく音声認識は、近年の標準方式として採用されている。これは図 1 4 に示すように、入力部から取り込んだ話者の音声波形を、特徴量抽出部により周波数軸に対するパワースペクトルの関係に変換して、特徴パラメータ系列に数値化し、パワースペクトル  $S_1 \sim S_5$  の形状の違いによって、次の数 1 に示す定式化されたクラス  $c$  を認識するものである。

【 0 0 0 4 】

【 数 1 】

$$c = \underset{c}{\operatorname{argmax}} P(x|w_c) \times P(w_c)$$

【 0 0 0 5 】

ここで、 $x$  は入力音声（特徴パラメータ系列）、 $w_c$  は単語列を表わし、右辺第 1 項の  $P(x|w_c)$  は音響モデル、右辺第 2 項の  $P(w_c)$  は言語モデルである。コンピュータに備えた認識部（デコーダ）は、入力音声  $x$  が観測されると、上記数式 1 に基づいて、単語列  $w_c$  から入力音声  $x$  が生成される確率と、単語列  $w_c$  が生成される先験的な確率との積が最大となるクラス  $c$  の単語列  $w_c$  を決定し、これを認識結果として出力する。

【 0 0 0 6 】

しかし、上述した特徴パラメータ系列を入力として、例えば HMM（隠れマルコフモデル）でモデル化した音響モデルと、N - g r a m でモデル化した言語モデルとを用いて認識部がデコーディングしても、話者音声を完全には認識できない。その理由は、話者の違いや、音素環境の違いや、周囲騒音や反射物により音が重畳するなどの影響があるからである。

【 0 0 0 7 】

そこで従来の音声認識装置では、話者の違いにより生じる音響変動（話者変動）に対処する目的で、例えば非特許文献 1 に開示されるような、話者適応の技術が一般的に用いられている。話者適応とは、図 1 5（A）にその概略を示すように、特定話者の少量の発話音声データを用いて、不特定話者の音響モデルから当該話者の音響モデルへの写像を推定する手法である。話者適応は、特徴量抽出部から得られる音声の特徴量を音響モデルに合わせる正規化法と、音響モデルのパラメータを変更して音声の特徴量に合わせるモデル適応法に大別でき、何れも事前のデータ収集を必要とする。

【 0 0 0 8 】

前者の正規化法では、代表的な話者正準化手法として、話者による声道長の違いを正規化する声道長正規化法（V T L N : Vocal Tract Length Normalization）が、例えば非特許文献 2 で提案されている。声道長は個人差が大きく、特に男女では大きな差があることが知られており、V T L N では認識対象となる話者の声道長をスペクトルから求め、これを標準的な声道長を持つ話者のスペクトルに周波数軸上で変換することで、声道長の違いに起因する認識性能の低下を防いでいる。

【 0 0 0 9 】

V T L N は、ごく少量のデータで正規化が可能という特徴があるものの、実環境では発声変形や周囲雑音の影響から、声道長の推定精度が低い。そこで従来は、スペクトル変換の際に周波数軸上の伸縮を行なうのに用いるワーピング関数について、そのワーピング関

10

20

30

40

50

数の概形を定義付けるワーピングパラメータが予め複数用意され、話者ごとに観測系列を最大化する最適なワーピングパラメータを選択する  $ML-VTLN$  と呼ばれる手法が、例えば非特許文献 3 で提案されている。

【0010】

図 15 (B) は、 $VTLN$  方式の概略を示したもので、図中、 $f_0$  は変換前の周波数、 $f_1$  (数式以外では、記号の上下に記されたアクセントを、対応する記号の後に併記する) は変換後の周波数、 $\alpha$  はワーピングパラメータである。ここでは、例として声道長を短くしたことに相当する  $\alpha = 0.5$  と、声道長を長くしたことに相当する  $\alpha = -0.25$  のワーピングパラメータだけが図示されているが、 $ML-VTLN$  では  $-0.25$  から  $0.5$  の間に 5 種類程度のワーピングパラメータが用意されており、観測された音声パワースペクトルの周波数軸に対して、5 種類のワーピングパラメータによる伸縮を施し、全てのワーピングパラメータで尤度計算を実行した後に、尤度の高い最適な伸縮を確定する。これにより、図 15 (B) で示すような、人によるパワースペクトル  $S_A$ 、 $S_B$  の形状の違いひいては声道長のずれを矯正する構成となっている。

10

【0011】

$ML-VTLN$  は、同一の話者音声に対して、予め用意されたワーピングパラメータの数の尤度計算が必要となり、声道長の推定精度が高くなる半面、計算量が多くなる。また、用意されたワーピングパラメータが必ずしも話者音声に対して最適なものとして存在するとは限らないことも問題となる。

【0012】

こうした問題に対処するために、例えば非特許文献 4 では、 $HMM$  のパラメータ推定時に計算される占有度数を用いて、話者毎に最適なワーピングパラメータを推定することで、 $ML-VTLN$  よりも計算量を減らした声道長正規化の手法が提案されている。

20

【0013】

しかし、上述した手法の何れも、話者毎に定めたパラメータによって線形の伸縮写像を行っており、母音も子音も全て同一の変換関数によって変換する。一方、非特許文献 5 によれば、声道長の一様な正規化だけでは声道伝達関数の話者間分散を吸収しきれないと言われており、また最適なパラメータは母音毎に異なっている。

【0014】

後者のモデル適応法では、非特許文献 6 のような  $SAT$  (Speaker Adaptive Training : 話者正規化学習) が広く知られている。 $SAT$  は、話者適応のための初期モデルとして、話者間の変動を含む不特定話者モデルではなく、標準的な 1 名の話者モデルを利用して、先ず各学習話者のデータを標準話者のデータに線形変換し、この変換後のデータを用いて学習を行なうもので、作成されたモデルは話者変動が抑制されている。

30

【0015】

以上に説明したように、従来からの音声認識では、声道を一様な断面を持つ管として扱い、かつ話者が同じなら声道形状は同じという仮定の下で考えられている。しかしながら、同じ話者でも声道長は母音により異なり(「い」では長く、「お」では短いなど)、このために最適なワーピングパラメータは異なることになるため、認識性能に限界を持つ。さらに、発話による声道変形は一様ではないため、線形なワーピング関数だけでは表現できないことから、認識性能改善は限定されていた。

40

【先行技術文献】

【非特許文献】

【0016】

【非特許文献 1】篠田浩一、「音声認識における転移学習：話者適応」、人工知能学会誌、2012 年 7 月 1 日、vol. 27、no. 4、pp.359-364

【非特許文献 2】E. Eide (イー・エイデ)、H. Gish (エイチ・ギッシュ)、「A parametric approach to vocal tract length normalization,」「声道長正規化へのパラメトリック手法」、Proc. ICASSP (「信号処理とその応用」に関する国際会議 議事録)、1996 年、pp.346-348

50

【非特許文献3】L. Welling (エル・ ウェリング), S. Kanthak (エス・ カンタック), H. Ney (エイチ・ ネイ), “Improved methods for vocal tract normalization,” 「声道正規化の改良方法」、Proc. ICASSP (「信号処理とその応用」に関する国際会議 議事録)、1999年、pp.1436

【非特許文献4】江森 正, 篠田 浩一, 「音声認識のための高速最尤推定を用いた声道長正規化」、電子情報通信学会論文誌、2000年11月25日、Vol.J83-D2、No.11、pp.2108-2117

【非特許文献5】北村 達也, 竹本 浩典, 足立 整治, 「声道の局所的伸縮による話者正規化」、電子情報通信学会技術研究報告、2010年2月、Vol.109、No.451、pp.57-62

【非特許文献6】T. Anastasakos (ティー・ アナスタサコス), J. McDonough (ジェイ・ マクドナウ), R. Schwartz (アール・ シュワルツ), J. Makhoul (ジェイ・ マコウル), “A compact model for speaker-adaptive training,” 「話者正規化学習向けコンパクトモデル」、ICSLP (音声言語処理国際会議)、1996年、pp.1137-1140

【発明の概要】

【発明が解決しようとする課題】

【0017】

上述したVTLNに基づく音声認識技術は、周波数軸上でのスペクトルの変換が何れも母音に依らず、かつ線形伸縮で行われており、また特にML-VTLNでは、複数の伸縮計算を並行して行った後に、最適なパラメータを選択しなければならず、認識能力がさほど向上しないのにも拘らず、計算量は多くなるという問題があった。

【0018】

一方、SATは標準話者を想定している点がVTLNと共通しているものの、これも母音に依らず線形変換を行なっているために認識能力が低い。

【0019】

そこで本発明は上記問題点に鑑み、母音スペクトルの違いによりワーピングを与える写像関数を変えると共に、一度の計算で話者変動を効果的に抑圧して、全ての入力音声に対して高い認識性能を実現できる音声認識装置、音声認識方法及び音声認識プログラムを提供することを、主な目的とする。

【課題を解決するための手段】

【0020】

本発明の音声認識装置は、任意話者の音声スペクトルを、標準話者の音声スペクトルに変換する話者正準化手段を備えた音声認識装置において、前記標準話者の音響特徴を教師データとして、前記任意話者の音声スペクトルを、前記標準話者の音声スペクトルに周波数軸上で、スペクトル形状の違いに応じて非線形に変換するニューラルネットワークを含む写像関数学習手段を、前記話者正準化手段に備えたことを特徴とする。

【0021】

この場合の前記写像関数学習手段は、前記ニューラルネットワークからの出力を、前記標準話者の音声スペクトルとしてそのまま音声認識に用いる周波数領域を制限する構成とするのが好ましい。

【0022】

また前記写像関数学習手段は、制限された前記周波数領域以外の周波数領域で、前記ニューラルネットワークからの出力と、前記任意話者の音声スペクトルのそれぞれを重み付けして合成し、前記標準話者の音声スペクトルを変換出力する構成とするのが好ましい。

【0023】

さらに、前記多数話者の音声スペクトルをクラスター分析することにより、当該多数話者の中から前記標準話者を特定する標準話者確定手段をさらに備えるのが好ましい。

【0024】

本発明の音声認識方法は、任意話者の音声スペクトルを、標準話者の音声スペクトルに変換して話者正準化を行なう音声認識方法において、前記話者正準化では、ニューラルネットワークを用い、前記標準話者の音響特徴を教師データとして、前記任意話者の音声ス

10

20

30

40

50

ペクトルを、前記標準話者の音声スペクトルに周波数軸上で非線形に変換することを特徴とする。

【0025】

この場合、前記ニューラルネットワークからの出力は、前記標準話者の音声スペクトルとしてそのまま音声認識に用いる周波数領域が制限されているのが好ましい。

【0026】

また、制限された前記周波数領域以外の周波数領域で、前記ニューラルネットワークからの出力と、前記任意話者の音声スペクトルのそれぞれを重み付けして合成し、前記標準話者の音声スペクトルに変換出力するのが好ましい。

【0027】

さらに、前記多数話者の音声スペクトルをクラスター分析することにより、当該多数話者の中から前記標準話者を特定する特徴をさらに備えるのが好ましい。

【0028】

本発明の音声認識プログラムは、上記音声認識方法をコンピュータに実行させることを特徴とする。

【発明の効果】

【0029】

話者正準化のために、任意話者の音声スペクトルを一人の標準話者の音声スペクトルに変換するので、従来のような複数の伸縮計算を並行して行った後に、最適なパラメータを選択する手間を解消できると同時に、決められたパラメータに限定された音声スペクトルの伸縮方法では、性能向上に限界があった点を解消でき、一度の計算でスペクトルが異なることに応じた異なるワーピング関数を実現することにより、話者変動を効果的に抑圧することが可能になる。また、ニューラルネットワーク含む写像関数学習手段を用いて、任意話者の音声スペクトルを標準話者の音声スペクトルに周波数軸上で非線形に変換することから、話者変動の大きい周波数帯域を制限して変換を行なうことで、全ての入力音声に対して高い認識性能を得ることができる。

【0030】

また、標準話者の音声スペクトルを変換出力するのに、ニューラルネットワークからの出力を全周波数に渡ってそのまま音声認識に用いるのではなく、ニューラルネットワークからの全出力の中で、制限された特定の周波数領域の出力だけをそのまま音声認識に用いることで、音声の認識精度をより高めることが可能になる。

【0031】

また、制限された周波数領域以外の周波数領域では、ニューラルネットワークからの出力と、任意話者の音声スペクトルのそれぞれを重み付けして合成し、その結果を標準話者の音声スペクトルとすることで、任意話者の音声スペクトルを恒等写像で標準話者の音声スペクトルに変換するものよりも、音声の認識精度を高めることが可能になる。

【0032】

また、多数話者の音声スペクトルをクラスター分析して、各多数話者の音声スペクトル間の距離を比較することで、多数話者の中から標準話者を精度よく特定できる。

【0033】

さらに、上述した音声認識装置や音声認識方法としての作用効果を、音声認識プログラムにそのまま適用できる。

【図面の簡単な説明】

【0034】

【図1】本発明の第一の実施の形態に係る話者正準化手段を導入した音声認識装置の概略構成図である。

【図2】同上、標準話者確定手段の機能を説明するための模式図である。

【図3】同上、音声データを対数パワースペクトルに変換する機能を説明する模式図である。

【図4】同上、ワーピングを与える写像関数学習手段の内部構成を示す概略図である。

10

20

30

40

50

【図 5】同上、音声認識装置の動作手順を示すフローチャートである。

【図 6】同上、母音「あ」について、任意話者の音声スペクトルを標準話者の音声スペクトルへ正準化することを説明する模式図である。

【図 7】話者正準化によるスペクトル変換の効果を観察する実験の結果として、(A)は正準化前の全スペクトルの分散を示し、(B)は正準化後の全スペクトルの分散を示す波形図である。

【図 8】同上、MLPによる変換前後でのスペクトルの分散を示すグラフである。

【図 9】不特定話者の連続数字認識実験の結果として、単語認識精度を示すグラフである。

【図 10】同上、単語正解率を示すグラフである。

10

【図 11】本発明の第二の実施の形態に係る写像関数学習手段の内部主要構成を示す概略図である。

【図 12】同上、各重み付け関数  $w_{out}(k)$ 、 $w_{in}(k)$  の概念図である。

【図 13】数字音声認識精度を、バイアス値の関数として示すグラフである。

【図 14】従来例において、ベーズ判定に基づく音声認識の手法を示す説明図である。

【図 15】従来手法の概略を説明する図であり、(A)は話者適応方式、(B)は声道長正規化(VTLN)方式を示している。

【発明を実施するための形態】

【0035】

以下、本発明における音声認識装置、音声認識方法及び音声認識プログラムの各実施形態について、添付図面を参照して説明する。なお、添付図面は本発明の技術的特徴を説明するのに用いられており、記載されている装置の構成、各種処理の手順などは、特に特定の記載がない限り、そのみに限定する趣旨ではない。

20

【0036】

先ず、図 1 を参照して、第一の実施の形態に係る音声認識装置 1 の構成を説明する。同図において、音声認識装置 1 は、話者正準化変換前の学習 DB (データベース) 2 および評価データ記憶部 3 と、正準化器 4 と、話者正準化変換後の学習 DB 5 および評価データ記憶部 6 と、音響モデル記憶部 7 と、認識部 8 とにより構成される。なお、学習 DB と評価データ (音声認識の段階にはマイクロホンから入力される) は、図に示されていない音響分析部において 24 チャンネル程度の帯域通過フィルタ (BPF) によってパワースペクトルの時系列として予め変換されている。図の中で、正準化器 4 と認識部 8 は、コンピュータの中央演算処理装置 11 が、以下に説明する処理手順に従い数値演算や制御などの処理を実行することで構成される。また、学習 DB 2、5 や、評価データ記憶部 3、6 や、音響モデル記憶部 7 は、何れも前記中央演算処理装置 11 と接続するコンピュータの記憶装置 12 に設けられる。記憶装置 12 はその他に、中央演算処理装置 11 によって実行される処理手順に対応した正準化器 4 や認識部 8 の音声認識プログラムを格納している。

30

【0037】

中央演算処理装置 11 は、例えば入出力インターフェースを備えた CPU などが使用可能である。また記憶装置 12 は、例えば ROM (リード・オンリー・メモリ) や、RAM (ランダム・アクセス・メモリ) や、HDD (ハードディスクドライブ) などが使用可能である。ここには図示しないが、話者音声の入力を可能にするにマイクロホンなどの入力装置や、例えば認識部 8 で得られた認識結果などの出力を可能にするディスプレイやスピーカなどの出力装置を、中央演算処理装置 11 の入出力インターフェースと接続してもよい。

40

【0038】

なお、本発明における音声認識装置 1 のハードウェア構成は、図 1 に示すものに限定されない。従って、インターネットなどの通信ネットワークを介して、音声認識装置 1 の一部の構成を接続しても構わない。

【0039】

また、本実施形態の音声認識装置 1 と音声認識プログラムは、他のシステムから独立し

50

て設けられているが、本発明はこの構成に限定されない。従って、他の装置の一部として組込まれた構成や、他のプログラムの一部として組込まれた構成とすることも可能である。また、その場合における入出力は、上述の他の装置やプログラムを介して間接的に行われることになる。

#### 【0040】

学習DB2は、学習対象となる多数話者の音声データを含む学習データを格納するもので、ここでの学習データは、例えば前述した入力装置や通信ネットワークを通して収集される。この処理はオフラインで行うため、他の計算システム上で行なうことも可能である。入力装置を利用する場合、多数話者の音声を入力装置に入力する毎に、その音声が入力装置でアナログ電気信号に変換され、コンピュータに備えたA/D変換部（図示せず）に出力される。これを受けて、A/D変換部で変換されたデジタル電気信号を中央演算処理装置11に取り込むことで、中央演算処理装置11で生成された話者毎の音声データを、話者毎の学習データに含めて学習DB2に記憶保存する構成となっている。また、通信ネットワークを利用する場合、音声認識装置1の外部に設置したファイルサーバに、中央演算処理装置11が通信ネットワークを通じてアクセスし、当該ファイルサーバに格納される話者毎の学習データをダウンロードすることで、その学習データを学習DB2に記憶保存する構成となっている。

#### 【0041】

評価データ記憶部3は、認識対象となる話者の音声データを含む評価データを格納するものである。この評価データも、前述した多数話者の学習データを収集する場合と同様の構成で、入力装置や通信ネットワークを通して話者毎に収集される。音声認識装置1としては、評価データに対する処理がオンラインシステムとして使用される。

#### 【0042】

本発明の特徴的部分となる正準化器4は、学習DB2に格納される多数話者の学習データから、セントロイドとなる標準話者を確定する標準話者確定手段21と、標準話者確定手段21で求めた標準話者以外の任意話者について、その音響特徴を標準話者に写像する関数を備えると共に、ニューラルネットワークとしてMLP（多層パーセプトロン）を用いてその写像関数を学習して、話者に対する正準化を実現する写像関数学習手段22と、を備えている。

#### 【0043】

標準話者確定手段21は、学習DB2に格納される全ての多数話者の学習データについて、個々の音声データのスペクトル間距離を比較することで、話者正準化のための標準話者を決定するものである。図2には、標準話者確定手段21の機能を模式的に示しているが、例えば385人分の話者1，話者2，...話者385について、音声データの5種類の母音（例えば、「a（あ）」，「i（い）」，「u（う）」，「e（え）」，「（お）」）の音声データを含む学習データが学習DB2にそれぞれ格納されているとすると、標準話者確定手段21は、各話者1，話者2，...話者385の音声データについて、母音毎に24チャンネルの対数パワースペクトル $X_1, X_2, \dots, X_{24}$ を算出し、Modified K-mean法を用いた対数パワースペクトル $X_1, X_2, \dots, X_{24}$ のクラスター分析によって、母音のそれぞれについて、セントロイドとなる標準話者を決定する機能を有する。ここで決定した標準話者の学習データは、教師データとして写像関数学習手段22に与えられる。

#### 【0044】

図3は、学習DB2や評価データ記憶部3内に存在する任意話者の音声データを、上述した24チャンネルの対数パワースペクトル $X_1, X_2, \dots, X_{24}$ に変換するまでの音声分析機能を模式的に示している。同図において、標準話者確定手段21はまず、学習DB2から読み込んだ時間軸に対する振幅の関係を表す音声データを、所定のサンプリング周波数でサンプリングした後、決められたハミング窓長とフレームシフトで、周波数軸に対するパワースペクトルの関係にフーリエ変換する。フーリエ変換の後、パワースペクトルは中心周波数を聴覚のメルスケールで設定した24チャンネル程度の帯域（メルフィルタバンク）に分割され（各帯域内ではパワースペクトルは加算され）、この後で対数演算したも

10

20

30

40

50



のが音声データのパワースペクトル  $S$  として音声認識に利用される(図 3 (A) 参照)。

【0045】

次に、オフラインで処理される標準話者確定手段 21 は、複数個の BPF (バンドパスフィルタ) を組み合わせたメルフィルタバンク (図示せず) を用いて前述のパワースペクトル  $S$  を分析し、その分析結果を対数化して、前述の対数パワースペクトル  $X_1, X_2, \dots, X_{24}$  を算出する。メルフィルタバンクは、個々の BPF のフィルタ特性を示す三角窓  $T$  が、メル尺度上で等間隔に配置されるように設計されており、本実施形態では 24 チャンネルの BPF からなる対数メルフィルタバンクにパワースペクトル  $S$  を通すことにより、個々のチャンネルで周波数軸に対する対数パワースペクトル  $X_1, X_2, \dots, X_{24}$  が得られる。図 3 (B) は、例として 24 番目のチャンネルの対数パワースペクトル  $X_{24}$  を示している。

10

【0046】

なお、上述した音声分析の機能は、標準話者確定手段 21 のみならず写像関数学習手段 22 も同様に備えている。また、対数メルフィルタバンクを構成する BPF の数は、上述した 24 個に限定されない。

【0047】

図 4 は、ニューラルネットワークによるスペクトル変換を行なう写像関数学習手段 22 の内部構成を示している。同図において、写像関数学習手段 22 は、MLP 31 の入力層 32 と出力層 33 に前述の対数メルフィルタバンクをそれぞれ接続して構成される。ここでは MLP 31 として、入力層 32 と出力層 33 との間に中間層としての隠れ層 34 を接続した一般的な 3 層パーセプトロンを例示しているが、それ以外のニューラルネットワーク構造を採用してもよい。

20

【0048】

MLP 31 は、入力側の対数メルフィルタバンクを通して、任意話者 B の対数パワースペクトル  $X_1, X_2, \dots, X_{24}$  が入力層 32 にそれぞれ与えられる毎に、標準話者 A の学習データから得られる音響特徴を教師データとした写像関数を用いて、標準話者 A に正準化した対数パワースペクトル  $y_1, y_2, \dots, y_{24}$  を算出し、これを出力層 33 から出力側の対数メルフィルタバンクに送出すると共に、当該写像関数を学習する機能を有している。ここでは、前述した 24 チャンネルの対数パワースペクトル  $X_1, X_2, \dots, X_{24}$  を中心フレームとして、その中心フレームの前後  $N$  フレームを結合したものを、MLP 31 の入力層 32 と結合して用い、また同様に、24 チャンネルの対数パワースペクトル  $y_1, y_2, \dots, y_{24}$  を中心フレームとして、その中心フレームの前後  $N$  フレームを結合したものを、MLP 31 の出力層 33 としている。したがって、例えば  $N = 2$  であれば、MLP 31 の入力層 32 と出力層 33 は、24 チャンネル  $\times$  3 フレーム = 72 次元となる。

30

【0049】

話者正準器 4 の出力は、対数メルフィルタバンクの低域と高域を除く帯域で使用し、他の帯域は MLP 31 の入力を帯域境界での平滑処理後に使用する。つまりここでの写像関数学習手段 22 は、変動の大きい周波数帯域である対数メルフィルタバンクの低域と高域で、MLP 31 の入力層 32 に与えられる任意話者 B の音声スペクトルを、恒等変換にてそのまま出力層 33 から出力させる一方で、それ以外の母音スペクトルの話者間変動を反映する周波数帯域である中間域で、MLP 31 の入力層 32 に与えられる任意話者 B の音声スペクトルを、標準話者 A に正準化した音声スペクトルに線形写像して出力層 33 から出力させるような、周波数領域に応じて音声スペクトルを全体で非線形に変換する構成を備えている。これら恒等写像は、母音認識に重要な帯域を除いて、子音認識を安定に行なうことを意図している。実際、このような構成をとることで、音声認識性能は全体として数 % 向上する。

40

【0050】

写像関数学習手段 22 は、MLP 31 の出力側に DCT (離散コサイン変換) による特徴量生成部 (図示せず) が接続される。これにより、学習 DB 2 に格納した音声データについて、標準話者 A に正準化した対数パワースペクトル  $y_1, y_2, \dots, y_{24}$  の MFC

50

C (Mel-Frequency Cepstrum Coefficients : メル周波数ケプストラム係数) 特徴量が、特徴量生成部で生成されて学習DB5に格納され、評価データ記憶部3に格納した音声データについて、標準話者Aに正準化した対数パワースペクトル  $y_1, y_2, \dots, y_{24}$  のMFCC特徴量が、特徴量生成部で生成されて評価データ記憶部6に格納される。

【0051】

学習DB5や評価データ記憶部6は、特徴量生成部で変換された後の音声の特徴量を記憶保持する特徴量記憶部に相当する。本実施形態では、特徴量生成部の構成を適宜変更することで、MFCC特徴量の他に、例えば MFCCや、Pや、Pなどの各種特徴量を得ることができる。

【0052】

音響モデル記憶部7は、学習DB5や評価データ記憶部6に格納した音声の特徴量に基づいて生成される音響モデルを記憶保持するものである。音響モデルとしては、例えば一般的に知られるHMMなどを用いることができる。また認識部8は、音響モデル記憶部7に格納した音響モデルを用い、評価データ記憶部6から読み出した音声の特徴量に対する認識結果を、出力装置に送出するものである。認識部8として、通信ネットワークを介して入手可能な各種の音声認識ソフトウェアを用いてもよい。

【0053】

続いて、上記構成の音声認識装置1について、その動作手順を図5のフローチャートに沿って説明する。

【0054】

まず、オフラインで処理される(A)の学習過程から説明すると、予め学習DB2に多数話者の学習データを格納し、音声認識装置1の動作を開始させる。するとステップT1の手順に移行し、正準化器4に組み込まれた標準話者確定手段21は、学習DB2に格納した学習データに含まれる音声データを個々に読み出し、多数話者の全てについて、音声データから得られる音声スペクトル(対数パワースペクトル  $X_1, X_2, \dots, X_{24}$ )をModified K-mean法によりクラスター分析する。そして、このクラスター分析により、各音声スペクトル間の距離を比較することで、多数話者の中から一人の標準話者を選び出す(ステップT2)。

【0055】

標準話者確定手段21により標準話者を確定すると、音声認識装置1はステップT3の手順に移行し、写像関数学習手段22に組み込まれたニューラルネットワークとしてのMLP31により、任意話者の音声スペクトルを標準話者の音声スペクトル(対数パワースペクトル  $y_1, y_2, \dots, y_{24}$ )に変換する。ここでは特に、最終的な音声の認識性能を向上させるために、任意話者の音声スペクトルの中で、低周波数領域と高周波数領域は恒等変換を行なう一方で、中間の周波数領域は線形変換を行なうような写像関数がMLP31に組み込まれており、任意話者の音声スペクトルは、その全体が標準話者の音声スペクトルに周波数軸上で非線形変換される。

【0056】

ステップT3の手順で、MLP31を通して音声スペクトルの変換が行われると、ステップT4の手順に移行して、MLP31に組み込まれた写像関数が学習され、これを多数話者の全ての学習データについて繰り返し行なうことで、一連の学習過程が終了する。

【0057】

次に、オンラインで処理される(B)の認識過程を説明すると、ステップT11の手順で、マイクロホンから入力された音声、音響分析部により分析され、パワースペクトルの時系列としての評価データが評価データ記憶部3に格納される。次のステップT12では、ステップT11で得られた分析結果が、学習済みのニューラルネットワークであるMLP31に入力され、標準話者のスペクトルに写像変換される。その際、高域周波数帯と低域周波数帯はMLP31の入力をそのまま出力として使用し、中間の周波数帯のみMLP31を通した出力を使用する。こうして非線形に変換された音声スペクトルは、MFCC特徴量に変換され、そのMFCC特徴量が評価データ記憶部6からHMMの音声分類器

10

20

30

40

50

に入力する（ステップ T 1 3）。これにより音響モデルが音響モデル記憶部 7 に格納される。

【 0 0 5 8 】

続くステップ T 1 4 の手順では、評価データ記憶部 6 に格納した M F C C 特徴量を用い、認識対象となる話者音声の認識結果を認識部 8 から出力する。そして、音声認識装置 1 による一連の認識過程を終了する。

【 0 0 5 9 】

ところで、従来は人の声道の部分の大きさ（長さ）の違いに基づいて、標準話者を確定する話者正準化の手法が知られていたが、その手法では、老若男女の違いに基づく 5 ~ 6 人程度の標準話者を設定し、入力する音声スペクトルをその標準話者の何れかの音声スペクトルに合わせて音声認識を行なう必要があった。しかし本実施形態の音声認識装置 1 では、話者正準化手段としてニューラルネットワークによる音声スペクトルの変換を実現しているので、標準話者を一人に設定することができる。

【 0 0 6 0 】

また、話者の音声データとして、特に日本語の発音で母音部分と子音部分では周波数が異なり、一般的には子音部分の周波数領域が高い。一方、男性と女性とを比較した場合には、女性の方が周波数領域は高くなる。入力する音声スペクトルを標準話者の音声スペクトルに当てはめを行なう場合、本実施形態のようなニューラルネットワークを用いることを考えると、人が発音できる周波数領域の中で、子音部分や女性の周波数領域である高音領域と、低音領域では、標準話者の音声スペクトルに上手く当てはめることができない。そこで、こうした高音領域や低音領域では、入力データをそのまま用い、所定の高音領域と低音領域を削除した中間の領域（300 Hz ~ 3 kHz）では、ニューラルネットワークを用いて標準話者の音声スペクトルへの当てはめを行なうことで、最終的な音声の認識率を向上することが可能になる。

【 0 0 6 1 】

次に、話者正準化におけるスペクトル変換の効果を観察した実験について、以下詳しく説明する。この実験では、学習 DB 2 に含まれない話者の音声を認識対象として用い、正準化器 4 によるスペクトル変換の前後で、周波数 bin 毎に分散の平均がどのようになるのかを求める。

【 0 0 6 2 】

実験試料として、学習 DB 2 に格納する多数話者の学習データは、身体情報付きの男・女・子供の母音音声データベース（JVPD）に基づいており、話者数は 385 名である。また、評価データ記憶部 3 に記憶する評価データは、男女 1 名ずつの ATRSetB を用いる。

【 0 0 6 3 】

実験条件として、サンプリング周波数は 16 kHz で、ハミング窓長は 30 ms、フレームシフトは 10 ms である。このとき、24 チャンネルのメルフィルタバンク分析を行なうことで得た対数パワースペクトルを、正準化器 4 の入出力として用いる。

【 0 0 6 4 】

学習 DB 2 に格納した JVPD からの母音音声データから、各母音の中心フレームに前後 2 フレーム分を加えて特徴量を算出する。教師データはスペクトルのピークとして現れるフォルマントを動かすことを考え、メルフィルタバンクの中間周波数領域に相当する 6 ~ 18 bin に、入力データと対応する正準者話者の母音を用いた線形写像変換を行ない、それ以外の低周波数領域と高周波数領域では、入力データと同じものを用いた恒等変換を行なう。

【 0 0 6 5 】

標準話者確定手段 2 1 は、学習 DB 2 に格納した 385 名の話者から、セントロイドとなる話者を求めて標準話者とし、残りの 384 名の話者を学習話者として確定する。MLP 3 1 を学習するための学習データに含まれる音声データの数は、話者の数である 385 に母音の数である 5 を積算した 1920 である。

## 【 0 0 6 6 】

M L P 3 1として使用した3層パーセプトロンは、入出力となる入力層3 2と出力層3 3が、何れも2 4チャンネル×3フレーム＝7 2次元で、隠れ層3 4が1 4 4次元である。

## 【 0 0 6 7 】

こうして、多数話者の学習データを正準化器4に通して標準話者への正準化を行ない、学習DB5に格納した変換後の学習データを用いて音響モデル記憶部7を学習した後、評価データに含まれる音声データについて、母音ラベルの中心フレームに前後2フレーム分を加えたものを切り出し、これをM L P 3 1で変換する。図6は、母音「a（あ）」について、任意話者の音声スペクトルを標準話者の音声スペクトルへ正準化する動作を示している。

10

## 【 0 0 6 8 】

図7（A）は、母音「a（あ）」について、標準話者への正準化を行なう前の全スペクトル（正規化済み）の分散を示し、図7（B）は、同じく母音「a」について、標準話者への正準化を行なった後の全スペクトル（正規化済み）の分散を示している。この実験結果によれば、M L P 3 1による正準化により、スペクトルの分散が小さくなっていることが確認できる。また図8は、M L P 3 1による変換前後でのスペクトルの分散を、母音「a」、「i」、「u」、「e」、「」のそれぞれと、全母音の平均で示したものである。この場合も、正準化後では全ての母音について、スペクトルの分散が低減されたことが分かる。

## 【 0 0 6 9 】

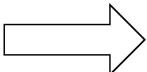
20

また次の表1は、正準化の有無と認識部8による母音正解率との関係を示している。

## 【 0 0 7 0 】

## 【表1】

## 正準化の有無と母音正解率(%)

正準化なし						正準化あり							
		認識結果							認識結果				
		あ	い	う	え	お			あ	い	う	え	お
正 解 ラ ベ ル	あ	93.5	0.0	4.2	0.0	2.3	正 解 ラ ベ ル	あ	97.9	0.0	0.5	0.0	1.6
	い	0.0	95.6	3.4	1.0	0.0		い	0.0	98.2	0.0	1.8	0.0
	う	0.3	0.0	99.0	0.8	0.0		う	0.3	0.0	96.1	3.1	0.5
	え	0.8	2.9	11.5	84.1	0.8		え	0.0	1.0	0.8	97.9	0.3
	お	4.7	0.0	2.9	0.0	92.4		お	2.1	0.0	1.3	0.0	96.6
Ave: 92.9								Ave: 97.3					
						4.4%上昇							

声認識に関しては4.4%の上昇が確認できた。

【0072】

続いて、音声認識装置1の有効性を、不特定話者の連続数字認識実験から評価する。

【0073】

実験試料として、学習DB2に格納する多数話者の学習データは、新聞読み上げコーパスと、ATRSetバランス文(ASJ/JNAS)に基づいており、話者数は男女370名、音声データの数に相当する発話数は70800である。また、評価データ記憶部3に記憶する評価データは、連続数字(CENSREC-4)に基づいており、話者数は男女110名、発話数は8440である。

【0074】

本実験では、学習DB2に格納した学習データや、評価データ記憶部3に格納した評価データに含まれる音声データを、共に標準化器4のMLP31で標準化を行ない、MLP31で変換されたスペクトルのうち、中心フレームの24次元を切り出したものを、前述の特徴量生成部でのDCTによって、音声の特徴量となるMFCCに変換する。学習DB5に格納される特徴量はMFCC、MFCC、P、Pの計26次元を用い、この特徴量からモノフォンの音響モデルを作成して、認識部8による連続数字の認識実験を行った。

【0075】

実験結果として、混合分布数に対する単語認識の正解精度を図9に示し、混合分布数に対する単語正解率を図10に示す。その際、単語挿入ペナルティは方式毎に最適な値を選定した。図9に示すように、「MFCC/HMM 話者標準化MLP」に対応する本実施形態の音声認識装置1では、従来手法である「MFCC/HMM VTLN無し」や「MFCC/HMM VTLN有」よりも、単語認識の正解精度が向上したことが判明した。

【0076】

こうして、本実施形態の音声認識装置1は、話者の違いに伴うスペクトルの分散を低減し、これまで多く用いられてきたVTLNによる話者標準化と比較して、認識部8での音声の認識性能を改善できることが判明した。同時に、本実施形態の音声認識装置1は、VTLNで必要なパラメータの尤度選択の手間を省くことができ、一度の計算でありながら、全ての入力音声に対して高い認識性能を実現できる。

【0077】

以上のように本実施形態では、任意話者の音声スペクトルである対数パワースペクトル $X_1, X_2, \dots, X_{24}$ を、標準話者に標準化した音声スペクトルである対数パワースペクトル $y_1, y_2, \dots, y_{24}$ に変換する話者標準化手段として、標準化器4を備えた音声認識装置1において、標準話者の音響特徴を教師データとして、任意話者の対数パワースペクトル $X_1, X_2, \dots, X_{24}$ を、標準話者に標準化した対数パワースペクトル $y_1, y_2, \dots, y_{24}$ に周波数軸上で、スペクトル形状の違いに応じて非線形に変換するニューラルネットワークとしてのMLP31を含む写像関数学習手段22を、話者標準化器4に備えている。

【0078】

この場合、話者標準化のために、任意話者の対数パワースペクトル $X_1, X_2, \dots, X_{24}$ を、一人の標準話者に標準化した対数パワースペクトル $y_1, y_2, \dots, y_{24}$ に変換するので、従来のような複数の伸縮計算を並行して行った後に、最適なパラメータを選択する手間を解消できると同時に、決められたパラメータに限定された音声スペクトルの伸縮方法では、性能向上に限界があった点を解消でき、一度の計算でスペクトルが異なることに応じた異なるワーピング関数を実現することにより、話者変動を効果的に抑圧することが可能になる。また本実施形態では、ニューラルネットワークとしてのMLP31を含む写像関数学習手段22を用いて、任意話者の対数パワースペクトル $X_1, X_2, \dots, X_{24}$ を、標準話者に標準化した対数パワースペクトル $y_1, y_2, \dots, y_{24}$ に周波数軸上で非線形に変換することから、話者変動の大きい周波数帯域を制限して変換を行なうことで、全ての入力音声に対して高い認識性能を得ることができる。

10

20

30

40

50

## 【 0 0 7 9 】

そしてこれは、コンピュータに組み込まれた正準化器 4 により、任意話者の音声スペクトルである対数パワースペクトル  $X_1, X_2, \dots, X_{24}$  を、標準話者に正準化した音声スペクトルである対数パワースペクトル  $y_1, y_2, \dots, y_{24}$  に変換して話者正準化を行なう音声認識方法において、話者正準化では、任意話者の対数パワースペクトル  $X_1, X_2, \dots, X_{24}$  を、標準話者に正準化した対数パワースペクトル  $y_1, y_2, \dots, y_{24}$  に周波数軸上で、スペクトル形状の違いに応じて非線形に変換するような M L P 3 1 を含む写像関数学習手段 2 2 を用いることでも実現する。

## 【 0 0 8 0 】

なお本実施形態では、母音スペクトラムで学習した M L P 3 1 を含む写像関数学習手段 2 2 を、話者の正準化器 4 として組み込んだことを特徴としているが、M L P 3 1 からの出力 ( B P F の出力でパワースペクトルを表現 ) を、24 チャンネル全てで使用すると、音声の母音部に対しては性能が高いものの、子音部では認識が低下する。そこで、上記非線形変換の具体的な手法として、母音認識に貢献する中間帯域として、好ましくは 3 0 0 H z から 3 k H z の範囲に対応したチャンネルでは、M L P 3 1 の出力を使用し、それ以外の範囲である低域と高域のチャンネルでは、M L P 3 1 の入力をそのまま出力して使用するという恒等写像を用いることで、母音部のみならず子音部に対しても高い認識性能を得ることが可能になる。

## 【 0 0 8 1 】

また、本実施形態の音声認識装置 1 は、学習 D B 2 に格納した学習データから得られる多数話者の音声スペクトルとしての対数パワースペクトル  $X_1, X_2, \dots, X_{24}$  を、好ましくは Modified K-mean 法を用いてクラスター分析することにより、当該多数話者の中から標準話者を特定する標準話者確定手段 2 1 を正準化器 4 に備えている。

## 【 0 0 8 2 】

この場合、多数話者の対数パワースペクトル  $X_1, X_2, \dots, X_{24}$  をクラスター分析して、各多数話者の対数パワースペクトル  $X_1, X_2, \dots, X_{24}$  間の距離を比較することで、多数話者の中から標準話者を精度よく特定できる。

## 【 0 0 8 3 】

そしてこれは、学習 D B 2 に格納した学習データから得られる多数話者の音声スペクトルとしての対数パワースペクトル  $X_1, X_2, \dots, X_{24}$  を、好ましくは Modified K-mean 法を用いてクラスター分析することにより、当該多数話者の中から標準話者を特定する音声認識方法でも実現する。

## 【 0 0 8 4 】

さらに本実施形態では、上述した音声認識方法を、コンピュータの中央演算処理装置 1 1 に実行させるための音声認識プログラムを記憶装置 1 2 に格納している。

## 【 0 0 8 5 】

この場合、上述した音声認識方法としての作用効果を、記憶装置 1 2 に格納した音声認識プログラムにそのまま適用することが可能になる。

## 【 0 0 8 6 】

次に、図 1 1 ~ 図 1 3 を参照しながら、本発明の第二の実施の形態について説明する。第二の実施の形態に係る音声認識装置 1 の構成は、図 1 に示す第一の実施の形態に係る音声認識装置 1 とほぼ同様なものを用いることができる。第一の実施の形態では、正準化器 4 における話者正準化の精度を上げるため、話者変動の大きい周波数帯域を制限している。その方法として、低音領域と高音領域では音声スペクトルの入力データをそのまま用い、それ以外の中間の領域では、M L P 3 1 のようなニューラルネットワークを用いて、標準話者の音声スペクトルへの当てはめを行った。第二の実施の形態では、話者変動の大きい周波数帯域を制限する方法として、低音領域と高音領域における音声スペクトルの入力データをそのまま用いるのではなく、低音領域と高音領域において、音声スペクトルの入力データとニューラルネットワークの出力データとを重み付けして合成する形で処理を行なう。つまりここでは、話者変動の大きい低音領域と高音領域で、M L P 3 1 で変換し出

10

20

30

40

50

力された音声スペクトルと、MLP31への入力をそのまま出力した音声スペクトルとを、両方重み付けした上で重ね合わせて出力するように正準化器4の写像関数学習手段22を構成し、中間音領域以外の低音領域や高音領域の音声スペクトルを、音声の正準化や認識に効果的に利用する。

【0087】

図11は、第二の実施の形態における写像関数学習手段22の内部主要構造を示す概略図である。同図において、本実施形態の写像関数学習手段22は、第一の実施の形態で説明したMLP31に、重み付け加算手段40を追加した構成となっている。また $x_t$ は、第一の実施の形態で説明した特徴量抽出部となるメルフィルタバンクから得られ、写像関数学習手段22への音声スペクトルの入力データとなる対数パワースペクトル、 $y_t$ はMLP31から出力された対数パワースペクトル、 $z_t$ は重み付け加算手段40から出力され、写像関数学習手段22からの音声スペクトルの出力データとなる対数パワースペクトルである。ここでの添え字 $t$ はフレームを示し、それぞれのフレーム $t$ では、前述のような24チャンネルの対数パワースペクトルが存在する。

【0088】

本実施形態では、MLP31と重み付け加算手段40に、メルフィルタバンクからの対数パワースペクトル $x_t$ がそれぞれ入力される構成となっている。MLP31は第一の実施の形態と同様に、入力層32に対数パワースペクトル $x_t$ が与えられると、MLP31からの写像出力となる対数パワースペクトル $y_t$ を、出力層33から重み付け加算手段40に出力する。重み付け加算手段40は、入力する対数パワースペクトル $x_t$ 、 $y_t$ に対し、周波数に応じて異なる重み付けを行ない、それぞれの重み付け出力を加算した後に、メルフィルタバンクを通して最終的に正準化された対数パワースペクトル $z_t$ を出力する。そして重み付け加算手段40の出力側には、第一の実施の形態で説明したDCTが接続され、このDCTで対数パワースペクトル $z_t$ のMFCC特徴量が生成される。なお、図11に示す写像関数学習手段22以外の構成は、第一の実施の形態と共通する。

【0089】

前記MLP31を用いた写像出力は、中間周波数領域の正準化スペクトラムとなる対数パワースペクトル $y_t$ として利用できるが、このMLP31を用いた写像出力による高周波領域や低周波領域での悪影響は、同時に減らさなければならない。そのため本実施形態では、MLP31を用いた写像の出力に利用する重み付け関数が、こうした悪影響を低減できるように、重み付け加算手段40を構成しなければならない。その代わりに、これらの周波数領域における重み付けされた出力は、重み付け加算手段40でMLP31を用いない対数パワースペクトル $x_t$ と統合される。この場合の重み付け加算手段40は、低周波領域と高周波領域の各信号の両方を通過するフィルタとしての周波数重み付け関数を、MLP31を用いない写像である対数パワースペクトル $x_t$ にも利用しなければならない。こうして、メルフィルタバンクからの各チャンネルの正準化された出力スペクトル $z_t(k)$ は、以下の数式で表現することができる。

【0090】

【数2】

$$z_t(k) = w^{\text{out}}(k) \cdot y_t(k) + w^{\text{in}}(k) \cdot x_t(k)$$

$$w^{\text{out}}(k) + w^{\text{in}}(k) = 1.0 \text{ for } \forall k$$

【0091】

数式2において、 $t$ と $k$ はそれぞれフレームと対数メルフィルタバンクを表す添え字である。例えば $x_t(k)$ は、図11の対数パワースペクトル $x_t$ に相当するもので、フレーム $t$ で抽出された24次元対数メルフィルタバンクの $k$ 番目のチャンネル成分を表している。同様に $y_t(k)$ は、図11の対数パワースペクトル $y_t$ に相当するもので、フレーム $t$ で抽出された対数メルフィルタバンクの $k$ 番目のチャンネル成分を表す。

【0092】

重み付け加算手段40からの出力スペクトル $z_t(k)$ は、図11の対数パワースペク

10

20

30

40

50

トル  $z_t$  に相当するもので、これは M L P 3 1 で変換された出力スペクトル  $y_t(k)$  と、M L P 3 1 で変換されない出力スペクトル  $x_t(k)$  とを、それぞれ周波数重み付け関数  $w^{out}(k)$  と  $w^{in}(k)$  で重み付けして加算したものである。ここでの  $w^{out}(k)$  は、M L P 3 1 を用いた対数メルフィルタバンクからの各チャネルの出力スペクトル  $y_t(k)$  に対する周波数重み付け関数の  $k$  番目の成分を表す。また  $w^{in}(k)$  は、M L P 3 1 を用いない対数メルフィルタバンクからの各チャネルの出力スペクトル  $x_t(k)$  に対する周波数重み付け関数の  $k$  番目の成分を表す。

#### 【0093】

図 1 2 は、前記周波数重み付け関数  $w^{out}(k)$  および  $w^{in}(k)$  の概念図である。図中、実線は周波数重み付け関数  $w^{out}(k)$  を示し、破線は周波数重み付け関数  $w^{in}(k)$  を示す。周波数重み付け関数  $w^{out}(k)$  および  $w^{in}(k)$  は、周波数重み付け関数  $w^{out}(k)$  および  $w^{in}(k)$  の形状を形成するために、パラメータ  $k_L, k_H, g$  がそれぞれ含まれる。これらのパラメータ  $k_L, k_H, g$  は、重み付け加算手段 40 に組み込まれたフィルタの性能を決定するものである。

#### 【0094】

パラメータ  $k_L$  は、高周波チャネルおよび低周波チャネルにおいて、M L P 3 1 からの出力に対する重みのバイアス値である。パラメータ  $k_L$  および  $k_H$  は、それぞれ周波数中間領域の下限と上限の代表値である。パラメータ  $g$  は、線形補間する代表値  $k_L$  および  $k_H$  において、周波数重み付け関数  $w^{out}(k)$  および  $w^{in}(k)$  の勾配を表す。ここで、バイアス値として用いる記号  $k_L$  は、前述のワーピングパラメータを示す記号  $k_L$  と同一であるが、この 2 つを混同しないように注意されたい。

#### 【0095】

同図で示されている例において、周波数の中間領域（チャネル番号 1 2 の付近）では、M L P 3 1 で変換された出力スペクトル  $y_t(k)$  を、すべて正準化された出力スペクトル  $z_t(k)$  として用いている（ $w^{out}(k) = 1, w^{in}(k) = 0$ ）。これに対して、周波数の低音領域と高音領域においては、数 2 に示す重み付け出力  $w^{out}(k) \cdot y_t(k)$  の割合を減らし、その減った分を補うように、数 2 に示す重み付け出力  $w^{in}(k) \cdot x_t(k)$  の割合を増やす。そして、これらの重み付け出力  $w^{out}(k) \cdot y_t(k)$  と重み付け出力  $w^{in}(k) \cdot x_t(k)$  を加算した値を、重み付け加算手段 40 からの出力スペクトル  $z_t(k)$  として出力する。周波数の低音領域および高音領域において、正準化された出力スペクトル  $z_t(k)$  に含まれる重み付け出力  $w^{out}(k) \cdot y_t(k)$  と重み付け出力  $w^{in}(k) \cdot x_t(k)$  の割合は、バイアス値  $k_L, k_H$  によって調節可能である。

#### 【0096】

図 1 2 に示すように、重み付け加算手段 40 で設定されるバイアス値  $k_L, k_H$  は、最低周波数のチャネルおよび最高周波数のチャネルにおける周波数重み付け関数  $w^{out}(k)$  の値として定義することができる。このため、 $k_L = 0$  とした場合は、第一の実施の形態に相当するものとなり、M L P 3 1 を用いない対数メルフィルタバンクからの出力スペクトル  $x_t(k)$  が、そのまま恒等写像により正準化された出力スペクトル  $z_t(k)$  として出力される。一方、 $k_L = 1$  とした場合は、すべて M L P 3 1 の出力を用いた形態に相当するものとなり、M L P 3 1 を用いて変換された出力スペクトル  $y_t(k)$  が、すべて正準化された出力スペクトル  $z_t(k)$  として出力される。つまり、本実施形態は第一の実施の形態の拡張となっており、本実施形態の特別な場合（ $k_L = 0$ ）が、第一の実施の形態となっていることに注目されたい。

#### 【0097】

なお図 1 2 に示す例では、周波数の中間領域と低音領域との境界  $k_L$  と、高音領域との境界  $k_H$  のそれぞれについて、前後 2 チャネルを境界部分の幅（ $k_L - 2 \sim k_L + 2, k_H - 2 \sim k_H + 2$ ）とし、当該幅で線形補間を行っている。このため、境界  $k_L$  および  $k_H$  を含む所定の幅において、前述の勾配  $g$  は以下のように設定される。

#### 【0098】

10

20

30

40

50



【数 3】

$$g = \pm(1-\alpha)/4$$

【0099】

図13は、本実施形態を採用した不特定話者連続数字認識システムの実験結果となる数字音声認識精度を、バイアス値の関数として表したものである。このバイアス値は、図12で示されているように、低周波チャンネルや高周波チャンネルにおける $w_{out}(k)$ の値を表現したものになっている。前述のように、 $\alpha = 0$ は、周波数の低音領域と高音領域で、MLP31を全く用いていないことを意味し、 $\alpha = 1$ は、すべての周波数領域でMLP31を用いたことを意味する。

10

【0100】

実験の結果から、重み付け加算手段40でバイアス値を0.1付近に設定した場合に、数字音声の認識精度が最も高くなることがわかる。つまり本実施形態において、バイアス値を0.1付近に調整した場合、その認識精度は、第一の実施の形態( $\alpha = 0$ )よりも向上することがわかる。この結果から、周波数の重み付けを行なうことによって、子音データを学習させることなしに話者の正準化を正確に行なうことが可能であることがわかり、また子音などの周波数が低音や高音の領域で、MLP31の出力を少し混ぜること( $\alpha = 0.1$ の場合に相当)で、数字音声の認識精度にさらなる改善が望めることがわかった。

20

【0101】

但し、音声の認識精度が最も良好なバイアス値は、図12に示す周波数重み付け関数 $w_{out}(k)$ や、周波数重み付け関数 $w_{in}(k)$ モデルのバリエーションによって変化するものであり、 $\alpha = 0.1$ 付近に限られたものではない。また、境界 $k_L$ および $k_H$ を含む境界部分の幅の決め方や、その境界部分での補間の方法は、図12で示されているものに限られない。その他、各パラメータ $\alpha$ 、 $k_L$ 、 $k_H$ 、 $g$ も、音声の認識精度を高めるのに適宜最適な値に設定してよい。

【0102】

以上のように、本実施形態の音声認識装置1も図11に示すように、標準話者の音響特徴を教師データとして、任意話者の対数パワースペクトル $x_t$ を、標準話者に正準化した対数パワースペクトル $z_t$ に周波数軸上で、スペクトル形状の違いに応じて非線形に変換するニューラルネットワークとしてのMLP31を含む写像関数学習手段22を、話者正準化器4に備えている。それに加えて、ここでの写像関数学習手段22は、MLP31からの出力である対数パワースペクトル $y_t$ を、正準化された対数パワースペクトル $z_t$ としてそのまま音声認識に用いる周波数領域を制限する構成を有し、その周波数領域はパラメータ $k_L$ および $k_H$ に基づき中間領域として設定している。

30

【0103】

この場合、写像関数学習手段22で対数パワースペクトル $z_t$ を変換出力するのに、MLP31からの対数パワースペクトル $y_t$ を全周波数に渡ってそのまま音声認識に用いるのではなく、MLP31からの全ての対数パワースペクトル $y_t$ の中で、制限された特定の周波数中間領域の対数パワースペクトル $y_t$ だけをそのまま音声認識に用いることで、音声の認識精度をより高めることが可能になる。

40

【0104】

そしてこれは、MLP31からの対数パワースペクトル $y_t$ を、標準話者に正準化した対数パワースペクトル $z_t$ としてそのまま音声認識に用いる周波数領域が制限されており、その周波数領域をパラメータ $k_L$ および $k_H$ に基づき中間領域として設定する音声認識方法でも実現する。

【0105】

また、本実施形態の写像関数学習手段22は、パラメータ $k_L$ および $k_H$ に基づき制限された周波数中間領域以外の高音領域や低音領域で、MLP31からの出力である対数パワースペクトル $y_t$ と、任意話者の対数パワースペクトル $x_t$ のそれぞれを、周波数に応

50

じて設定される周波数重み付け関数  $w^{out}(k)$  および  $w^{in}(k)$  で重み付けして、それらの出力を合成し、標準話者に正準化した対数パワースペクトル  $z_t$  を変換出力する構成となっている。

【0106】

この場合、制限された周波数中間領域以外の高音領域や低音領域では、MLP31からの出力である対数パワースペクトル  $y_t$  と、任意話者の対数パワースペクトル  $x_t$  のそれぞれを、周波数重み付け関数  $w^{out}(k)$  および  $w^{in}(k)$  で重み付けした上で合成し、その結果を標準話者に正準化した対数パワースペクトル  $z_t$  をとすることで、任意話者の対数パワースペクトル  $x_t$  を恒等写像で標準話者に正準化した対数パワースペクトル  $z_t$  に変換するものよりも、音声の認識精度を高めることが可能になる。

10

【0107】

そしてこれは、制限された周波数中間領域以外の高音領域や低音領域で、MLP31からの出力である対数パワースペクトル  $y_t$  と、任意話者の対数パワースペクトル  $x_t$  のそれぞれを、周波数重み付け関数  $w^{out}(k)$  および  $w^{in}(k)$  で重み付けして、それらの出力を合成し、標準話者に正準化した対数パワースペクトル  $z_t$  に変換出力する音声認識方法でも実現する。

【0108】

以上、本発明の実施形態について説明したが、当該実施形態はあくまでも例として提示したに過ぎず、発明の範囲を限定することを意図していない。ここに提示した実施形態は、その他の様々な形態で実施可能であり、発明の要旨を逸脱しない範囲で、種々の省略、置換、変更が可能である。例えば、話者正準化を実現するMLPへの入力は、実施形態では出力の標準話者と同じ24チャンネルとしているが、出力のチャンネル数は同じ24チャンネルとしてMLPの入力を24チャンネルよりも分析を細かくすることで(40チャンネルなど)、性能が向上することなどが確認されている。また、図1では次単語予測用の言語モデルの構成を省略しているが、この言語モデルを音声認識装置1に付加することで、認識部8の認識結果として文字列を適宜出力することができる。また本発明を、文音声認識など他のタスクに適用してもよい。

20

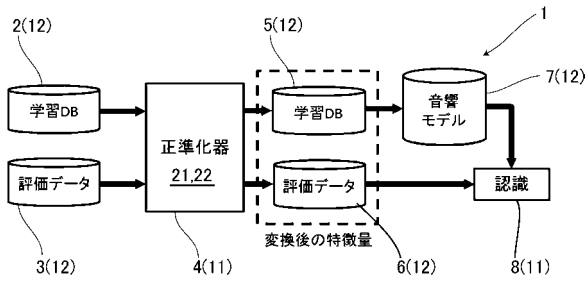
【符号の説明】

【0109】

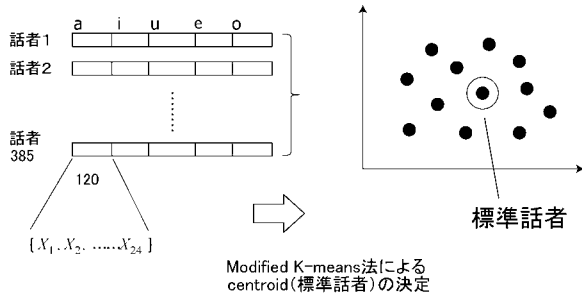
- 1 音声認識装置
- 4 正準化器(話者正準化手段)
- 21 標準話者確定手段
- 22 写像関数学習手段
- 31 MLP(ニューラルネットワーク)

30

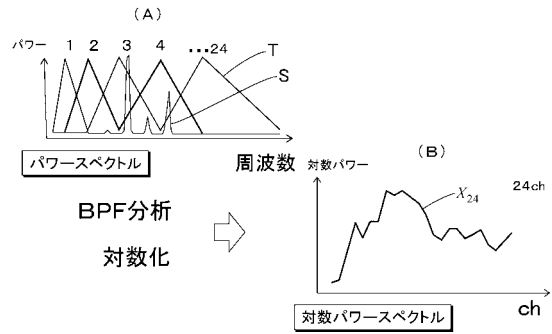
【図1】



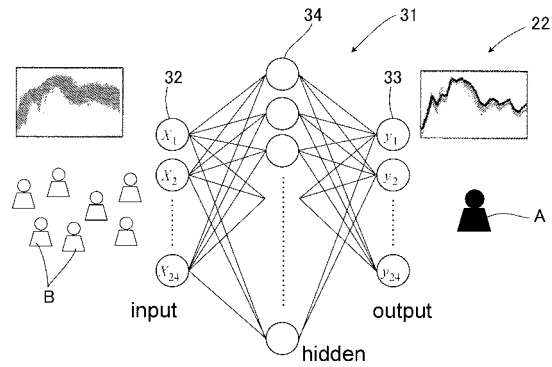
【図2】



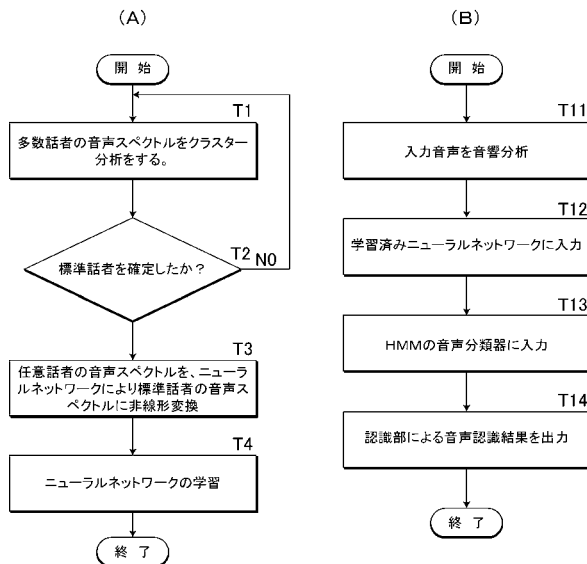
【図3】



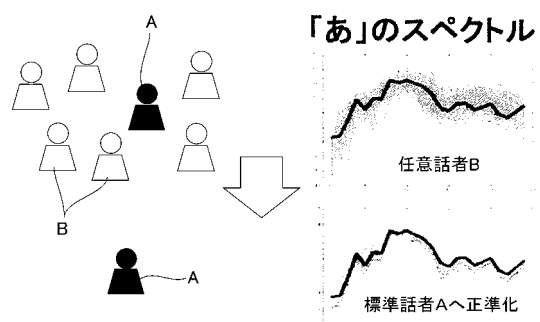
【図4】



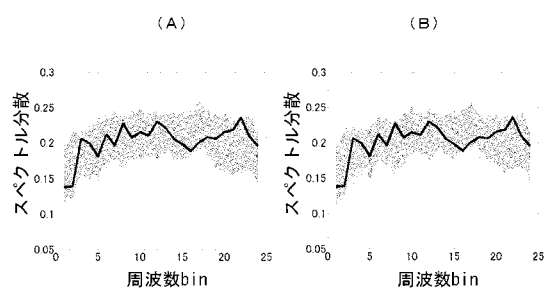
【図5】



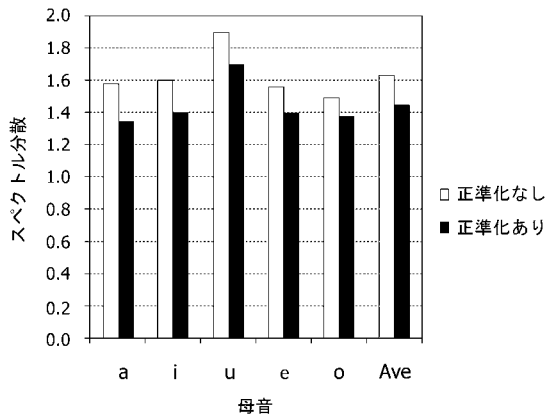
【図6】



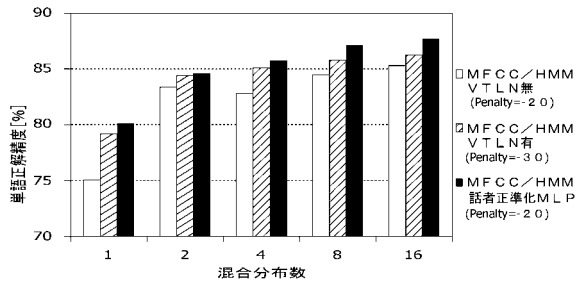
【図7】



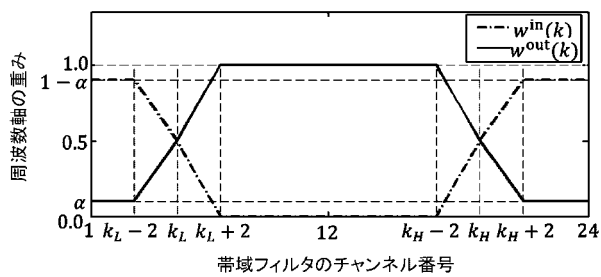
【図 8】



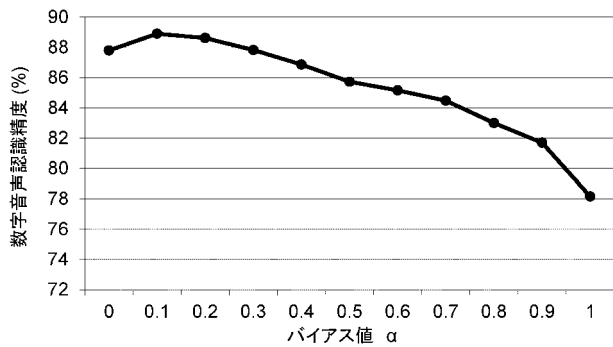
【図 9】



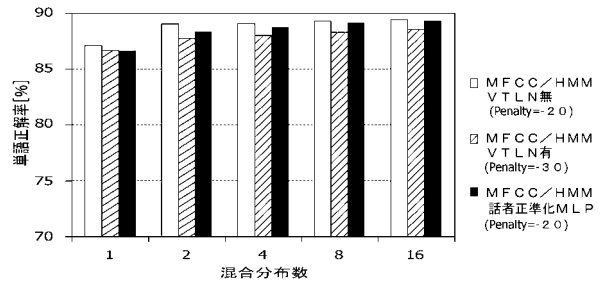
【図 12】



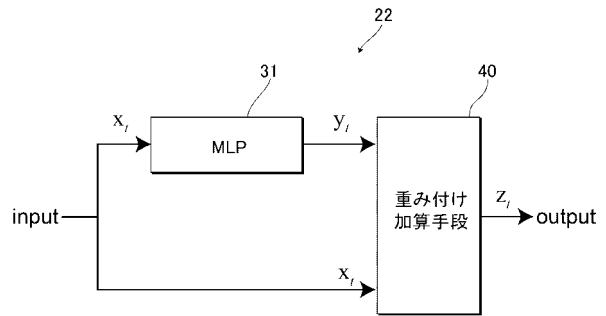
【図 13】



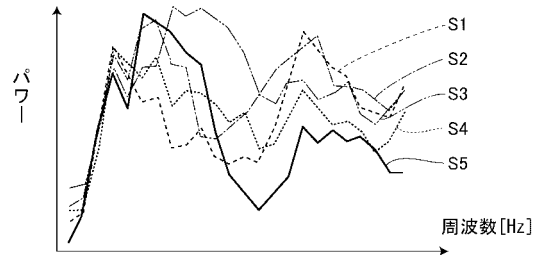
【図 10】



【図 11】



【図 14】



【図 15】

