

(51) Int.Cl.	F I	テーマコード(参考)
G 1 0 L 15/06 (2013.01)	G 1 0 L 15/06	3 0 0 Y
G 1 0 L 15/16 (2006.01)	G 1 0 L 15/16	
G 1 0 L 15/10 (2006.01)	G 1 0 L 15/10	5 0 0 N
G 1 0 L 25/63 (2013.01)	G 1 0 L 25/63	

審査請求 未請求 請求項の数8 O L (全16頁)

(21)出願番号	特願2019-21332(P2019-21332)	(71)出願人	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31年2月8日(2019.2.8)	(74)代理人	100121706 弁理士 中尾 直樹
		(74)代理人	100128705 弁理士 中村 幸雄
		(74)代理人	100147773 弁理士 義村 宗洋
		(72)発明者	安藤 厚志 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
		(72)発明者	神山 歩相名 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】 パラ言語情報推定モデル学習装置、パラ言語情報推定装置、およびプログラム

(57) 【要約】

【課題】パラ言語情報を特定することが困難である発話をモデル学習に用いた場合であっても、高精度にパラ言語情報を推定する。

【解決手段】音響特徴抽出部11は、発話から音響特徴を抽出する。逆教師決定部12は、複数の聴取者が発話ごとに付与したパラ言語情報の判定結果を表すパラ言語情報ラベルに基づいて、その発話のパラ言語情報の正解ではない逆教師を表す逆教師ラベルを決定する。逆教師推定モデル学習部13は、発話から抽出した音響特徴と逆教師ラベルとに基づいて、入力された音響特徴に対する逆教師の事後確率を出力する逆教師推定モデルを学習する。

【選択図】図2

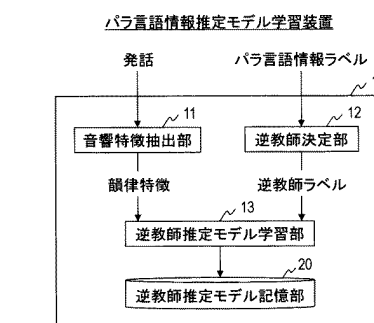


図2

1

## 【特許請求の範囲】

## 【請求項 1】

複数の聴取者が発話ごとに付与したパラ言語情報の判定結果を表すパラ言語情報ラベルに基づいて、その発話のパラ言語情報の正解ではない逆教師を表す逆教師ラベルを決定する逆教師決定部と、

上記発話から抽出した音響特徴と上記逆教師ラベルとに基づいて、入力された音響特徴に対する逆教師の事後確率を出力する逆教師推定モデルを学習する逆教師推定モデル学習部と、

を含むパラ言語情報推定モデル学習装置。

## 【請求項 2】

請求項 1 に記載のパラ言語情報推定モデル学習装置であって、

上記パラ言語情報ラベルに基づいて、その発話のパラ言語情報の正解である従来教師を表す従来教師ラベルを決定する従来教師決定部と、

上記発話から抽出した音響特徴と上記従来教師ラベルとに基づいて、入力された音響特徴に対する従来教師の事後確率を出力する従来教師推定モデルを学習する従来教師推定モデル学習部と、

をさらに含むパラ言語情報推定モデル学習装置。

## 【請求項 3】

複数の聴取者が発話ごとに付与したパラ言語情報の判定結果を表すパラ言語情報ラベルに基づいて、その発話のパラ言語情報の正解ではない逆教師を表す逆教師ラベルを決定する逆教師決定部と、

上記パラ言語情報ラベルに基づいて、その発話のパラ言語情報の正解である従来教師を表す従来教師ラベルを決定する従来教師決定部と、

上記発話から抽出した音響特徴と上記逆教師ラベルと上記従来教師ラベルとに基づいてマルチタスク学習を行い、入力された音響特徴に対する逆教師の事後確率と従来教師の事後確率とを出力するマルチタスク推定モデルを学習するマルチタスク推定モデル学習部と、  
を含むパラ言語情報推定モデル学習装置。

## 【請求項 4】

請求項 3 に記載のパラ言語情報推定モデル学習装置であって、

上記マルチタスク推定モデルは、従来教師の事後確率を出力する深層学習に基づくパラ言語情報推定モデルに対して逆教師の事後確率を出力する分岐構造を加えたモデルである、

パラ言語情報推定モデル学習装置。

## 【請求項 5】

請求項 1 に記載のパラ言語情報推定モデル学習装置が学習した逆教師推定モデルを記憶する逆教師推定モデル記憶部と、

入力発話から抽出した音響特徴を上記逆教師推定モデルに入力して得られる逆教師の事後確率に基づいて上記入

2

力発話のパラ言語情報を推定するパラ言語情報推定部と

、  
を含むパラ言語情報推定装置。

## 【請求項 6】

請求項 5 に記載のパラ言語情報推定装置であって、請求項 2 に記載のパラ言語情報推定モデル学習装置が学習した従来教師推定モデルを記憶する従来教師推定モデル記憶部をさらに含み、

上記パラ言語情報推定部は、上記音響特徴を上記逆教師推定モデルに入力して得られる逆教師の事後確率と上記音響特徴を上記従来教師推定モデルに入力して得られる従来教師の事後確率との重み付き差に基づいて上記入力発話のパラ言語情報を推定する、  
パラ言語情報推定装置。

## 【請求項 7】

請求項 3 または 4 に記載のパラ言語情報推定モデル学習装置が学習したマルチタスク推定モデルを記憶するマルチタスク推定モデル記憶部と、

入力発話から抽出した音響特徴を上記マルチタスク推定モデルに入力して得られる従来教師の事後確率に基づいて上記入力発話のパラ言語情報を推定するパラ言語情報推定部と、  
を含むパラ言語情報推定装置。

## 【請求項 8】

請求項 1 から 4 のいずれかに記載のパラ言語情報推定モデル学習装置もしくは請求項 5 から 7 のいずれかに記載のパラ言語情報推定装置としてコンピュータを機能させるためのプログラム。

## 【発明の詳細な説明】

## 【技術分野】

## 【0001】

本発明は、音声からパラ言語情報を推定する技術に関する。

## 【背景技術】

## 【0002】

音声からパラ言語情報（例えば、感情が喜び・悲しみ・怒り・平静のいずれか）を推定する技術が求められている。パラ言語情報は、音声対話における話し相手の感情を考慮した対話制御（例えば、相手が怒っていれば話題を変える、など）や、音声をを用いたメンタルヘルス診断（例えば、毎日の音声を収録し、悲しみや怒り音声の頻度からメンタルヘルス状況を予測する、など）に応用可能である。

## 【0003】

従来技術として、機械学習に基づくパラ言語情報推定技術が非特許文献 1 に示されている。非特許文献 1 では、図 1 に示すように、音声から抽出した短時間ごとの音響特徴（例えば、声の高さなど）の時系列情報を入力とし、話者のパラ言語情報を推定する。このとき、再帰型ニューラルネットワーク（Recurrent Neural Network: RN

10

20

30

40

50

N)と注意機構と呼ばれる機能を組み合わせた**深層学習**に基づく推定モデルを用いており、音声の部分的な特性に基づいてパラ言語情報を推定することが可能となっている(例えば、発話末尾で声の大きさが急激に低くなったことから、悲しみ感情であると推定することができる)。近年では、非特許文献1のような**深層学習**に基づくパラ言語情報推定モデルが主流となっている。

#### 【0004】

なお、従来技術では、ある音声を複数の聴取者が聴取し、聴取者の過半数がその音声に対して特定の**パラ言語情報**を感じた場合にのみ、その特定の**パラ言語情報**を正解の**パラ言語情報**とみなす。従来技術では、この正解の**パラ言語情報**を推定するように学習を行う。

#### 【先行技術文献】

#### 【非特許文献】

#### 【0005】

【非特許文献1】S. Mirsamadi, E. Barsoum, and C. Zhang, "Automatic speech emotion recognition using recurrent neural networks with local attention," in Proc. of ICASSP, 2017, pp. 2227-2231.

#### 【発明の概要】

#### 【発明が解決しようとする課題】

#### 【0006】

しかしながら、従来技術を用いても、**パラ言語情報推定**の精度が不十分な場合がある。

これは、従来技術では音声に対してただ一つの正解となる**パラ言語情報**を特定しようと推定モデルを学習するが、正解となる**パラ言語情報**の特定は人間でも困難な課題であるためである。例えば、**パラ言語情報推定**のうち感情推定(例えば、喜び・悲しみ・怒り・平静のいずれかを推定する問題)では、従来技術は、ある音声と正解感情との組を学習データとして感情推定モデルを学習する。しかしながら、実際には、正解感情の特定が困難である発話が多く存在する。例えば、3名が聴取した際に、2名が「喜び」、1名が「平静」と判定するような発話が存在する(この場合、従来技術では「喜び」を正解感情とする)。このような発話から正解感情(すなわち「喜び」)に固有の特性を学習することは困難である。この結果、推定モデルを正しく学習することが困難となり、**パラ言語情報推定精度**が低下するおそれがある。

#### 【0007】

本発明の目的は、上述のような技術的課題に鑑みて、**パラ言語情報**を特定することが困難である発話をモデル学習に用いた場合であっても、高精度に**パラ言語情報**を推定することである。

#### 【課題を解決するための手段】

#### 【0008】

本発明の第一の態様の**パラ言語情報推定モデル学習装置**は、複数の聴取者が発話ごとに付与した**パラ言語情報**の判定結果を表す**パラ言語情報ラベル**に基づいて、その発

話の**パラ言語情報**の正解ではない逆教師を表す逆教師ラベルを決定する逆教師決定部と、発話から抽出した音響特徴と逆教師ラベルとに基づいて、入力された音響特徴に対する逆教師の事後確率を出力する逆教師推定モデルを学習する逆教師推定モデル学習部と、を含む。

#### 【0009】

本発明の第二の態様の**パラ言語情報推定モデル学習装置**は、複数の聴取者が発話ごとに付与した**パラ言語情報**の判定結果を表す**パラ言語情報ラベル**に基づいて、その発話の**パラ言語情報**の正解ではない逆教師を表す逆教師ラベルを決定する逆教師決定部と、**パラ言語情報ラベル**に基づいて、その発話の**パラ言語情報**の正解である従来教師を表す従来教師ラベルを決定する従来教師決定部と、発話から抽出した音響特徴と逆教師ラベルと従来教師ラベルとに基づいてマルチタスク学習を行い、入力された音響特徴に対する逆教師の事後確率と従来教師の事後確率とを出力するマルチタスク推定モデルを学習するマルチタスク推定モデル学習部と、を含む。

#### 【0010】

本発明の第三の態様の**パラ言語情報推定装置**は、第一の態様の**パラ言語情報推定モデル学習装置**が学習した逆教師推定モデルを記憶する逆教師推定モデル記憶部と、入力発話から抽出した音響特徴を逆教師推定モデルに入力して得られる逆教師の事後確率に基づいて入力発話の**パラ言語情報**を推定する**パラ言語情報推定部**と、を含む。

#### 【0011】

本発明の第四の態様の**パラ言語情報推定装置**は、第二の態様の**パラ言語情報推定モデル学習装置**が学習したマルチタスク推定モデルを記憶するマルチタスク推定モデル記憶部と、入力発話から抽出した音響特徴をマルチタスク推定モデルに入力して得られる従来教師の事後確率に基づいて入力発話の**パラ言語情報**を推定する**パラ言語情報推定部**と、を含む。

#### 【発明の効果】

#### 【0012】

本発明によれば、**パラ言語情報**を特定することが困難である発話をモデル学習に用いた場合であっても、高精度に**パラ言語情報**を推定することができる。

#### 【図面の簡単な説明】

#### 【0013】

【図1】図1は、従来の**パラ言語情報推定モデル**を説明するための図である。

【図2】図2は、第一実施形態の**パラ言語情報推定モデル学習装置**の機能構成を例示する図である。

【図3】図3は、第一実施形態の**パラ言語情報推定モデル学習方法**の処理手順を例示する図である。

【図4】図4は、第一実施形態の**パラ言語情報推定装置**の機能構成を例示する図である。

【図5】図5は、第一実施形態の**パラ言語情報推定方法**の処理手順を例示する図である。

5

【図6】図6は、第二実施形態のパラ言語情報推定モデル学習装置の機能構成を例示する図である。

【図7】図7は、第二実施形態のパラ言語情報推定モデル学習方法の処理手順を例示する図である。

【図8】図8は、第二実施形態のパラ言語情報推定装置の機能構成を例示する図である。

【図9】図9は、第二実施形態のパラ言語情報推定方法の処理手順を例示する図である。

【図10】図10は、第三実施形態のパラ言語情報推定モデル学習装置の機能構成を例示する図である。

【図11】図11は、第三実施形態のパラ言語情報推定モデル学習方法の処理手順を例示する図である。

【図12】図12は、第三実施形態のパラ言語情報推定装置の機能構成を例示する図である。

【図13】図13は、第三実施形態のパラ言語情報推定方法の処理手順を例示する図である。

【図14】図14は、第三実施形態のパラ言語情報推定モデルを説明するための図である。

【発明を実施するための形態】

【0014】

文中で使用される記号「 $\wedge$ 」は、本来直後の文字の真上に記載されるべきものであるが、テキスト記法の制限により、当該文字の直前に記載する。数式中においてはこれらの記号は本来の位置、すなわち文字の真上に記述している。例えば、「 $\wedge c$ 」は数式中では次式で表される。

【0015】

【数1】

$$\hat{c}$$

【0016】

以下、本発明の実施の形態について詳細に説明する。なお、図面中において同じ機能を有する構成部には同じ番号を付し、重複説明を省略する。

【0017】

[発明のポイント]

本発明のポイントは、あえて「絶対に正解ではないパラ言語情報」を推定することで、正解となるパラ言語情報の特定に貢献する点にある。人間にとってパラ言語情報は、ただ一つの正解の特定は困難である一方で、絶対に正解ではないものの推定は容易であると考えられる。例えば、人間がある音声を聴取した際に、喜びか平静かを特定することは困難な場合があるが、そのような音声でも「怒りではない」「悲しみでもない」と判断することは可能であることが多い。このことから、絶対に正解ではないパラ言語情報の推定は正解のパラ言語情報の特定に比べて容易である可能性があり、高い精度で絶対に正解ではないパラ言語情報の推定を行うことができると考えられる。また、消去法のような枠組みを用いることで、絶対に正解ではないパラ言語情報がわかることは、ただ一つの正解となるパラ言語情報の特定にも貢献できる。以降では、ただ一つの正解となるパラ言語情報を「従

6

来教師」、絶対に正解ではないパラ言語情報を「逆教師」と呼ぶ。

【0018】

上記の発明のポイントを実現するために、後述する実施形態は下記のように構成される。

【0019】

1. 複数名の聴取者によるパラ言語情報の判定結果に基づき、逆教師を決定する。本発明では、逆教師とは、推定対象となるパラ言語情報のうち、一定以下（例えば1割以下）の聴取者が判定したパラ言語情報を指すものとする。例えば、喜び・悲しみ・怒り・平静の4クラスの感情推定において、ある音声に対して3名の聴取者が「喜び」「喜び」「平静」と判定した場合、その音声の逆教師は「悲しみ」「怒り」の2クラスを指す。

【0020】

2. 逆教師の推定モデルを学習する。この推定モデルは、入力特徴量と推定モデル構造が従来技術と同様であるが、最終的な推定部がマルチラベル分類（一つの音声は複数のクラスに同時に分類され得る）の構造を持つモデルを利用することで実現される。

【0021】

3. 逆教師の推定モデル単独、または逆教師の推定モデルと従来教師の推定モデルの両方を用いてパラ言語情報を推定する。逆教師の推定モデルを単独で用いる場合であれば、逆教師の推定モデルによる逆教師推定を行い、その出力確率が最も小さい（すなわち、絶対に正解ではないパラ言語情報の確率が最も小さい）クラスを正解のパラ言語情報推定結果とみなす。逆教師の推定モデルと従来教師の推定モデルの両方を用いる場合であれば、従来教師の推定モデルの出力確率から、逆教師の推定モデルの出力確率を減算した値（すなわち、正解であるパラ言語情報の確率から正解でないパラ言語情報の確率を減算した値）が最も大きいクラスを正解のパラ言語情報推定結果とみなす。

【0022】

[第一実施形態]

第一実施形態では、逆教師の推定モデルを単独で用いてパラ言語情報を推定する。

【0023】

<パラ言語情報推定モデル学習装置1>

第一実施形態のパラ言語情報推定モデル学習装置は、複数の発話と、発話ごとに複数の聴取者が付与したパラ言語情報の判定結果を表すパラ言語情報ラベルとからなる学習データから、逆教師推定モデルを学習する。第一実施形態のパラ言語情報推定モデル学習装置1は、図2に例示するように、音響特徴抽出部11、逆教師決定部12、逆教師推定モデル学習部13、および逆教師推定モデル記憶部10を備える。このパラ言語情報推定モデル学習装置1が、図3に例示する各ステップの処理を行うことにより第一実施形態のパラ言語情報推定モデル学習

7

方法が実現される。

【0024】

パラ言語情報推定モデル学習装置1は、例えば、中央演算処理装置(CPU: Central Processing Unit)、主記憶装置(RAM: Random Access Memory)などを有する公知又は専用のコンピュータに特別なプログラムが読み込まれて構成された特別な装置である。パラ言語情報推定モデル学習装置1は、例えば、中央演算処理装置の制御のもとで各処理を実行する。パラ言語情報推定モデル学習装置1に入力されたデータや各処理で得られたデータは、例えば、主記憶装置に格納され、主記憶装置に格納されたデータは必要に応じて中央演算処理装置へ読み出されて他の処理に利用される。パラ言語情報推定モデル学習装置1の各処理部は、少なくとも一部が集積回路等のハードウェアによって構成されていてもよい。パラ言語情報推定モデル学習装置1が備える各記憶部は、例えば、RAM(Random Access Memory)などの主記憶装置、ハードディスクや光ディスクもしくはフラッシュメモリ(Flash Memory)のような半導体メモリ素子により構成される補助記憶装置、またはリレーショナルデータベースやキーバリューストアなどのミドルウェアにより構成することができる。

【0025】

ステップS11において、音響特徴抽出部11は、学習データの発話から韻律特徴を抽出する。韻律特徴は、基本周波数、短時間パワー、メル周波数ケプストラム係数(Mel-frequency Cepstral Coefficients: MFCC)、ゼロ交差率、Harmonics-to-Noise-Ratio(HNR)、メルフィルタバンク出力、のいずれか一つ以上の特徴量を含むベクトルである。また、これらの時間ごと(フレームごと)の系列ベクトルであってもよいし、これらの一定時間ごとまたは発話全体の統計量(平均、分散、最大値、最小値、勾配など)のベクトルであってもよい。音響特徴抽出部11は、抽出した韻律特徴を逆教師推定モデル学習部13へ出力する。

【0026】

ステップS12において、逆教師決定部12は、学習データのパラ言語情報ラベルから逆教師ラベルを決定する。逆教師とは、推定対象となるパラ言語情報のうち、予め定めた閾値(以下、「逆教師閾値」と呼ぶ)以下(例えば、1割以下)の聴取者が判定したパラ言語情報を指すものとする。逆教師ラベルは、逆教師のパラ言語情報クラスが1、それ以外は0となるベクトルを指す。すなわち、従来教師のようにいずれか一つのパラ言語情報クラスが1、それ以外が0となるベクトルではなく、少なくとも一つ以上のパラ言語情報クラスが1となるベクトルが逆教師ラベルである。例えば、喜び・悲しみ・怒り・平静の4クラスの感情推定において、逆教師閾値を0.1とし、3名の聴取者がある音声に対して「喜び」「喜び」「平静」と判定した場合、その音声の逆教師ラベル

8

は「悲しみ」「怒り」の2クラスが1、「喜び」「平静」の2クラスが0となる4次元のベクトルを指す。

【0027】

逆教師ラベルは、具体的には以下のように表される。

【0028】

【数2】

$$t^* = \begin{bmatrix} t_1^* \\ \vdots \\ t_K^* \end{bmatrix},$$

$$t_k^* = \begin{cases} 1 & \text{if } \frac{1}{N} \sum_{n=1}^N h_k^n \leq \beta \\ 0 & \text{otherwise} \end{cases}$$

【0029】

ここで、 $h_k^n$ はn番目の聴取者がk番目のパラ言語情報クラスを感じたか否か(感じた場合は1、感じなかった場合は0)を表す。Kはパラ言語情報クラスの総数である。Nは聴取者の総数である。βは0以上1以下の逆教師閾値である。

【0030】

逆教師ラベルは、いずれの聴取者も判定しなかったパラ言語情報クラスを逆教師としてもよい。これは、逆教師閾値を0に設定した場合に相当する。

【0031】

逆教師決定部12は、決定した逆教師ラベルを逆教師推定モデル学習部13へ出力する。

【0032】

ステップS13において、逆教師推定モデル学習部13は、音響特徴抽出部11が出力する韻律特徴と逆教師決定部12が出力する逆教師ラベルとに基づいて、逆教師推定モデルを学習する。推定モデルにはマルチラベル分類問題(一つの音声複数のクラスに同時に分類される分類問題)を扱うことができるモデルを用いるものとする。これは一つの音声に対して複数のクラスに逆教師が表れることがあるためである。推定モデルは例えば従来技術のような深層学習に基づくモデルであってもよいし、多クラスロジスティック回帰であってもよいが、出力が確率値(あるパラ言語情報クラスが1である確率)として表現可能であるものとする。逆教師推定モデル学習部13は、学習した逆教師推定モデルを逆教師推定モデル記憶部20へ記憶する。

【0033】

<パラ言語情報推定装置2>

第一実施形態のパラ言語情報推定装置は、学習済みの逆教師推定モデルを用いて、入力された発話のパラ言語情報を推定する。第一実施形態のパラ言語情報推定装置2は、図4に例示するように、音響特徴抽出部11、逆教

師推定モデル記憶部 2 0、およびパラ言語情報推定部 2 1 を備える。このパラ言語情報推定装置 2 が、図 5 に例示する各ステップの処理を行うことにより第一実施形態のパラ言語情報推定方法が実現される。

#### 【 0 0 3 4 】

パラ言語情報推定装置 2 は、例えば、中央演算処理装置 (CPU: Central Processing Unit)、主記憶装置 (RAM: Random Access Memory) などを有する公知又は専用のコンピュータに特別なプログラムが読み込まれて構成された特別な装置である。パラ言語情報推定装置 2 は、例えば、中央演算処理装置の制御のもとで各処理を実行する。パラ言語情報推定装置 2 に入力されたデータや各処理で得られたデータは、例えば、主記憶装置に格納され、主記憶装置に格納されたデータは必要に応じて中央演算処理装置へ読み出されて他の処理に利用される。パラ言語情報推定装置 2 の各処理部は、少なくとも一部が集積回路等のハードウェアによって構成されている。パラ言語情報推定装置 2 が備える各記憶部は、例えば、RAM (Random Access Memory) などの主記憶装置、ハードディスクや光ディスクもしくはフラッシュメモリ (Flash Memory) のような半導体メモリ素子により構成される補助記憶装置、またはリレーショナルデータベースやキーバリューストアなどのミドルウェアにより構成することができる。

#### 【 0 0 3 5 】

逆教師推定モデル記憶部 2 0 には、パラ言語情報推定モデル学習装置 1 が学習した逆教師推定モデルが記憶されている。

#### 【 0 0 3 6 】

ステップ S 1 1 において、音響特徴抽出部 1 1 は、入力された発話から韻律特徴を抽出する。韻律特徴の抽出は、パラ言語情報推定モデル学習装置 1 と同様に行えばよい。音響特徴抽出部 1 1 は、抽出した韻律特徴をパラ言語情報推定部 2 1 へ出力する。

#### 【 0 0 3 7 】

ステップ S 2 1 において、パラ言語情報推定部 2 1 は、逆教師推定モデル記憶部 2 0 に記憶されている逆教師推定モデルに基づいて、音響特徴抽出部 1 1 が出力した韻律特徴からパラ言語情報を推定する。推定においては、ある韻律特徴に対し、逆教師推定モデルの出力が最も低いクラスをパラ言語情報推定結果とみなす。これは、逆教師の可能性が最も低いパラ言語情報、すなわち「絶対に正解ではないパラ言語情報」ではないと考えられるパラ言語情報を選択することに相当する。パラ言語情報推定部 2 1 は、パラ言語情報の推定結果をパラ言語情報推定装置 2 の出力とする。

#### 【 0 0 3 8 】

##### [ 第二実施形態 ]

第二実施形態では、逆教師の推定モデルに加えて、従来教師の推定モデルを用いてパラ言語情報を推定する。こ

のとき、各推定モデルの出力結果の重み付け差によりパラ言語情報の推定を行う。「あるパラ言語情報が正解である確率」と「あるパラ言語情報が正解でない確率」の両方を考慮してパラ言語情報推定を行うことに相当する。その結果、どちらかの確率のみを考慮する場合 (すなわち、従来技術や第一実施形態のそれぞれ) に比べてパラ言語情報の推定精度が向上する。

#### 【 0 0 3 9 】

< パラ言語情報推定モデル学習装置 3 >

第二実施形態のパラ言語情報推定モデル学習装置は、第一実施形態と同様の学習データから逆教師推定モデルと従来教師推定モデルを学習する。第二実施形態のパラ言語情報推定モデル学習装置 3 は、図 6 に例示するように、第一実施形態の音響特徴抽出部 1 1、逆教師決定部 1 2、逆教師推定モデル学習部 1 3、および逆教師推定モデル記憶部 2 0 に加えて、従来教師決定部 3 1、従来教師推定モデル学習部 3 2、および従来教師推定モデル記憶部 4 0 をさらに備える。このパラ言語情報推定モデル学習装置 3 が、図 7 に例示する各ステップの処理を行うことにより第二実施形態のパラ言語情報推定モデル学習方法が実現される。

#### 【 0 0 4 0 】

以下、第一実施形態のパラ言語情報推定モデル学習装置 1 との相違点を中心に、第二実施形態のパラ言語情報推定モデル学習装置 3 を説明する。

#### 【 0 0 4 1 】

ステップ S 3 1 において、従来教師決定部 3 1 は、学習データのパラ言語情報ラベルから従来教師ラベルを決定する。従来教師ラベルは、従来技術と同様に、ある音声に対して全聴取者のうち過半数が同じパラ言語情報を判定した場合に、そのパラ言語情報クラスを 1、それ以外は 0 としたベクトルである。過半数が同じパラ言語情報を判定しなかった場合、その音声は正解なしとしてモデル学習には利用しない。例えば、喜び・悲しみ・怒り・平静の 4 クラスの感情推定において、ある音声に対して 3 名の聴取者が「喜び」「喜び」「平静」と判定した場合、その音声の従来教師ラベルは「喜び」クラスが 1、残りの「悲しみ」「怒り」「平静」の 3 クラスが 0 となる 4 次元のベクトルを指す。従来教師決定部 3 1 は、決定した従来教師ラベルを従来教師推定モデル学習部 3 2 へ出力する。

#### 【 0 0 4 2 】

ステップ S 3 2 において、従来教師推定モデル学習部 3 2 は、音響特徴抽出部 1 1 が出力する韻律特徴と従来教師決定部 3 1 が出力する従来教師ラベルとに基づいて、従来教師推定モデルを学習する。推定モデルには多クラス分類問題 (一つの音声から一つのクラスを分類する分類問題) を扱うことができるモデルを用いるものとする。推定モデルは例えば従来技術のような深層学習に基づくモデルであってもよいし、多クラスロジスティック回

帰であってもよいが、出力が確率値（あるパラ言語情報クラスが1である確率）として表現可能であるものとする。従来教師推定モデル学習部32は、学習した従来教師推定モデルを従来教師推定モデル記憶部40へ記憶する。

## 【0043】

<パラ言語情報推定装置4>

第二実施形態のパラ言語情報推定装置は、学習済みの逆教師推定モデルと従来教師推定モデルの両方を用いて、入力された発話のパラ言語情報を推定する。第二実施形態のパラ言語情報推定装置4は、図8に例示するように、第一実施形態の音響特徴抽出部11および逆教師推定モデル記憶部20に加えて、従来教師推定モデル記憶部40およびパラ言語情報推定部41をさらに備える。このパラ言語情報推定装置4が、図9に例示する各ステップの処理を行うことにより第二実施形態のパラ言語情報推定方法が実現される。

## 【0044】

以下、第一実施形態のパラ言語情報推定装置2との相違点を中心に、第二実施形態のパラ言語情報推定装置4を説明する。

## 【0045】

従来教師推定モデル記憶部40には、パラ言語情報推定モデル学習装置3が学習した従来教師推定モデルが記憶されている。

## 【0046】

ステップS41において、パラ言語情報推定部41は、逆教師推定モデル記憶部20に記憶されている逆教師推定モデルと従来教師推定モデル記憶部40に記憶されている従来教師推定モデルの両方に基づいて、音響特徴抽出部11が出力した韻律特徴からパラ言語情報を推定する。推定においては、ある韻律特徴に対する従来教師推定モデルの出力と逆教師推定モデルの出力との重み付け差により決定する。これは、「あるパラ言語情報が正解である確率」と「あるパラ言語情報が正解でない確率」の両方を考慮してパラ言語情報推定を行うことに相当する。

## 【0047】

パラ言語情報の推定は、具体的には以下のように表される。

## 【0048】

【数3】

$$\hat{c}_k = \arg \max_{c_k} ((1 - \alpha)p(c_k) - \alpha q(c_k))$$

## 【0049】

ここで、 $\hat{c}_k$  はパラ言語情報の推定結果を表す。 $c_k$  はk番目のパラ言語情報クラスを表す。

$p(c_k)$  はk番目のパラ言語情報クラスが正解である確率を表し、従来教師推定モデルの出力である。 $q(c_k)$  はk番目

のパラ言語情報クラスが正解でない確率を表し、逆教師推定モデルの出力である。 $\alpha$  は推定重みを表す。

## 【0050】

推定重み  $\alpha$  は0から1までの連続値のいずれかの値とする。推定重みが0に近いほど「あるパラ言語情報が正解である確率」を重視し、1に近いほど「あるパラ言語情報が正解でない確率」を重視した推定を行うこととなる。例えば、推定重みは0.3とする。

## 【0051】

[第二実施形態の変形例]

第二実施形態では、逆教師推定モデル学習部と従来教師推定モデル学習部を完全に別々の処理として構成する例を説明したが、一方で学習した推定モデルを他方の推定モデル初期値として利用しながら各推定モデルの学習を行うことも可能である。すなわち、従来教師推定モデル学習部32に逆教師推定モデル記憶部20で記憶されている逆教師推定モデルを入力し、従来教師推定モデル学習部32では逆教師推定モデルを従来教師推定モデルの初期値として、音響特徴抽出部11が出力する韻律特徴と従来教師決定部31が出力する従来教師ラベルとに基づいて、従来教師推定モデルを学習する。または、逆教師推定モデル学習部13に従来教師推定モデル記憶部40で記憶されている従来教師推定モデルを入力し、逆教師推定モデル学習部13では従来教師推定モデルを逆教師推定モデルの初期値として、音響特徴抽出部11が出力する韻律特徴と逆教師決定部12が出力する逆教師ラベルとに基づいて、逆教師推定モデルを学習する。従来教師推定モデルと逆教師推定モデルはそれぞれ韻律特徴と従来教師または韻律特徴と逆教師との関連性を学習しており、一方で学習された推定基準は他方でも利用可能であると考えられるため、これらの処理を加えることでパラ言語情報推定精度がさらに向上する可能性がある。

## 【0052】

[第三実施形態]

第三実施形態では、従来教師の推定と逆教師の推定とを同時に行うマルチタスク推定モデルを用いてパラ言語情報を推定する。このとき、モデル学習において、従来教師の推定と逆教師の推定をマルチタスク学習として同時に学習する。マルチタスク学習は、異なる課題を単一のモデルで解くことにより、課題間で共通する知識を獲得でき、各課題の推定精度が向上することが知られている（下記参考文献1参照）。第三実施形態は、第二実施形態と同様に従来教師と逆教師の両方を用いたパラ言語情報推定であるが、マルチタスク学習として学習することで推定モデル自体を改善することができるため、推定精度がより向上する。

【参考文献1】R. Caruana, "Multitask Learning", Machine Learning, vol. 28, pp.41-75, 1997.

## 【0053】

<パラ言語情報推定モデル学習装置5>

第三実施形態のパラ言語情報推定モデル学習装置は、第一実施形態と同様の学習データからマルチタスク推定モデルを学習する。第三実施形態のパラ言語情報推定モデル学習装置 5 は、図 10 に例示するように、第一実施形態の音響特徴抽出部 11 および逆教師決定部 12 と、第二実施形態の従来教師決定部 31 とに加えて、マルチタスク推定モデル学習部 51 およびマルチタスク推定モデル記憶部 60 をさらに備える。このパラ言語情報推定モデル学習装置 5 が、図 11 に例示する各ステップの処理を行うことにより第三実施形態のパラ言語情報推定モデル学習方法が実現される。

#### 【0054】

以下、第一実施形態のパラ言語情報推定モデル学習装置 1 および第二実施形態のパラ言語情報推定モデル学習装置 3 との相違点を中心に、第三実施形態のパラ言語情報推定モデル学習装置 5 を説明する。

#### 【0055】

ステップ S51 において、マルチタスク推定モデル学習部 51 は、音響特徴抽出部 11 が出力する韻律特徴と逆教師決定部 12 が出力する逆教師ラベルと従来教師決定部 31 が出力する従来教師ラベルとを用いてマルチタスク学習を行い、マルチタスク推定モデルを学習する。マルチタスク学習では、一般的にニューラルネットワークに基づく推定モデルが用いられるため、本実施形態での推定モデルはニューラルネットワークに基づく推定モデルとする。例えば、図 10 に示すように、従来技術の深層学習に基づく推定モデルに対して、逆教師の推定を行う分岐構造を加えた推定モデルとする。マルチタスク推定モデル学習部 51 は、学習したマルチタスク推定モデルをマルチタスク推定モデル記憶部 60 へ記憶する。

#### 【0056】

<パラ言語情報推定装置 6 >

第三実施形態のパラ言語情報推定装置は、学習済みのマルチタスク推定モデルを用いて、入力された発話のパラ言語情報を推定する。第三実施形態のパラ言語情報推定装置 6 は、図 12 に例示するように、第一実施形態の音響特徴抽出部 11 に加えて、マルチタスク推定モデル記憶部 60 およびパラ言語情報推定部 61 をさらに備える。このパラ言語情報推定装置 6 が、図 13 に例示する各ステップの処理を行うことにより第三実施形態のパラ言語情報推定方法が実現される。

#### 【0057】

以下、第一実施形態のパラ言語情報推定装置 2 および第二実施形態のパラ言語情報推定装置 4 との相違点を中心に、第三実施形態のパラ言語情報推定装置 6 を説明する。

#### 【0058】

マルチタスク推定モデル記憶部 60 には、パラ言語情報推定モデル学習装置 5 が学習したマルチタスク推定モデルが記憶されている。

#### 【0059】

ステップ S61 において、パラ言語情報推定部 61 は、マルチタスク推定モデル記憶部 60 に記憶されているマルチタスク推定モデルに基づいて、音響特徴抽出部 11 が出力した韻律特徴からパラ言語情報を推定する。推定においては、ある韻律特徴に対し、従来教師の推定出力が最も高いクラスをパラ言語情報推定結果とみなす。推定モデルの学習においてマルチタスク学習を用いているため、逆教師の影響を考慮した（すなわち、逆教師を間違えないようにしつつ従来教師を推定している）パラ言語情報推定を行うことができ、パラ言語情報推定精度が向上する。

#### 【0060】

[変形例]

上述の実施形態では、パラ言語情報推定モデル学習装置とパラ言語情報推定装置とを別々の装置として構成する例を説明したが、パラ言語情報推定モデル学習する機能と学習済みのパラ言語情報推定モデルを用いてパラ言語情報を推定する機能とを兼ね備えた 1 台のパラ言語情報推定装置として構成することも可能である。すなわち、第一実施形態の変形例のパラ言語情報推定装置は、音響特徴抽出部 11、逆教師決定部 12、逆教師推定モデル学習部 13、逆教師推定モデル記憶部 20、およびパラ言語情報推定部 21 を備える。

また、第二実施形態の変形例のパラ言語情報推定装置は、音響特徴抽出部 11、逆教師決定部 12、逆教師推定モデル学習部 13、および逆教師推定モデル記憶部 20 に加えて、従来教師決定部 31、従来教師推定モデル学習部 32、従来教師推定モデル記憶部 40、およびパラ言語情報推定部 41 をさらに備える。そして、第三実施形態の変形例のパラ言語情報推定装置は、音響特徴抽出部 11、逆教師決定部 12、および従来教師決定部 31 に加えて、マルチタスク推定モデル学習部 51、マルチタスク推定モデル記憶部 60、およびパラ言語情報推定部 61 をさらに備える。

#### 【0061】

以上、本発明の実施の形態について説明したが、具体的な構成は、これらの実施の形態に限られるものではなく、本発明の趣旨を逸脱しない範囲で適宜設計の変更等があっても、本発明に含まれることはいうまでもない。実施の形態において説明した各種の処理は、記載の順に従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。

#### 【0062】

[プログラム、記録媒体]

上記実施形態で説明した各装置における各種の処理機能をコンピュータによって実現する場合、各装置が有すべき機能の処理内容はプログラムによって記述される。そして、このプログラムをコンピュータで実行することに



15

より、上記各装置における各種の処理機能がコンピュータ上で実現される。

## 【0063】

この処理内容を記述したプログラムは、コンピュータで読み取り可能な記録媒体に記録しておくことができる。コンピュータで読み取り可能な記録媒体としては、例えば、磁気記録装置、光ディスク、光磁気記録媒体、半導体メモリ等どのようなものでもよい。

## 【0064】

また、このプログラムの流通は、例えば、そのプログラムを記録したDVD、CD-ROM等の可搬型記録媒体を販売、譲渡、貸与等することによって行ふ。さらに、このプログラムをサーバコンピュータの記憶装置に格納しておき、ネットワークを介して、サーバコンピュータから他のコンピュータにそのプログラムを転送することにより、このプログラムを流通させる構成としてもよい。

## 【0065】

このようなプログラムを実行するコンピュータは、例えば、まず、可搬型記録媒体に記録されたプログラムもしくはサーバコンピュータから転送されたプログラムを、一旦、自己の記憶装置に格納する。そして、処理の実行時、このコンピュータは、自己の記憶装置に格納されたプログラムを読み取り、読み取ったプログラムに従った処理を実行する。また、このプログラムの別の実行形態として、コンピュータが可搬型記録媒体から直接プログラムを読み取り、そのプログラムに従った処理を実行することとしてもよく、さらに、このコンピュータにサーバコンピュータからプログラムが転送されるたびに、逐次、受け取ったプログラムに従った処理を実行すること

16

としてもよい。また、サーバコンピュータから、このコンピュータへのプログラムの転送は行わず、その実行指示と結果取得のみによって処理機能を実現する、いわゆるASP (Application Service Provider) 型のサービスによって、上述の処理を実行する構成としてもよい。なお、本形態におけるプログラムには、電子計算機による処理の用に供する情報であってプログラムに準ずるもの(コンピュータに対する直接の指令ではないがコンピュータの処理を規定する性質を有するデータ等)を含むものとする。

## 【0066】

また、この形態では、コンピュータ上で所定のプログラムを実行させることにより、本装置を構成することとしたが、これらの処理内容の少なくとも一部をハードウェア的に実現することとしてもよい。

## 【符号の説明】

## 【0067】

- 1, 3, 5 パラ言語情報推定モデル学習装置
- 11 音響特徴抽出部
- 12 逆教師決定部
- 13 逆教師推定モデル学習部
- 20 逆教師推定モデル記憶部
- 31 従来教師決定部
- 32 従来教師推定モデル学習部
- 40 従来教師推定モデル記憶部
- 51 マルチタスク推定モデル学習部
- 60 マルチタスク推定モデル記憶部
- 2, 4, 6 パラ言語情報推定装置
- 21, 41, 61 パラ言語情報推定部

【図1】

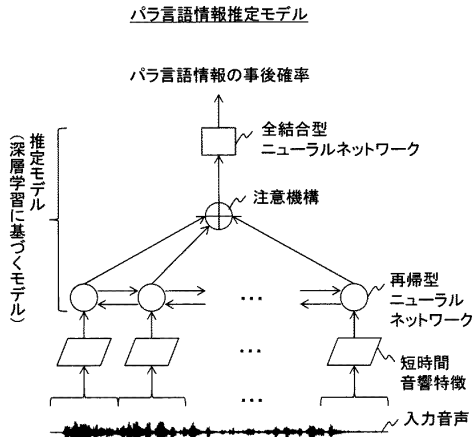


図1

【図2】

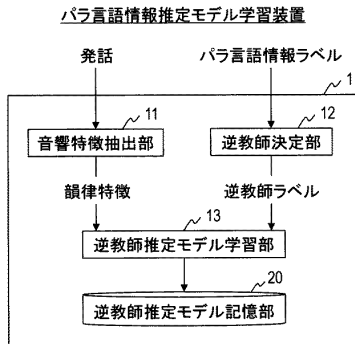


図2

【図3】

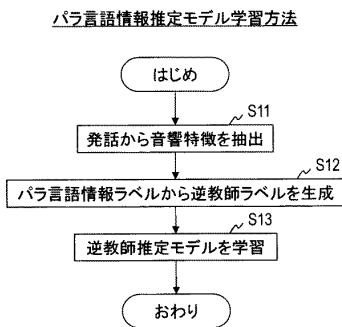


図3

【図4】

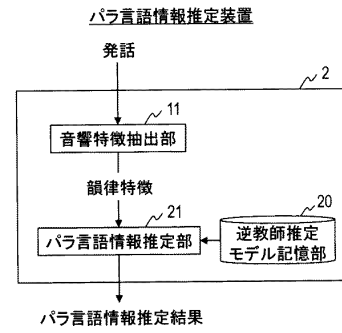


図4

【図5】

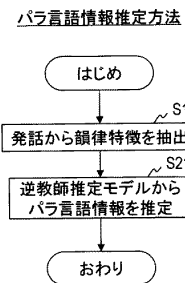


図5

【図6】

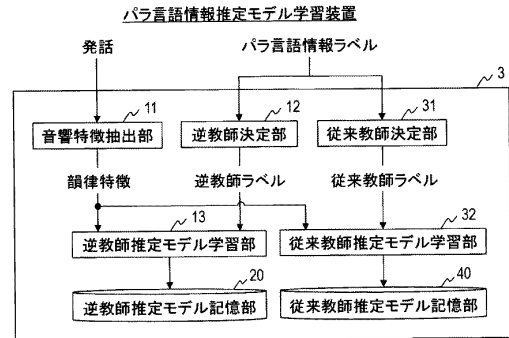


図6

【図7】

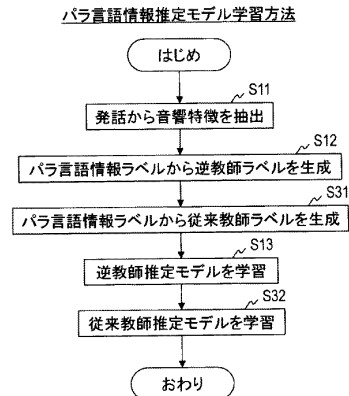


図7

【図 8】

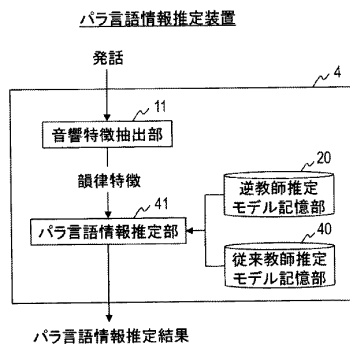


図 8

【図 9】

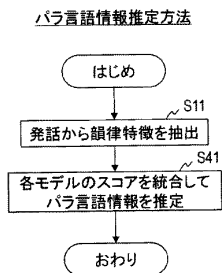


図 9

【図 10】

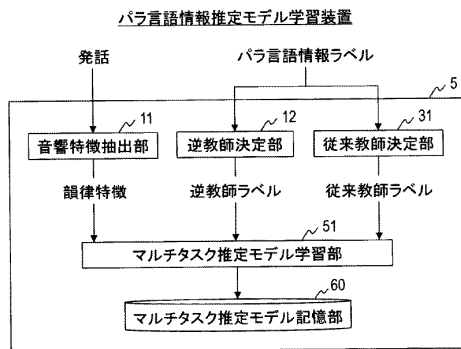


図 10

【図 11】

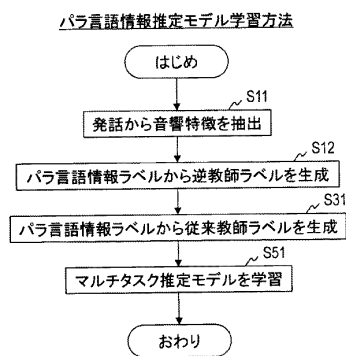


図 11

【図 12】

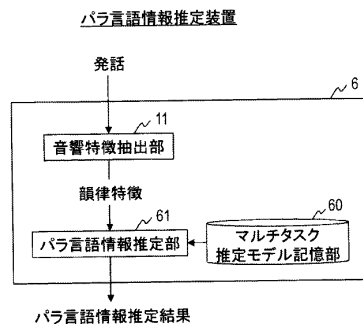


図 12

【図 13】

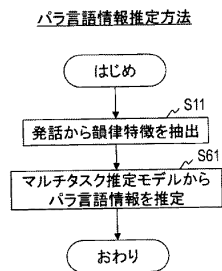


図 13

【図 14】

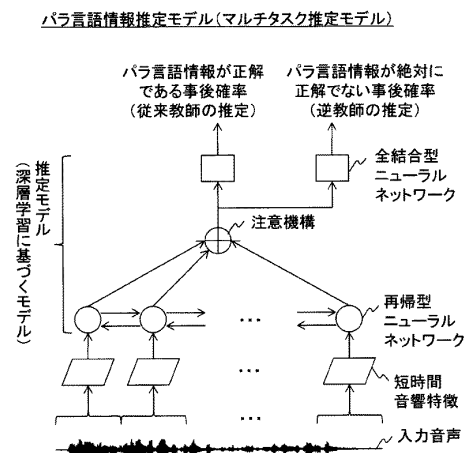


図 14

(72)発明者 小橋川 哲  
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