

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号

特許第6468519号
(P6468519)

(45) 発行日 平成31年2月13日 (2019. 2. 13)

(24) 登録日 平成31年1月25日 (2019. 1. 25)

(51) Int. Cl.		F I			
G 1 O L 25/90	(2013. 01)	G 1 O L	25/90		
G 1 O L 13/10	(2013. 01)	G 1 O L	13/10	1 1 1 C	
		G 1 O L	13/10	1 1 3 Z	

請求項の数 8 (全 39 頁)

(21) 出願番号	特願2016-32412 (P2016-32412)	(73) 特許権者	000004226
(22) 出願日	平成28年2月23日 (2016. 2. 23)		日本電信電話株式会社
(65) 公開番号	特開2017-151224 (P2017-151224A)		東京都千代田区大手町一丁目5番1号
(43) 公開日	平成29年8月31日 (2017. 8. 31)	(73) 特許権者	504143441
審査請求日	平成29年12月8日 (2017. 12. 8)		国立大学法人 奈良先端科学技術大学院大学
			奈良県生駒市高山町8916-5
		(74) 代理人	110001519
			特許業務法人太陽国際特許事務所
		(72) 発明者	亀岡 弘和
			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	田中 宏
			奈良県生駒市高山町8916-5 国立大
			学法人奈良先端科学技術大学院大学内
			最終頁に続く

(54) 【発明の名称】 基本周波数パターン予測装置、方法、及びプログラム

(57) 【特許請求の範囲】

【請求項1】

学習サンプルのソース音声の時系列データとターゲット音声の時系列データとからなるパラレルデータを入力として、前記ソース音声の時系列データから抽出される各時刻のスペクトル特徴量ベクトルと、前記ターゲット音声の時系列データから抽出される、各時刻の基本周波数とに基づいて、前記ソース音声の各時刻のスペクトル特徴量ベクトルと、前記ターゲット音声の各時刻の基本周波数との間の関係をモデル化した第1確率分布のパラメータを学習する第1モデルパラメータ学習部と、

前記ターゲット音声の各時刻の基本周波数に基づいて、基本周波数パターン生成過程をモデル化した第2確率分布のパラメータを学習する第2モデルパラメータ学習部と、

予測対象のソース音声の時系列データを入力として、前記予測対象のソース音声の時系列データから抽出される各時刻のスペクトル特徴量ベクトルと、前記第1モデルパラメータ学習部によって学習された前記第1確率分布のパラメータと、前記第2モデルパラメータ学習部によって学習された前記第2確率分布のパラメータとに基づいて、前記第1確率分布と前記第2確率分布とを用いて表される規準を大きくするように、前記予測対象のソース音声に対応する前記ターゲット音声の各時刻の基本周波数を予測する基本周波数予測部と、

を含む基本周波数パターン予測装置。

【請求項2】

前記第1確率分布を、前記ソース音声の各時刻のスペクトル特徴量ベクトルと前記ター

ゲット音声の各時刻の基本周波数と前記基本周波数の動的成分との同時確率分布を表す混合正規分布とし、

前記第 2 確率分布を、

各時刻の基本周波数と、

隠れマルコフモデルの各時刻の状態からなる状態系列、又は各時刻における甲状軟骨の平行移動運動によって生じる基本周波数パターンを表すフレーズ指令及び甲状軟骨の回転運動によって生じる基本周波数パターンを表すアクセント指令のペアからなる指令関数との組み合わせの確率分布とした請求項 1 記載の基本周波数パターン予測装置。

【請求項 3】

前記規準を、

前記第 1 確率分布と前記第 2 確率分布との積を用いて表される、各時刻の基本周波数と、前記状態系列との組み合わせに応じた関数、または

前記第 1 確率分布と前記第 2 確率分布との積を用いて表される、各時刻の基本周波数と、各時刻の前記指令関数との組み合わせに応じた関数とした請求項 2 記載の基本周波数パターン予測装置。

【請求項 4】

前記第 1 モデルパラメータ学習部は、E M (Expectation-Maximization) アルゴリズムにより、前記第 1 確率分布から求められる前記ソース音声の各時刻のスペクトル特徴量ベクトル及び前記ターゲット音声の各時刻の基本周波数の尤もらしさが大きくなるように、前記第 1 確率分布のパラメータを学習する請求項 1 ~ 請求項 3 の何れか 1 項記載の基本周波数パターン予測装置。

【請求項 5】

前記第 2 モデルパラメータ学習部は、E M (Expectation-Maximization) アルゴリズムにより、前記第 2 確率分布から求められる、各時刻の基本周波数と、隠れマルコフモデルの各時刻の状態からなる状態系列との尤もらしさが大きくなるように、前記第 2 確率分布のパラメータとして、前記状態系列における状態遷移確率、及び各時刻における状態に応じたフレーズ指令の振幅及び各アクセント指令の振幅を表すパラメータ群を学習するか、又は

各時刻の基本周波数が与えられたときの、各時刻のフレーズ指令及びアクセント指令のペアからなる指令関数及び前記パラメータ群の対数事後確率を目的関数として、前記目的関数を増加させるように、前記指令関数及び前記パラメータ群を、前記第 2 確率分布のパラメータとして学習する請求項 1 ~ 請求項 4 の何れか 1 項記載の基本周波数パターン予測装置。

【請求項 6】

前記基本周波数予測部は、

前記第 1 確率分布と前記第 2 確率分布との積を用いて表される、各時刻の基本周波数と、隠れマルコフモデルの各時刻の状態からなる状態系列との組み合わせに応じた関数が大きくなるように、各時刻の基本周波数と前記状態系列とを推定することにより、前記予測対象のソース音声に対応する前記ターゲット音声の各時刻の基本周波数を予測するか、又は

前記第 1 確率分布と前記第 2 確率分布との積を用いて表される、各時刻の基本周波数、各時刻のフレーズ指令及びアクセント指令のペアからなる指令関数との組み合わせに応じた関数が大きくなるように、各時刻の基本周波数と、各時刻の前記指令関数を推定することにより、前記予測対象のソース音声に対応する前記ターゲット音声の各時刻の基本周波数を予測する請求項 1 ~ 請求項 5 の何れか 1 項記載の基本周波数パターン予測装置。

【請求項 7】

第 1 モデルパラメータ学習部と、第 2 モデルパラメータ学習部と、基本周波数予測部とを含む基本周波数パターン予測装置における基本周波数パターン予測方法であって、

前記第 1 モデルパラメータ学習部が、学習サンプルのソース音声の時系列データとターゲット音声の時系列データとからなるパラレルデータを入力として、前記ソース音声の時

10

20

30

40

50

系列データから抽出される各時刻のスペクトル特徴量ベクトルと、前記ターゲット音声の時系列データから抽出される、各時刻の基本周波数とに基づいて、前記ソース音声の各時刻のスペクトル特徴量ベクトルと、前記ターゲット音声の各時刻の基本周波数との間の関係をモデル化した第1確率分布のパラメータを学習し、

前記第2モデルパラメータ学習部が、前記ターゲット音声の各時刻の基本周波数に基づいて、基本周波数パターン生成過程をモデル化した第2確率分布のパラメータを学習し、

前記基本周波数予測部が、予測対象のソース音声の時系列データを入力として、前記予測対象のソース音声の時系列データから抽出される各時刻のスペクトル特徴量ベクトルと、前記第1モデルパラメータ学習部によって学習された前記第1確率分布のパラメータと、前記第2モデルパラメータ学習部によって学習された前記第2確率分布のパラメータとに基づいて、前記第1確率分布と前記第2確率分布とを用いて表される規準を大きくするように、前記予測対象のソース音声に対応する前記ターゲット音声の各時刻の基本周波数を予測する

基本周波数パターン予測方法。

【請求項8】

請求項1～請求項6の何れか1項に記載の基本周波数パターン予測装置の各部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、基本周波数パターン予測装置、方法、及びプログラムに係り、特に、ソース音声から、ターゲット音声の基本周波数パターンを予測する基本周波数パターン予測装置、方法、及びプログラムに関する。

【背景技術】

【0002】

他者とのコミュニケーションにおいて音声は利便性に優れた手段ではあるが、時として物理的制約により様々な障壁が必然的にもたらされる。例えば、発声器官の内、わずかに所でも正常に動作しなくなると、深刻な発声障害を患い、音声コミュニケーションに支障をきたす。また、音声生成という物理的行為は、秘匿性の高い意思伝達には不向きであるし、周囲の騒音に脆弱である。これらの障壁を無くすためには、身体的制約を超えて発声器官を動作させて音声を生成したり、適切な発音動作を指定して音声を生成したり、聴取困難なほど微かな音声発声時の発声器官動作から通常音声を生成するなど、物理的・身体的制約を超えた音声生成機能の拡張が必要である。

【0003】

例えば、喉頭癌などで喉頭を失った喉頭摘出者に対して、残存器官を用いた代替発声法により生成される自然性に乏しい音声を、より自然な音声へと変換する発声補助技術が提案されている（非特許文献1～非特許文献3を参照）。この他にも、非可聴つぶやき音声を自然な音声に変換する技術も提案されており、秘匿性に優れた通話技術としての応用が期待されている。上述の技術はいずれも音声のスペクトル特徴量系列から自然音声の基本周波数(F_0)パターンを予測する問題を扱っている点で共通しており、学習処理と変換処理で構成される。学習処理では、対象音声（前者であれば電気音声、後者であれば非可聴つぶやき音声）と通常音声の同一発話データを用いる。まず各離散時刻（以後、フレーム）において、前後数フレームから得られる対象音声のスペクトル特徴量と、通常音声の対数 F_0 とその動的成分（時間微分または時間差分）を抽出し、スペクトル距離尺度に基づく動的時間伸縮によりこれらに対応付けた結合ベクトルを得る。これをパラレルデータと呼ぶ。各フレームのパラレルデータを用い、対象音声のスペクトル特徴量と通常音声の対数 F_0 の静的・動的成分の結合確率密度関数を混合正規分布モデル（Gaussian Mixture Model; GMM）で表現する。GMMのパラメータはExpectation-Maximizationアルゴリズムにより学習することができる。変換処理では、学習されたGMMを用いて、系列内変動を考慮した最尤系列変換法により、対象音声のスペクトル特徴量系列から通常音声の F_0 パターン

10

20

30

40

50

へと変換することができる。

【先行技術文献】

【非特許文献】

【0004】

【非特許文献1】Keigo Nakamura, Tomoki Toda, Hiroshi Saruwatari, Kiyohiro Shikano, "Speaking-aid systems using GMM-based voice conversion for electrolaryngeal speech," Speech Communication, vol. 54, no. 1, pp. 134-146, 2012.

【非特許文献2】Kou Tanaka, Tomoki Toda, Graham Neubig, Sakriani Sakti, Satoshi Nakamura, "A hybrid approach to electrolaryngeal speech enhancement based on noise reduction and statistical excitation generation," IEICE Transactions on Information and Systems, vol. E97-D, no. 6, pp. 1429-1437, Jun. 2014. 10

【非特許文献3】Kou Tanaka, Tomoki Toda, Graham Neubig, Sakriani Sakti, Satoshi Nakamura, "Direct F0 control of an electrolarynx based on statistical excitation feature prediction and its evaluation through simulation," Proc. INTERSPEECH, pp. 31-35, Sep. 2014.

【非特許文献4】Hirokazu Kameoka, Jonathan Le Roux, Yasunori Ohishi, "A statistical model of speech F0 contours," ISCA Tutorial and Research Workshop on Statistical And Perceptual Audition (SAPA 2010), pp. 43-48, Sep. 2010.

【非特許文献5】Kota Yoshizato, Hirokazu Kameoka, Daisuke Saito, Shigeki Sagayama, "Hidden Markov convolutive mixture model for pitch contour analysis of speech," in Proc. The 13th Annual Conference of the International Speech Communication Association (Interspeech 2012), Sep. 2012. 20

【発明の概要】

【発明が解決しようとする課題】

【0005】

従来技術では、学習処理や変換処理において音声の F_0 パターンの物理的な生成過程を考慮したモデルが用いられていなかったため、物理的に人間が発声しえないような不自然な F_0 パターンを生成することが起こりえた。この問題に対し、 F_0 パターンの物理的な生成過程を考慮した予測を行うことで、より自然な F_0 パターンを生成できる可能性がある。

【0006】

F_0 パターンは声帯に張力を与える甲状軟骨の運動によって生み出されており、非特許文献4、5ではその制御機構の確率モデルに基づき、フレーズ・アクセント指令と呼ぶ甲状軟骨の運動に関係するパラメータを推定する技術が提案されている。この技術では、フレーズ・アクセント指令の時系列の生成プロセスを隠れマルコフモデル(HMM)により表現した点がポイントの一つであり、HMMのトポロジーの設計や遷移確率の学習を通して、指令列に関する言語学的ないし先験的な知識をパラメータ推定に組み込むことが可能である。

【0007】

本発明は、上記事情を鑑みてなされたものであり、 F_0 パターンの物理的な生成過程の制約を考慮しながらスペクトル特徴量系列に対応する最適な F_0 パターンを推定することができる基本周波数パターン予測装置、方法、及び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】

【0008】

上記の目的を達成するために本発明に係る基本周波数パターン予測装置は、学習サンプルのソース音声の時系列データとターゲット音声の時系列データとからなるパラレルデータを入力として、前記ソース音声の時系列データから抽出される各時刻のスペクトル特徴量ベクトルと、前記ターゲット音声の時系列データから抽出される、各時刻の基本周波数とに基づいて、前記ソース音声の各時刻のスペクトル特徴量ベクトルと、前記ターゲット音声の各時刻の基本周波数との間の関係をモデル化した第1確率分布のパラメータを学習 50

する第1モデルパラメータ学習部と、前記ターゲット音声の各時刻の基本周波数に基づいて、基本周波数パターン生成過程をモデル化した第2確率分布のパラメータを学習する第2モデルパラメータ学習部と、予測対象のソース音声の時系列データを入力として、前記予測対象のソース音声の時系列データから抽出される各時刻のスペクトル特徴量ベクトルと、前記第1モデルパラメータ学習部によって学習された前記第1確率分布のパラメータと、前記第2モデルパラメータ学習部によって学習された前記第2確率分布のパラメータとに基づいて、前記第1確率分布と前記第2確率分布とを用いて表される規準を大きくするように、前記予測対象のソース音声に対応する前記ターゲット音声の各時刻の基本周波数を予測する基本周波数予測部と、を含んで構成されている。

【0009】

本発明に係る基本周波数パターン予測方法は、第1モデルパラメータ学習部と、第2モデルパラメータ学習部と、基本周波数予測部とを含む基本周波数パターン予測装置における基本周波数パターン予測方法であって、前記第1モデルパラメータ学習部が、学習サンプルのソース音声の時系列データとターゲット音声の時系列データとからなるパラレルデータを入力として、前記ソース音声の時系列データから抽出される各時刻のスペクトル特徴量ベクトルと、前記ターゲット音声の時系列データから抽出される、各時刻の基本周波数とに基づいて、前記ソース音声の各時刻のスペクトル特徴量ベクトルと、前記ターゲット音声の各時刻の基本周波数との間の関係をモデル化した第1確率分布のパラメータを学習し、前記第2モデルパラメータ学習部が、前記ターゲット音声の各時刻の基本周波数に基づいて、基本周波数パターン生成過程をモデル化した第2確率分布のパラメータを学習し、前記基本周波数予測部が、予測対象のソース音声の時系列データを入力として、前記予測対象のソース音声の時系列データから抽出される各時刻のスペクトル特徴量ベクトルと、前記第1モデルパラメータ学習部によって学習された前記第1確率分布のパラメータと、前記第2モデルパラメータ学習部によって学習された前記第2確率分布のパラメータとに基づいて、前記第1確率分布と前記第2確率分布とを用いて表される規準を大きくするように、前記予測対象のソース音声に対応する前記ターゲット音声の各時刻の基本周波数を予測する。

【0010】

本発明に係るプログラムは、上記の基本周波数パターン予測装置の各部としてコンピュータを機能させるためのプログラムである。

【発明の効果】

【0011】

以上説明したように、本発明の基本周波数パターン予測装置、方法、及びプログラムによれば、ソース音声の各時刻のスペクトル特徴量ベクトルと、ターゲット音声の各時刻の基本周波数との間の関係をモデル化した第1確率分布のパラメータを学習し、基本周波数パターン生成過程をモデル化した第2確率分布のパラメータを学習し、予測対象のソース音声の時系列データから抽出される各時刻のスペクトル特徴量ベクトルから、第1確率分布と第2確率分布とを用いて表される規準を大きくするように、予測対象のソース音声に対応するターゲット音声の各時刻の基本周波数を予測することにより、 F_0 パターンの物理的な生成過程の制約を考慮しながらスペクトル特徴量系列に対応する最適な F_0 パターンを推定することができる、という効果が得られる。

【図面の簡単な説明】

【0012】

【図1】HMMの状態遷移ネットワークの一例を説明するための図である。

【図2】HMMの状態遷移ネットワークの一例を説明するための図である。

【図3】HMMの状態遷移ネットワークの一例を説明するための図である。

【図4】状態の分割を説明するための図である。

【図5】本発明の第1の実施の形態に係る基本周波数パターン予測装置の構成を示す概略図である。

【図6】本発明の第1の実施の形態に係る基本周波数パターン予測装置の学習部の構成を

10

20

30

40

50

示す概略図である。

【図7】本発明の第1の実施の形態に係る基本周波数パターン予測装置の変換処理部の構成を示す概略図である。

【図8】本発明の第1の実施の形態に係る基本周波数パターン予測装置における学習処理ルーチンの内容を示すフローチャートである。

【図9】本発明の第1の実施の形態に係る基本周波数パターン予測装置における基本周波数パターン予測処理ルーチンの内容を示すフローチャートである。

【図10】本発明の第2の実施の形態に係る基本周波数パターン予測装置の変換処理部の構成を示す概略図である。

【図11】本発明の第2の実施の形態に係る基本周波数パターン予測装置における学習処理ルーチンの内容を示すフローチャートである。 10

【図12】本発明の第2の実施の形態に係る基本周波数パターン予測装置における基本周波数パターン予測処理ルーチンの内容を示すフローチャートである。

【図13】実験データに用いた音声の F_0 パターンを示す図である。

【図14】実験結果を示す図である。

【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実施の形態を詳細に説明する。本発明で提案する技術は、音声の特徴量系列から基本周波数パターンを予測し、原音声の基本周波数パターンを予測した基本周波数パターンに置き換えることで音声の自然性を向上させることを目的とした音声処理技術である。 20

【0014】

<関連技術1：スペクトル特徴量系列からの F_0 パターン予測方法>

まず、スペクトル特徴量系列からの F_0 パターン予測方法について説明する。

【0015】

上記非特許文献1～非特許文献3では、スペクトル特徴量系列から F_0 パターンを予測する方法が提案されている。当該従来手法はスペクトル特徴量系列と F_0 パターンの同時確率分布モデルのパラメータを学習する処理と学習した当該モデルを用いて所与のスペクトル特徴量系列から F_0 パターンに変換する処理からなる。 30

【0016】

<学習処理>

ソース音声（例えば電気音声）とターゲット音声（例えば自然音声）の平行データが与えられているものとする。ソース音声のスペクトル特徴量ベクトルを $c[k]$ とし、ターゲット音声の対数 F_0 とその動的成分（時間微分または時間差分）の結合ベクトル（ F_0 特徴量と呼ぶ。）を $q[k] = (y[k]; y[k])^T$ とする。ここで k は離散時刻のインデックスである。音声特徴量 $c[k]$ としては例えば時刻 k を中心とした前後数フレーム分のメルケプストラム（ベクトル）の系列を連結したベクトルに対し主成分分析により次元圧縮を行ったものを用いる。本手法では $c[k]$ と $q[k]$ の同時確率分布を混合正規分布モデル（Gaussian Mixture Model; GMM） 40

【0017】

【数1】

$$P(c[k], q[k]|\gamma) = \sum_m \alpha_m \mathcal{N} \left(\begin{bmatrix} c[k] \\ q[k] \end{bmatrix}; \begin{bmatrix} \mu_m^{(c)} \\ \mu_m^{(q)} \end{bmatrix}, \begin{bmatrix} \Sigma_m^{(cc)} & \Sigma_m^{(cq)} \\ \Sigma_m^{(qc)} & \Sigma_m^{(qq)} \end{bmatrix} \right) \quad (1)$$

$$\sum_m \alpha_m = 1 \quad (2)$$

【0018】

でモデル化し、学習処理では所与の平行データ $\{c[k]; q[k]\}_{k=1}^K$ から当該GMMの 50

パラメータ（各正規分布の重み，平均、分散共分散行列）を学習する。ただし、 $N(x; \mu, \Sigma)$ は、 x の確率密度関数が平均が μ 、分散共分散行列が Σ の正規分布で与えられることを意味する。

【 0 0 1 9 】

GMM のパラメータはExpectation-Maximization (EM) アルゴリズムにより推定することができる。学習したGMM パラメータを $\hat{\gamma}$ とすると、条件付分布 $P(q[k] | c[k], \hat{\gamma})$ はスペクトル特徴量 $c[k]$ から F_0 特徴量 $q[k]$ を予測するための分布と見ることができ、

【 0 0 2 0 】

【数 2】

$$P(q[k] | c[k], \hat{\gamma}) = \sum_m P(m | c[k], \hat{\gamma}) P(q[k] | c[k], m, \hat{\gamma}) \quad (3) \quad 10$$

【 0 0 2 1 】

のように $P(c[k], q[k] | \hat{\gamma})$ と同様にGMM で与えられる。ただし、

【 0 0 2 2 】

【数 3】

$$P(m | c[k], \hat{\gamma}) = \frac{\hat{\alpha}_m \mathcal{N}(c[k]; \hat{\mu}_m^{(e)}, \hat{\Sigma}_m^{(ee)})}{\sum_n \hat{\alpha}_n \mathcal{N}(c[k]; \hat{\mu}_n^{(e)}, \hat{\Sigma}_n^{(ee)})} \quad (4) \quad 20$$

$$P(q[k] | c[k], m, \hat{\gamma}) = \mathcal{N}(q[k]; e_m^{(q|e)}, D_m^{(q|e)}) \quad (5)$$

【 0 0 2 3 】

であり $e_m^{(q|e)}$ および $D_m^{(q|e)}$ は

【 0 0 2 4 】

【数 4】

$$e_m^{(q|e)} = \hat{\mu}_m^{(q)} + \hat{\Sigma}_m^{(qe)} \hat{\Sigma}_m^{(ee)^{-1}} (c[k] - \hat{\mu}_m^{(e)}) \quad (6)$$

$$D_m^{(q|e)} = \hat{\Sigma}_m^{(qq)} - \hat{\Sigma}_m^{(qe)} \hat{\Sigma}_m^{(ee)^{-1}} \hat{\Sigma}_m^{(eq)} \quad (7) \quad 30$$

【 0 0 2 5 】

で与えられる。

【 0 0 2 6 】

<変換処理>

変換処理では、所与のスペクトル特徴量系列

【 0 0 2 7 】

【数 5】

$$\mathbf{c} = (c[1]^T, \dots, c[K]^T)^T \quad 40$$

【 0 0 2 8 】

の下で、最尤の F_0 パターン

【 0 0 2 9 】

【数 6】

$$\mathbf{y} = (y[1], \dots, y[K])^T$$

【 0 0 3 0 】

を以下の式(8)により求めることが目的である。

【0031】

【数7】

$$\hat{\mathbf{y}} = \underset{\mathbf{y}}{\operatorname{argmax}} \log P(\mathbf{q}|\mathbf{c}, \gamma) \quad \text{subject to } \mathbf{q} = \mathbf{W}\mathbf{y} \quad (8)$$

【0032】

ただし、

【0033】

【数8】

$$\mathbf{q} = (\mathbf{q}[1]^\top, \dots, \mathbf{q}[K]^\top)^\top$$

【0034】

であり、Wはyとqの関係を示した変換行列(定数)である。ここで、 $P(\mathbf{q}|\mathbf{c}, \gamma)$ は学習処理により学習したパラメータ γ で与えられるGMMであり、

【0035】

【数9】

$$P(\mathbf{q}|\mathbf{c}, \hat{\gamma}) = \prod_k P(\mathbf{q}[k]|\mathbf{c}[k], \hat{\gamma}) \quad (9) \quad 20$$

$$= \prod_k \sum_{m_k} P(m_k|\mathbf{c}[k], \hat{\gamma}) P(\mathbf{q}[k]|\mathbf{c}[k], m_k, \hat{\gamma}) \quad (10)$$

【0036】

で与えられる。ただし、 $\mathbf{m} = (m_1, \dots, m_K)$ であり、 m_k は時刻kにおけるGMMの成分インデックスを表す。ここで、 $P(\mathbf{q}[k]|\mathbf{c}[k], \mathbf{m})$ が

【0037】

【数10】

$$P(\mathbf{q}[k]|\mathbf{c}[k], \hat{\gamma}) \simeq P(\hat{m}_k|\mathbf{c}[k], \hat{\gamma}) P(\mathbf{q}[k]|\mathbf{c}[k], \hat{m}_k, \hat{\gamma}) \quad (11) \quad 30$$

$$\hat{m}_k = \underset{m_k}{\operatorname{argmax}} P(m_k|\mathbf{c}[k], \hat{\gamma}) \quad (12)$$

【0038】

により近似できるとする。式(12)より、 \hat{m}_k は、データ $\mathbf{c}[k]$ を生成したらしい確率が最も高い正規分布のインデックスを意味する。よって、式(9)より $P(\mathbf{q}|\mathbf{c}, \mathbf{m}, \hat{\gamma})$ はすべてのkについて

【0039】

【数11】

$$P(\mathbf{q}[k]|\mathbf{c}[k], \hat{m}_k, \hat{\gamma}) = \mathcal{N}(\mathbf{q}[k]; \mathbf{e}_{\hat{m}_k}^{(q|c)}, \mathbf{D}_{\hat{m}_k}^{(q|c)})$$

【0040】

の積をとったもので与えられる。式(11)の近似により $P(\mathbf{q}|\mathbf{c}, \mathbf{m}, \hat{\gamma})$ は

【0041】

【数12】

$$\mathbf{e}_{\hat{m}_1}^{(q|c)}, \dots, \mathbf{e}_{\hat{m}_K}^{(q|c)}$$

10

20

30

40

50

【 0 0 4 2 】

を連結したベクトル $e^{(q|c)}$ を平均、

【 0 0 4 3 】

【 数 1 3 】

$$D_{\hat{m}_1}^{(q|c)}, \dots, D_{\hat{m}_K}^{(q|c)}$$

【 0 0 4 4 】

を対角成分としたブロック対角行列 $D^{(q|c)}$ を分散共分散行列とした正規分布

【 0 0 4 5 】

【 数 1 4 】

$$P(q|c, \gamma) = \mathcal{N}(q; e^{(q|c)}, D^{(q|c)}) \quad (13)$$

【 0 0 4 6 】

となる。これに $q = Wy$ を代入し、 y の分布となるように正規化すると、

【 0 0 4 7 】

【 数 1 5 】

$$\begin{aligned} \mathcal{N}(Wy; e^{(q|c)}, D^{(q|c)}) \\ \propto \mathcal{N}(y; (W^T D^{(q|c)-1} W)^{-1} W^T D^{(q|c)-1} e^{(q|c)}, (W^T D^{(q|c)-1} W)^{-1}) \end{aligned} \quad (14)$$

【 0 0 4 8 】

となるため、

【 0 0 4 9 】

【 数 1 6 】

$$P(y|c, \gamma) = \mathcal{N}(y; (W^T D^{(q|c)-1} W)^{-1} W^T D^{(q|c)-1} e^{(q|c)}, (W^T D^{(q|c)-1} W)^{-1}) \quad (15)$$

【 0 0 5 0 】

のように y の条件付き分布を得る。よって、式(8)の解は、

【 0 0 5 1 】

【 数 1 7 】

$$\hat{y} = (W^T D^{(q|c)-1} W)^{-1} W^T D^{(q|c)-1} e^{(q|c)} \quad (16)$$

【 0 0 5 2 】

となる。

【 0 0 5 3 】

< 関連技術 2 : F_0 パターン生成過程モデル >

次に、 F_0 パターン生成過程の確率モデルについて説明する。

【 0 0 5 4 】

音声の F_0 パターンの生成過程を記述したモデルに、藤崎の基本周波数 (F_0) パターン生成過程モデル (藤崎モデル) が知られている (非特許文献 6)。

【 0 0 5 5 】

[非特許文献 6] : H. Fujisaki, "In Vocal Physiology: Voice Production, Mechanisms and Functions," Raven Press, 1988.

【 0 0 5 6 】

藤崎モデルとは、甲状軟骨の運動による F_0 パターンの生成過程を説明した物理モデルである。藤崎モデルでは、甲状軟骨の二つの独立な運動 (平行移動運動と回転運動) にそれぞれ伴う声帯の伸びの合計が F_0 の時間的变化をもたらすと解釈され、声帯の伸びと F_0 パ

10

20

30

40

50

ターンの対数値 $y(t)$ が比例関係にあるという仮定に基づいて F_0 パターンがモデル化される。甲状軟骨の平行移動運動によって生じる F_0 パターン $x_p(t)$ をフレーズ成分、回転運動によって生じる F_0 パターン $x_a(t)$ をアクセント成分と呼ぶ。藤崎モデルでは、音声の F_0 パターン $y(t)$ は、これらの成分に声帯の物理的制約によって決まるベースライン成分 b を足し合わせたものとして、

【 0 0 5 7 】

【数 1 8】

$$y(t) = x_p(t) + x_a(t) + \mu b \quad (17)$$

10

【 0 0 5 8 】

と表現される。これら二つの成分は二次の臨界制動系の出力と仮定され、

【 0 0 5 9 】

【数 1 9】

$$x_p(t) = g_p(t) * u_p(t), \quad (18)$$

$$g_p(t) = \begin{cases} \alpha^2 t e^{-\alpha t} & (t \geq 0) \\ 0 & (t < 0) \end{cases}, \quad (19)$$

$$x_a(t) = g_a(t) * u_a(t), \quad (20)$$

$$g_a(t) = \begin{cases} \beta^2 t e^{-\beta t} & (t \geq 0) \\ 0 & (t < 0) \end{cases}, \quad (21)$$

20

【 0 0 6 0 】

と表される（* は時刻 t に関する畳み込み演算）。ここで $u_p(t)$ はフレーズ指令関数と呼ばれ、デルタ関数（フレーズ指令）の列からなり、 $u_a(t)$ はアクセント指令関数と呼ばれ、矩形波（アクセント指令）の列からなる。これらの指令列には、発話の最初にはフレーズ指令が生起する、フレーズ指令は二連続で生起しない、異なる二つの指令は同時刻に生起しない、という制約条件がある。また α はそれぞれフレーズ制御機構、アクセント制御機構の固有角周波数であり、話者や発話内容によらず、おおよそ $\alpha = 3$ rad/s、 $\beta = 20$ rad/s 程度であることが経験的に知られている。

30

【 0 0 6 1 】

< 関連技術 3 : F_0 パターン生成過程モデルパラメータ推定法 >

上述の藤崎モデルは以下のような確率モデルで記述することができる（非特許文献 4、5、7 参照）。

【 0 0 6 2 】

[非特許文献 7]：石原達馬，吉里幸太，亀岡弘和，齋藤大輔，嵯峨山茂樹，"音声基本周波数の藤崎モデル指令列の統計的語彙モデル," 日本音響学会2013年春季研究発表会講演論文集, 1-7-9, pp. 283-286, Mar. 2013.

40

【 0 0 6 3 】

まずフレーズ、アクセント指令関数のペア $o[k] = (u_p[k], u_a[k])^T$ を出力するHMMを考える。ただし、 k は離散時刻のインデックスを表す。状態出力分布は正規分布とし、各時刻の状態が与えられた下で

【 0 0 6 4 】

【数 2 0】

$$o[k] \sim \mathcal{N}(o[k]; \rho[k], \Upsilon[k]) \quad (22)$$

$$\rho[k] = \begin{bmatrix} \mu_{s_k}^{(p)}[k] \\ \mu_{s_k}^{(a)} \end{bmatrix}, \quad \Upsilon[k] = \begin{bmatrix} v_{s_k}^{(p)2} & 0 \\ 0 & v_{s_k}^{(a)2} \end{bmatrix} \quad (23)$$

【 0 0 6 5】

により指令関数ペア $o[k]$ が生成されるものとする。ここで $\{s_k\}_{k=1}^K$ はHMMの状態系列であり、平均ベクトル $\rho[k]$ はHMMの状態遷移の結果として定まる値である。具体的なHM 10
Mの構成の例を図1～図3に示す。

【 0 0 6 6】

図1に示すHMMの状態遷移ネットワークの例では、状態 $t = r_0$ において $\mu^{(p)}_t[k]$ と $\mu^{(a)}_t$ はいずれも0である。状態 $t = r_0$ からは状態 p_0 にのみ遷移することができ、状態 $t = p_0$ において $\mu^{(p)}_t[k]$ は非負値 $A^{(p)}[k]$ をとり、 $\mu^{(a)}_t$ は0となる。状態 $t = p_0$ の次は状態 r_1 にのみ遷移することが許される。状態 $t = r_0$ 同様、状態 $t = r_1$ において $\mu^{(p)}_t[k]$ と $\mu^{(a)}_t$ はいずれも0である。状態 $t = r_1$ からは状態 a_0, \dots, a_{N-1} のいずれかにのみ遷移することができ、状態 $t = a_n$ において $\mu^{(a)}_t$ は非負値 $A^{(a)}_n$ をとり、 $\mu^{(p)}_t[k]$ は0となる。状態 $t = a_n$ の次は状態 r_0 または r_1 にのみ遷移することが許される。これより $\mu_a[k]$ が矩形パルス列となることが保証される。 20

【 0 0 6 7】

図2に示すHMMの状態遷移ネットワークの例では、状態 $t = r_0$ において $\mu^{(p)}_t[k]$ と $\mu^{(a)}_t$ はいずれも0である。状態 $t = r_0$ からは状態 p_0, \dots, p_{M-1} のいずれかにのみ遷移することができ、状態 $t = p_m$ において $\mu^{(p)}_t$ は非負値 $A^{(p)}_m$ をとり、 $\mu^{(a)}_t$ は0となる。状態 $t = p_m$ の次は状態 r_1 にのみ遷移することが許される。状態 $t = r_0$ 同様、状態 $t = r_1$ において $\mu^{(p)}_t[k]$ と $\mu^{(a)}_t$ はいずれも0である。状態 $t = r_1$ からは状態 a_0, \dots, a_{N-1} のいずれかにのみ遷移することができ、状態 $t = a_n$ において $\mu^{(a)}_t$ は非負値 $A^{(a)}_n$ をとり、 $\mu^{(p)}_t$ は0となる。状態 $t = a_n$ の次は状態 r_0 または r_1 にのみ遷移することが許される。これより $\mu_a[k]$ が矩形パルス列となることが保証される。 30

【 0 0 6 8】

図3に示すHMMの状態遷移ネットワークの例では、それぞれの終点と始点が連結された複数のLeft-to-Right型HMMからなる。図1、2と同様、状態 $t = r_1$ において $\mu^{(p)}_t[k]$ と $\mu^{(a)}_t$ はいずれも0である。また、状態 $t = p_m$ において $\mu^{(p)}_t$ は非負値 $A^{(p)}_m$ をとり、 $\mu^{(a)}_t$ は0となる。状態 $t = a_n$ において $\mu^{(a)}_t$ は非負値 $A^{(a)}_n$ をとり、 $\mu^{(p)}_t$ は0となる。 30

【 0 0 6 9】

$\rho[k]$ は、図1の例では、以下の式(24)で表わされる。

【 0 0 7 0】

【数 2 1】

$$\rho[k] = \begin{bmatrix} \mu_{s_k}^{(p)}[k] \\ \mu_{s_k}^{(a)} \end{bmatrix} = \begin{cases} (0, 0)^T & (s_k \in r_l) \\ (A^{(p)}[k], 0)^T & (s_k \in p_0) \\ (0, A_n^{(a)})^T & (s_k \in a_n) \end{cases} \quad (24)$$

【 0 0 7 1】

また、図2、3の例では、 $\rho[k]$ は、以下の式(25)で表わされる。

【 0 0 7 2】

10

20

30

40

【数 2 2】

$$\rho[k] = \begin{bmatrix} \mu_{s_k}^{(p)} \\ \mu_{s_k}^{(a)} \end{bmatrix} = \begin{cases} (0, 0)^T & (s_k \in r_l) \\ (A_m^{(p)}, 0)^T & (s_k \in p_m) \\ (0, A_n^{(a)})^T & (s_k \in a_n) \end{cases} \quad (25)$$

【0 0 7 3】

いずれの例においても、図 4 のようにそれぞれの状態を同じ出力分布をもついくつかの小状態に分割し、Left-to-Right 型の状態遷移経路を制約することで同一状態に停留する時間長の確率をパラメータ化することができる。図 4 は状態 a_n を分割した例である。例えばこの図のように全ての $m = 0$ に対して $a_{n,n}$ から $a_{n,n+1}$ への状態遷移確率を 1 に設定することで、 $a_{n,0}$ から $a_{n,n}$ への遷移確率が状態 a_n が n ステップだけ持続する確率に対応し、アクセント指令の持続長の確率を設定したり学習したりできるようになる。同様に p_m と r_l も小状態に分割することで、フレーズ指令の持続長と指令間の間隔の長さの分布をパラメータ化することが可能になる。以後、状態集合を

10

【0 0 7 4】

【数 2 3】

$$r_l = \{r_{l,0}, r_{l,1}, \dots\}, \quad p_m = \{p_{m,0}, p_{m,1}, \dots\}, \quad a_n = \{a_{n,0}, a_{n,1}, \dots\}$$

【0 0 7 5】

と表記する。上記の HMM の構成は次のように書ける。

20

【0 0 7 6】

【数 2 4】

出力値系列: $\{o[k]\}_{k=1}^K$
 状態集合: $\mathcal{S} = \{r_0, \dots, r_{L-1}, p_0, \dots, p_{M-1}, a_0, \dots, a_{N-1}\}$
 状態系列: $\{s_k\}_{k=1}^K$
 状態出力分布: $P(o[k]|s_k) = \mathcal{N}(o[k]; \rho[k], \Upsilon[k])$

$$\rho[k] = \begin{bmatrix} \mu_{s_k}^{(p)}[k] \\ \mu_{s_k}^{(a)} \end{bmatrix}$$

$$\Upsilon[k] = \begin{bmatrix} v_{s_k}^{(p)2} & 0 \\ 0 & v_{s_k}^{(a)2} \end{bmatrix}$$

状態遷移確率: $\phi_{i',i} = \log P(s_k = i | s_{k-1} = i')$

30

【0 0 7 7】

状態系列 $s = \{s_k\}_{k=1}^K$ が与えられたとき、この HMM はフレーズ指令関数 $u_p[k]$ とアクセント指令関数 $u_a[k]$ のペアを出力する。式(18) と式(20) で示した通り、 $u_p[k]$ と $u_a[k]$ にそれぞれ $g_p[k]$ と $g_a[k]$ が畳み込まれてフレーズ成分 $x_p[k]$ とアクセント成分 $x_a[k]$ が出力される。これを式で表すと、

40

【0 0 7 8】

【数 2 5】

$$x_p[k] = g_p[k] * u_p[k], \quad (26)$$

$$x_a[k] = g_a[k] * u_a[k], \quad (27)$$

【0 0 7 9】

50

と書ける（* は離散時刻 k に関する畳み込み演算）。このとき、 F_0 パターン $x[k]$ は
 【 0 0 8 0 】
 【 数 2 6 】

$$x[k] = x_p[k] + x_a[k] + \mu_b, \quad (28)$$

【 0 0 8 1 】

と三種類の成分の重ね合わせで書ける。ただし b は時刻によらないベースライン成分である。

【 0 0 8 2 】

また、実音声においては、いつも信頼のできる F_0 の値が観測できるとは限らない。藤崎モデルのパラメータ推定を行うにあたっては、信頼のおける観測区間の F_0 値のみを考慮に入れて、そうでない区間は無視することが望ましい。例えば音声の無声区間においては通常声帯の振動に伴う周期的な粗密波は観測されないため、仮に自動ピッチ抽出によって音声の無声区間から何らかの値が F_0 の推定値として得られたとしても、その値を声帯から発せられる信号の F_0 の値と見なすのは適当ではない。そこで、提案モデルに観測 F_0 値の時刻 k における不確かさの程度 $v_n^2[k]$ を導入する。具体的には、観測 F_0 値 $y[k]$ を、真の F_0 値 $x[k]$ とノイズ成分

【 0 0 8 3 】

【 数 2 7 】

$$x_n[k] \sim \mathcal{N}(0, v_n^2[k])$$

【 0 0 8 4 】

との重ね合わせで

【 0 0 8 5 】

【 数 2 8 】

$$y[k] = x[k] + x_n[k] \quad (29)$$

【 0 0 8 6 】

と表現することで、信頼のおける区間かどうかに関わらず全ての観測区間を統一的に扱える。

【 0 0 8 7 】

$x_n[k]$ を周辺化することで、出力値系列 $o = \{o[k]\}_{k=1}^K$ が与えられたときの $y = \{y[k]\}_{k=1}^K$ の確率密度関数

【 0 0 8 8 】

【 数 2 9 】

$$P(y|o) = \prod_{k=1}^K \mathcal{N}(y[k]; x[k], v_n^2[k]) \quad (30)$$

【 0 0 8 9 】

が得られる。状態系列 $s = \{s_k\}_{k=1}^K$ と指令の振幅を表すパラメータ

【 0 0 9 0 】

【 数 3 0 】

$$\theta = \begin{cases} \{\{A^{(p)}[k]\}_{k=1}^K, \{A_n^{(a)}\}_{n=1}^N\} & (\text{図 1 の場合}) \\ \{\{A_m^{(p)}\}_{m=1}^M, \{A_n^{(a)}\}_{n=1}^N\} & (\text{図 2, 3 の場合}) \end{cases} \quad (31)$$

10

20

30

40

50

【 0 0 9 1 】

および遷移確率行列 $\mathbf{A} = (a_{i,j})_{I \times I}$ が与えられたとき、出力値系列 \mathbf{o} は

【 0 0 9 2 】

【数 3 1】

$$P(\mathbf{o}|\mathbf{s}, \boldsymbol{\theta}) = \prod_{k=1}^K \mathcal{N}(\mathbf{o}[k]; \boldsymbol{\rho}[k], \boldsymbol{\Upsilon}[k]) \quad (32)$$

【 0 0 9 3 】

に従って生成される。また、 $P(\mathbf{s} | \cdot)$ は状態遷移確率の積として

10

【 0 0 9 4 】

【数 3 2】

$$\log P(\mathbf{s}) = \phi_{s_1} + \sum_{k=2}^K \phi_{s_{k-1}, s_k} \quad (33)$$

と書ける。ただし、

【 0 0 9 5 】

【数 3 3】

20

$$\phi_{s_1}$$

【 0 0 9 6 】

は初期状態が s_1 である確率をあらわす。式(30)、(32) および式(33) より $P(\mathbf{y}, \mathbf{o}, \mathbf{s} | \cdot, \cdot)$ は

【 0 0 9 7 】

【数 3 4】

$$P(\mathbf{y}, \mathbf{o}, \mathbf{s} | \boldsymbol{\theta}, \boldsymbol{\phi}) = P(\mathbf{y}|\mathbf{o})P(\mathbf{o}|\mathbf{s}, \boldsymbol{\theta})P(\mathbf{s}|\boldsymbol{\phi}) \quad (34) \quad 30$$

【 0 0 9 8 】

と書ける。これを \mathbf{o} に関して周辺化すると

【 0 0 9 9 】

【数 3 5】

$$P(\mathbf{y}, \mathbf{s} | \boldsymbol{\theta}, \boldsymbol{\phi}) = P(\mathbf{y}|\mathbf{s}, \boldsymbol{\theta})P(\mathbf{s}|\boldsymbol{\phi}) \quad (35)$$

$$P(\mathbf{y}|\mathbf{s}, \boldsymbol{\theta}) = \mathcal{N}(\mathbf{y}; \mathbf{G}_p \boldsymbol{\mu}_p + \mathbf{G}_a \boldsymbol{\mu}_a + \boldsymbol{\mu}_b \mathbf{1}, \mathbf{G}_p \boldsymbol{\Upsilon}_p \mathbf{G}_p^\top + \mathbf{G}_a \boldsymbol{\Upsilon}_a \mathbf{G}_a^\top + \boldsymbol{\Upsilon}_n) \quad (36)$$

【 0 1 0 0 】

が得られる。ただし、

40

【 0 1 0 1 】

【数 3 6】

$$G_p = \begin{bmatrix} g_p[1] & & & & O \\ g_p[2] & g_p[1] & & & \\ g_p[3] & g_p[2] & g_p[1] & & \\ \vdots & \ddots & \ddots & \ddots & \\ g_p[K] & \dots & g_p[3] & g_p[2] & g_p[1] \end{bmatrix} \quad (37)$$

$$G_a = \begin{bmatrix} g_a[1] & & & & O \\ g_a[2] & g_a[1] & & & \\ g_a[3] & g_a[2] & g_a[1] & & \\ \vdots & \ddots & \ddots & \ddots & \\ g_a[K] & \dots & g_a[3] & g_a[2] & g_a[1] \end{bmatrix} \quad (38)$$

$$\mu_p = \begin{cases} (\mu_{s_1}^{(p)}[1], \mu_{s_2}^{(p)}[2], \dots, \mu_{s_K}^{(p)}[K])^T & (\text{図 1 の場合}) \\ (\mu_{s_1}^{(p)}, \mu_{s_2}^{(p)}, \dots, \mu_{s_K}^{(p)})^T & (\text{図 2, 3 の場合}) \end{cases} \quad (39)$$

$$\mu_a = (\mu_{s_1}^{(a)}, \mu_{s_2}^{(a)}, \dots, \mu_{s_K}^{(a)})^T \quad (40)$$

$$\mathbf{1} = \underbrace{(1, 1, \dots, 1)}_K^T \quad (41)$$

$$Y_p = \begin{bmatrix} v_{s_1}^{(p)} & & O \\ & \ddots & \\ O & & v_{s_K}^{(p)} \end{bmatrix} \quad (42)$$

$$Y_a = \begin{bmatrix} v_{s_1}^{(a)} & & O \\ & \ddots & \\ O & & v_{s_K}^{(a)} \end{bmatrix} \quad (43)$$

$$Y_n = \begin{bmatrix} v_n^2[1] & & O \\ & \ddots & \\ O & & v_n^2[K] \end{bmatrix} \quad (44)$$

【0 1 0 2】

である。一方、s に関して周辺化すると

【0 1 0 3】

【数 3 7】

$$P(y, o | \theta, \phi) = P(y|o) \sum_s P(o|s, \theta) P(s|\phi) \quad (45)$$

が得られる。ただし、 \sum_s はあらゆる状態系列に関して和をとる操作を意味する。

【0 1 0 4】

< パラメータ推定アルゴリズム 1 >

y と o を完全データと見なすことで、式(35) を局所最大化する s と E を Expectation-

10

20

30

40

50

Maximization アルゴリズムにより探索することができる。導出は省略するが、

【 0 1 0 5 】

【数 3 8】

$$\begin{aligned}
 Q(s, \theta) = & -\frac{1}{2} \left[\log |G_p \Upsilon_p G_p^\top| + \log |G_p \Upsilon_p G_p^\top| \right. \\
 & + \text{tr}(\Upsilon_p^{-1} G_p^{-1} R_p G_p^{-\top}) + \text{tr}(\Upsilon_a^{-1} G_a^{-1} R_a G_a^{-\top}) \\
 & - 2\mu_p^\top \Upsilon_p^{-1} G_p^{-1} \bar{x}_p - 2\mu_a^\top \Upsilon_a^{-1} G_a^{-1} \bar{x}_a - 2\mu_b^\top \mathbf{1}^\top \Upsilon_b^{-1} \bar{x}_b \\
 & \left. + \mu_p^\top \Upsilon_p^{-1} \mu_p + \mu_a^\top \Upsilon_a^{-1} \mu_a + \mu_b^2 \mathbf{1}^\top \Upsilon_b^{-1} \mathbf{1} \right] + \log P(s)
 \end{aligned} \tag{46}$$

10

【 0 1 0 6 】

が大きくなるように s と を更新するステップと、更新した s と を用いて

【 0 1 0 7 】

【数 3 9】

$$\bar{\mathbf{x}} = (\bar{\mathbf{x}}_p^\top, \bar{\mathbf{x}}_a^\top, \bar{\mathbf{x}}_b^\top)^\top$$

【 0 1 0 8 】

と R を

20

【 0 1 0 9 】

【数 4 0】

$$\bar{\mathbf{x}} = \mathbf{a} + \Lambda H^\top (H \Lambda H^\top)^{-1} (\mathbf{y} - H \mathbf{a}) \tag{47}$$

$$R = \Lambda - \Lambda H^\top (H \Lambda H^\top)^{-1} H \Lambda + \bar{\mathbf{x}} \bar{\mathbf{x}}^\top \tag{48}$$

【 0 1 1 0 】

により更新するステップを繰り返すことで式(35) を単調増加させることができる (詳細は、上記非特許文献 4 参照)。

【 0 1 1 1 】

30

具体的には、以下の初期設定、Eステップ、及びMステップが実行される。

【 0 1 1 2 】

(初期設定)

s と を初期設定する。

【 0 1 1 3 】

(Eステップ)

フレーズ成分、アクセント成分、ベースライン成分の条件付き期待値

【 0 1 1 4 】

【数 4 1】

$$\bar{\mathbf{x}} = (\bar{\mathbf{x}}_p^\top, \bar{\mathbf{x}}_a^\top, \bar{\mathbf{x}}_b^\top)^\top$$

40

【 0 1 1 5 】

と条件付き分散共分散行列 R を

【 0 1 1 6 】

【数 4 2】

$$\bar{x} \leftarrow a + \Lambda H^T (H \Lambda H^T)^{-1} (y - Ha) \quad (49)$$

$$R \leftarrow \Lambda - \Lambda H^T (H \Lambda H^T)^{-1} H \Lambda + \bar{x} \bar{x}^T \quad (50)$$

により更新する。ただし、

【0 1 1 7】

【数 4 3】

$$a = \begin{bmatrix} G_p \mu_p \\ G_a \mu_a \\ \mu_b \mathbf{1} \end{bmatrix} \quad (51) \quad 10$$

$$\Lambda = \begin{bmatrix} G_p \Upsilon_p G_p^T & O & O \\ O & G_a \Upsilon_a G_a^T & O \\ O & O & \Upsilon_n \end{bmatrix} \quad (52)$$

20

$$H = [I \quad I \quad I] \quad (53)$$

【0 1 1 8】

である。また、Rにおける

【0 1 1 9】

【数 4 4】

$$\bar{x}_p, \bar{x}_a, \bar{x}_b$$

30

【0 1 2 0】

に対応するブロック対角成分を

【0 1 2 1】

【数 4 5】

$$R_p, R_a, R_b$$

【0 1 2 2】

とする。

【0 1 2 3】

すなわち、

【0 1 2 4】

【数 4 6】

$$R = \begin{bmatrix} R_p & * & * \\ * & R_a & * \\ * & * & R_b \end{bmatrix} \quad (54)$$

40

【0 1 2 5】

50

である（* は以後用いないブロック成分である）。

【 0 1 2 6 】

（ Mステップ）

$Q(s, \cdot)$ が最大となる状態系列 $s = (s_1, \dots, s_K)$ を探索する。 μ_p と μ_a は対角行

【 0 1 2 7 】

列なので、

【 0 1 2 8 】

【数 4 7】

$\log|\Upsilon_p|$ と $\log|\Upsilon_a|$, $\text{tr}(\Upsilon_p^{-1}G_p^{-1}R_pG_p^{-T})$ と $\text{tr}(\Upsilon_a^{-1}G_a^{-1}R_aG_a^{-T})$, $\mu_p^T\Upsilon_p^{-1}G_p^{-1}\bar{x}_p$

10

と $\mu_a^T\Upsilon_a^{-1}G_a^{-1}\bar{x}_a$, $\mu_p^T\Upsilon_p^{-1}\mu_p$ と $\mu_a^T\Upsilon_a^{-1}\mu_a$

【 0 1 2 9 】

はいずれも

【 0 1 3 0 】

【数 4 8】

$$\log|\Upsilon_p| = \sum_k \log v_{s_k}^{(p)2} \quad (55)$$

20

$$\log|\Upsilon_a| = \sum_k \log v_{s_k}^{(p)2} \quad (56)$$

$$\text{tr}(\Upsilon_p^{-1}G_p^{-1}R_pG_p^{-T}) = \sum_k \frac{[G_p^{-1}R_pG_p^{-T}]_{k,k}}{v_{s_k}^{(p)2}} \quad (57)$$

$$\text{tr}(\Upsilon_a^{-1}G_a^{-1}R_aG_a^{-T}) = \sum_k \frac{[G_a^{-1}R_aG_a^{-T}]_{k,k}}{v_{s_k}^{(p)2}} \quad (58)$$

$$\mu_p^T\Upsilon_p^{-1}G_p^{-1}\bar{x}_p = \sum_k \frac{\mu_{s_k}^{(p)}[k][G_p^{-1}\bar{x}_p]_k}{v_{s_k}^{(p)2}} \quad (59)$$

30

$$\mu_a^T\Upsilon_a^{-1}G_a^{-1}\bar{x}_a = \sum_k \frac{\mu_{s_k}^{(a)}[G_a^{-1}\bar{x}_a]_k}{v_{s_k}^{(a)2}} \quad (60)$$

$$\mu_p^T\Upsilon_p^{-1}\mu_p = \sum_k \frac{\mu_{s_k}^{(p)}[k]^2}{v_{s_k}^{(p)2}} \quad (61)$$

40

$$\mu_a^T\Upsilon_a^{-1}\mu_a = \sum_k \frac{\mu_{s_k}^{(a)2}}{v_{s_k}^{(a)2}} \quad (62)$$

【 0 1 3 1 】

のように k ごとの項の和の形で書ける。従って、 $Q(s, \cdot)$ は s に依らない項を除けば

【 0 1 3 2 】

【数 4 9】

$$\mathcal{J}(s) = -\frac{1}{2} \sum_{k=1}^K \left(\log v_{s_k}^{(p)2} + \log v_{s_k}^{(a)2} + \frac{[\mathbf{G}_p^{-1} \mathbf{R}_p \mathbf{G}_p^{-T}]_{k,k}}{v_{s_k}^{(p)2}} + \frac{[\mathbf{G}_a^{-1} \mathbf{R}_a \mathbf{G}_a^{-T}]_{k,k}}{v_{s_k}^{(a)2}} \right. \\ \left. - 2 \frac{\mu_{s_k}^{(p)} [k] [\mathbf{G}_p^{-1} \bar{\mathbf{x}}_p]_k}{v_{s_k}^{(p)2}} - 2 \frac{\mu_{s_k}^{(a)} [k] [\mathbf{G}_a^{-1} \bar{\mathbf{x}}_a]_k}{v_{s_k}^{(a)2}} + \frac{\mu_{s_k}^{(p)} [k]^2}{v_{s_k}^{(p)2}} + \frac{\mu_{s_k}^{(a)2}}{v_{s_k}^{(a)2}} \right) + \phi_{s_1} + \sum_{k=2}^K \phi_{s_{k-1}, s_k} \quad (63)$$

【0133】

と書ける。従って、 $Q(s, \cdot)$ を最大にする状態系列 $s = (s_1, \dots, s_K)$ は Viterbi アルゴリズムにより求めることができる（詳細は上記非特許文献 4 参照）。ただし、 $[\cdot]_{k,k}$ は行列の k 行 k 列の要素、 $[\cdot]_k$ はベクトルの第 k 要素を表す。

10

【0134】

続いて、 $Q(s, \cdot)$ を最大にするように \cdot を更新する。 $Q(s, \cdot)$ を最大にする \cdot は、 $Q(s, \cdot)$ の各変数に関する偏微分を 0 と置くことにより得られる（（詳細は上記非特許文献 4 参照））。

【0135】

また、推定された状態系列 s から、状態遷移確率 \cdot が求められる。

【0136】

[第 1 の実施の形態]

<本発明の実施の形態の概要>

20

本発明の実施の形態の技術は、上述した関連技術 1 のスペクトル特徴量系列からの F_0 パターン予測方法と同様、学習処理と変換処理からなるが、式 (8) の代わりに、上述した関連技術 1 のスペクトル特徴量系列からの F_0 パターン予測方法の確率分布と、上述した関連技術 3 の F_0 パターン生成過程モデルパラメータ推定法の確率分布の積を規準とすることにより、上述した F_0 パターン生成過程モデルにできるだけ即した F_0 パターンをスペクトル特徴量から統計的に予測することを可能にする技術である。

【0137】

学習処理ではパラレルデータの学習サンプル $\{c[k], q[k]\}_{k=1}^K$ が与えられた下で

【0138】

【数 5 0】

30

$$P(c, q | \gamma) = \prod_k P(c[k], q[k] | \gamma)$$

【0139】

ができるだけ大きくなるように \cdot を学習する。また、学習サンプルの基本周波数パターン $\{y[k]\}_{k=1}^K$ が与えられた下で $P(y, s | \cdot, \cdot)$ ができるだけ大きくなるように \cdot と \cdot を学習する。

【0140】

一方、変換処理では入力音声の $\{c[k]\}_{k=1}^K$ が与えられた下で $P(q | c, \cdot) P(y, s | \cdot, \cdot)$ またはこれらを近似する確率密度関数ができるだけ大きくなるように y を求める。

40

【0141】

<システム構成>

次に、ソース音声のスペクトル特徴量系列から、ターゲット音声の基本周波数パターンを予測する基本周波数パターン予測装置に、本発明を適用した場合を例にして、本発明の実施の形態を説明する。

【0142】

図 5 に示すように、本発明の第 1 の実施の形態に係る基本周波数パターン予測装置は、CPU と、RAM と、後述する学習処理ルーチン、及び基本周波数パターン予測処理ルーチンを実行するためのプログラムを記憶した ROM とを備えたコンピュータで構成され、

50

機能的には次に示すように構成されている。

【 0 1 4 3 】

図 5 に示すように、基本周波数パターン予測装置 1 0 0 は、入力部 1 0 と、演算部 2 0 と、出力部 9 0 とを備えている。

【 0 1 4 4 】

入力部 1 0 は、学習サンプルのソース音声（例えば電気音声）の時系列データとターゲット音声（例えば自然音声）の時系列データとからなる平行データを受け付ける。また、入力部 1 0 は、予測対象のソース音声の時系列データを受け付ける。

【 0 1 4 5 】

演算部 2 0 は、学習部 3 0 と、パラメータ記憶部 4 0 と、変換処理部 5 0 とを備えている。

10

【 0 1 4 6 】

図 6 に示すように、学習部 3 0 は、特徴量抽出部 3 2 と、基本周波数系列抽出部 3 4 と、第 1 モデルパラメータ学習部 3 6 と、第 2 モデルパラメータ学習部 3 8 とを備えている。

【 0 1 4 7 】

特徴量抽出部 3 2 は、入力部 1 0 によって受け付けた学習サンプルのソース音声の時系列データから、ソース音声のスペクトラム特徴量ベクトル $c[k]$ を抽出する。ここで k は離散時刻のインデックスである。例えば、非特許文献 1 ~ 3 と同様に時刻 k を中心とした前後数フレーム分のメルケプストラム（ベクトル）の系列を連結したベクトルに対し主成分分析により次元圧縮を行ったものを $c[k]$ として用いる。

20

【 0 1 4 8 】

基本周波数系列抽出部 3 4 は、入力部 1 0 によって受け付けた学習サンプルのターゲット音声の時系列データから、ターゲット音声の各時刻 k における基本周波数 $y[k]$ を抽出し、 $y = (y[1], \dots, y[K])^T$ とする。

【 0 1 4 9 】

この基本周波数の抽出処理は、周知技術により実現でき、例えば、非特許文献 8 (H. Kameoka, "Statistical speech spectrum model incorporating all-pole vocal tract model and F0 contour generating process model," in Tech. Rep. IEICE, 2010, in Japanese.) に記載の手法を利用して、8 ms ごとに基本周波数を抽出する。

30

【 0 1 5 0 】

また、 y とその動的成分（時間微分または時間差分）の結合ベクトル (F_0 特徴量と呼ぶ。) を $q[k] = (y[k], \dot{y}[k])^T$ とする。

【 0 1 5 1 】

以上より、 $\{c[k], q[k]\}_{k=1}^K$ というデータが得られる。

【 0 1 5 2 】

第 1 モデルパラメータ学習部 3 6 は、特徴量抽出部 3 2 によって抽出された各時刻 k のスペクトル特徴量ベクトル $c[k]$ と、基本周波数系列抽出部 3 4 によって抽出された各時刻 k の基本周波数の結合ベクトル $q[k]$ とに基づいて、ソース音声の各時刻のスペクトル特徴量ベクトル $c[k]$ とターゲット音声の各時刻 k の基本周波数の結合ベクトル $q[k]$ との同時確率分布を表す混合正規分布である第 1 確率分布のパラメータを学習する。

40

【 0 1 5 3 】

具体的には、第 1 モデルパラメータ学習部 3 6 は、上述したスペクトル特徴量系列からの F_0 パターン予測方法の学習処理と同様に、式(1) の GMM のパラメータ を学習する。学習した GMM パラメータを $\hat{\theta}$ とする。

【 0 1 5 4 】

第 2 モデルパラメータ学習部 3 8 は、基本周波数系列抽出部 3 4 によって抽出された各時刻 k の基本周波数 $y[k]$ に基づいて、各時刻 k の基本周波数 $y[k]$ と、隠れマルコフモデルの各時刻の状態からなる状態系列 s との組み合わせの確率分布である第 2 確率分布のパラメータを学習する。

50

【 0 1 5 5 】

具体的には、第2モデルパラメータ学習部38は、上述した関連技術3のF₀パターン生成過程モデルパラメータ推定法のパラメータ推定アルゴリズム1に従って、F₀パターン生成過程モデルのパラメータ、を学習する。

【 0 1 5 6 】

もし学習サンプルのフレーズ指令系列とアクセント指令系列のペアのデータ $o = \{o_k\}_{k=1}^K$ が入手できるのであれば、 o から、を学習しても良い（HMMの通常の学習に相当）。学習したF₀パターン生成過程モデルのパラメータを $\hat{\theta}$ 、 $\hat{\psi}$ とする。

【 0 1 5 7 】

変換処理部50は、予測対象のソース音声の時系列データを入力として、ソース音声の時系列データから抽出される各時刻のスペクトル特徴量ベクトルと、第1モデルパラメータ学習部36によって学習された第1確率分布のパラメータと、第2モデルパラメータ学習部38によって学習された第2確率分布のパラメータ、に基づいて、第1確率分布と第2確率分布との積を用いて表される規準を大きくするように、各時刻の基本周波数 y と、各時刻の状態からなる状態系列 s を推定することにより、予測対象のソース音声に対応するターゲット音声の各時刻の基本周波数 y を予測する。

10

【 0 1 5 8 】

ここで、予測対象のソース音声に対応するターゲット音声の各時刻の基本周波数 y を予測する原理について説明する。

【 0 1 5 9 】

学習処理で学習したパラメータ $\hat{\theta}$ 、 $\hat{\psi}$ 、 $\hat{\gamma}$ と、予測対象のソース音声の特徴量系列 $c = \{c[k]\}_{k=1}^K$ を用いて、式(9)と式(35)の積

20

【 0 1 6 0 】

【数51】

$$\mathcal{I}_1(\mathbf{y}, \mathbf{s}) := P(\mathbf{q}|\mathbf{c}, \hat{\gamma})P(\mathbf{y}, \mathbf{s}|\hat{\theta}, \hat{\psi})^\omega \quad (77)$$

ができるだけ大きくなるように y 、 s を推定する。ただし、 $\mathbf{q} = \mathbf{W}\mathbf{y}$ であり、

【 0 1 6 1 】

【数52】

$$P(\mathbf{q}|\mathbf{c}, \hat{\gamma}) = \prod_k \sum_{m_k} P(m_k|\mathbf{c}[k], \hat{\gamma})P(\mathbf{q}[k]|\mathbf{c}[k], m_k, \hat{\gamma}) \quad (78)$$

$$P(\mathbf{y}, \mathbf{s}|\hat{\theta}, \hat{\psi}) = P(\mathbf{s}|\hat{\psi})P(\mathbf{y}|\mathbf{s}, \hat{\theta}, \hat{\psi}) \quad (79)$$

$$P(\mathbf{q}[k]|\mathbf{c}[k], m_k, \hat{\gamma}) = \mathcal{N}(\mathbf{q}[k]; \mathbf{e}_{m_k}^{(q|c)}, \mathbf{D}_{m_k}^{(q|c)}) \quad (80)$$

$$P(\mathbf{y}|\mathbf{s}, \hat{\theta}, \hat{\psi}) = \mathcal{N}(\mathbf{y}; \mathbf{G}_p\boldsymbol{\mu}_p + \mathbf{G}_a\boldsymbol{\mu}_a + \boldsymbol{\mu}_b\mathbf{1}, \mathbf{G}_p\boldsymbol{\Upsilon}_p\mathbf{G}_p^\top + \mathbf{G}_a\boldsymbol{\Upsilon}_a\mathbf{G}_a^\top + \boldsymbol{\Upsilon}_n) \quad (81)$$

【 0 1 6 2 】

である。はF₀パターンの予測においてF₀パターン生成過程のモデルをどれだけ考慮に入れるかを意味した非負の定数である。

40

【 0 1 6 3 】

以下に、

【 0 1 6 4 】

【数53】

$$\mathcal{I}_1(\mathbf{y}, \mathbf{s})$$

【 0 1 6 5 】

を大きくするためのアルゴリズムについて述べる。上述した関連技術1のスペクトル特徴

50

量系列からの F_0 パターン予測方法と同様、

【 0 1 6 6 】

【数 5 4】

$$P(\mathbf{q}[k]|\mathbf{c}[k], \hat{\gamma}) \simeq P(\hat{m}_k|\mathbf{c}[k], \hat{\gamma})P(\mathbf{q}[k]|\mathbf{c}[k], \hat{m}_k, \hat{\gamma}) \quad (86)$$

$$\hat{m}_k = \underset{m_k}{\operatorname{argmax}} P(m_k|\mathbf{c}[k], \hat{\gamma}) \quad (87)$$

【 0 1 6 7 】

と近似することで以下の反復処理により y 、 s を推定することができる（ステップ 1 と 2 の実行順序は任意）。

10

【 0 1 6 8 】

（ステップ 1）上述した関連技術 1 のスペクトル特徴量系列からの F_0 パターン予測方法の変換処理により y を初期設定する。

【 0 1 6 9 】

（ステップ 2） c を用いて \hat{m} を式 (87) により求める。

【 0 1 7 0 】

（ステップ 3） y を固定し、上述した関連技術 3 の F_0 パターン生成過程モデルパラメータ推定法のパラメータ推定アルゴリズム 1 により s を推定する。

【 0 1 7 1 】

（ステップ 4） s と \hat{m} を固定して以下の式により y を更新し、ステップ 3 に戻る。

20

【 0 1 7 2 】

【数 5 5】

$$y \leftarrow (W^T D^{(q|c)-1} W + \omega(GY G^T)^{-1})^{-1} (W^T D^{(q|c)-1} e^{(q|c)} + \omega(GY G^T)^{-1} G\mu) \quad (88)$$

ただし、

$$G = [G_p \quad G_a \quad I] \quad (89)$$

$$\mu = \begin{bmatrix} \mu_p \\ \mu_a \\ \mu_b \mathbf{1} \end{bmatrix} \quad (90)$$

30

$$Y = \begin{bmatrix} Y_p & O & O \\ O & Y_a & O \\ O & O & Y_n \end{bmatrix} \quad (91)$$

【 0 1 7 3 】

$(G \quad G^T)^{-1}$ は大きなサイズの行列の逆行列であるが、以下に示すやり方で効率的に計算することができる。 $G \quad G^T$ は

【 0 1 7 4 】

【数 5 6】

$$GY G^T = G_p Y_p G_p^T + G_a Y_a G_a^T + Y_n \quad (92)$$

40

【 0 1 7 5 】

であること、 G^{-1}_p と G^{-1}_a が

【 0 1 7 6 】

【数57】

$$G_p^{-1} = \begin{bmatrix} f_0^{(p)} & & & & O \\ f_1^{(p)} & f_0^{(p)} & & & \\ f_2^{(p)} & f_1^{(p)} & f_0^{(p)} & & \\ & \ddots & \ddots & \ddots & \\ O & & f_2^{(p)} & f_1^{(p)} & f_0^{(p)} \end{bmatrix} =: F_p \quad (93)$$

$$G_a^{-1} = \begin{bmatrix} f_0^{(a)} & & & & O \\ f_1^{(a)} & f_0^{(a)} & & & \\ f_2^{(a)} & f_1^{(a)} & f_0^{(a)} & & \\ & \ddots & \ddots & \ddots & \\ O & & f_2^{(a)} & f_1^{(a)} & f_0^{(a)} \end{bmatrix} =: F_a \quad (94)$$

10

【0177】

のような下三角帯行列で近似できることより、 $(G \ G^T)^{-1}$ は

【0178】

20

【数58】

$$(GYG^T)^{-1} = (F_p^{-1}\Upsilon_p F_p^{-T} + F_a^{-1}\Upsilon_a F_a^{-T} + \Upsilon_n)^{-1} \quad (95)$$

【0179】

と書け、さらにWoodburyの公式

【0180】

【数59】

$$(A+B)^{-1} = A^{-1} - A^{-1}(B^{-1} + A^{-1})^{-1}A^{-1} \quad (96)$$

30

【0181】

を式(95)右辺に適用することで $(G \ G^T)^{-1}$ は

【0182】

【数60】

$$(GYG^T)^{-1} = \Upsilon_n^{-1} - \Upsilon_n^{-1}((F_p^{-1}\Upsilon_p F_p^{-T} + F_a^{-1}\Upsilon_a F_a^{-T})^{-1} + \Upsilon_n^{-1})^{-1}\Upsilon_n^{-1} \quad (97)$$

【0183】

と書ける。さらにWoodburyの公式より

【0184】

【数61】

40

$$(F_p^{-1}\Upsilon_p F_p^{-T} + F_a^{-1}\Upsilon_a F_a^{-T})^{-1}$$

【0185】

は

【0186】

【数 6 2】

$$(F_p^{-1} \Upsilon_p F_p^{-T} + F_a^{-1} \Upsilon_a F_a^{-T})^{-1} = F_p^T \Upsilon_p^{-1} F_p - F_p^T \Upsilon_p^{-1} F_p (F_p^T \Upsilon_p^{-1} F_p + F_a^T \Upsilon_a^{-1} F_a)^{-1} F_p^T \Upsilon_p^{-1} F_p \quad (98)$$

【0 1 8 7】

と書ける。式(93)、(94)より

【0 1 8 8】

【数 6 3】

$$F_p^T \Upsilon_p^{-1} F_p \quad \text{と} \quad F_a^T \Upsilon_a^{-1} F_a$$

10

【0 1 8 9】

はいずれも帯行列になるので

【0 1 9 0】

【数 6 4】

$$(F_p^T \Upsilon_p^{-1} F_p + F_a^T \Upsilon_a^{-1} F_a)^{-1} a \quad \text{または} \quad (F_p^T \Upsilon_p^{-1} F_p + F_a^T \Upsilon_a^{-1} F_a)^{-1} A$$

【0 1 9 1】

の形の計算はCholesky 分解により効率的に計算することができる。ただし、a は任意のベクトル、A は任意の行列である。

20

【0 1 9 2】

以上説明した原理を実現するために、本実施の形態では、図 7 に示すように、変換処理部 5 0 は、特徴量抽出部 5 2 と、基本周波数系列予測部 5 4 と、正規分布系列予測部 5 6 と、状態系列推定部 5 8 と、基本周波数系列更新部 6 0 と、収束判定部 6 2 とを備えている。

【0 1 9 3】

特徴量抽出部 5 2 は、入力部 1 0 によって受け付けた予測対象のソース音声の時系列データから、特徴量抽出部 3 2 と同様に、ソース音声の各時刻 k のスペクトラム特徴量ベクトル c[k] を抽出する。

【0 1 9 4】

基本周波数系列予測部 5 4 は、第 1 モデルパラメータ学習部 3 6 によって学習された第 1 確率分布のパラメータ と、特徴量抽出部 5 2 によって抽出された各時刻 k のスペクトル特徴量ベクトル c[k] とに基づいて、上述した F₀ パターン予測方法の変換処理と同様に、上記式 (16) に従って、各時刻 k の基本周波数 y [k] を推定することにより、各時刻 k の基本周波数 y [k] を初期設定する。

30

【0 1 9 5】

正規分布系列予測部 5 6 は、第 1 モデルパラメータ学習部 3 6 によって学習された第 1 確率分布のパラメータ と、特徴量抽出部 5 2 によって抽出された各時刻 k のスペクトル特徴量ベクトル c[k] とに基づいて、上記式 (87) に従って、各時刻 k のスペクトル特徴量ベクトル c[k] を生成したらしい確率が最も高い正規分布のインデックス \hat{m}_k を推定する。

40

【0 1 9 6】

状態系列推定部 5 8 は、基本周波数系列予測部 5 4 によって初期設定された、または状態系列推定部 5 8 によって前回更新された各時刻 k の基本周波数 y [k] を固定して、上述した関連技術 3 の F₀ パターン生成過程モデルパラメータ推定法のパラメータ推定アルゴリズム 1 と同様に、上記式 (35) を局所最大化する状態系列 s と各時刻 k における状態に応じたフレーズ指令の振幅及び各アクセント指令の振幅を表すパラメータ とを、EM アルゴリズムにより探索することにより、状態系列 s を推定する。

【0 1 9 7】

基本周波数系列更新部 6 0 は、状態系列推定部 5 8 によって推定された状態系列 s と、

50

正規分布系列予測部 56 によって推定された各時刻の正規分布のインデックス m_k とに基づいて、上記式 (88) に従って、各時刻 k の基本周波数 $y[k]$ を更新する。

【0198】

収束判定部 62 は、予め定められた収束判定条件を満たすまで、状態系列推定部 58 及び基本周波数系列更新部 60 による各処理を繰り返させる。収束判定条件としては、例えば、予め定められた繰り返し回数に到達することである。

【0199】

収束判定条件を満たしたときに、最終的に得られた各時刻 k の基本周波数 $y[k]$ を、予測対象のソース音声に対応するターゲット音声の各時刻の基本周波数の予測結果として、出力部 90 により出力する。

10

【0200】

<基本周波数パターン予測装置の作用>

次に、本実施の形態に係る基本周波数パターン予測装置 100 の作用について説明する。まず、学習サンプルのソース音声の時系列データ及びターゲット音声の時系列データからなるパラレルデータが、基本周波数パターン予測装置 100 に入力されると、基本周波数パターン予測装置 100 において、図 8 に示す学習処理ルーチンが実行される。

【0201】

まず、ステップ S101 において、入力されたソース音声の時系列データを読み込み、

【0202】

各時刻 k のスペクトル特徴量ベクトル $c[k]$ を抽出する。ステップ S102 において、入力されたターゲット音声の時系列データを読み込み、ターゲット音声の各時刻 k における基本周波数 $y[k]$ を抽出し、また、基本周波数 $y[k]$ とその動的成分の結合ベクトル $q[k]$ を抽出する。

20

【0203】

そして、ステップ S103 において、上記ステップ S101 で抽出された各時刻 k のスペクトル特徴量ベクトル $c[k]$ と、上記ステップ S102 で抽出された各時刻 k の基本周波数の結合ベクトル $q[k]$ とに基づいて、上記式 (1) の GMM のパラメータ を学習する。

【0204】

ステップ S104 では、状態系列 s と、各時刻における状態に応じたフレーズ指令の振幅及び各アクセント指令の振幅を表すパラメータ とを初期設定する。

30

【0205】

そして、ステップ S105 において、上記式 (49)、式 (50) に従って、フレーズ成分、アクセント成分、ベースライン成分の条件付き期待値 x と、条件付き分散共分散行列 R とを更新する。

【0206】

次のステップ S106 では、上記ステップ S104 で初期設定された、又は後述するステップ S107 で前回更新されたパラメータ と、上記ステップ S105 で更新されたフレーズ成分、アクセント成分、ベースライン成分の条件付き期待値 x と、条件付き分散共分散行列 R とに基づいて、上記式 (63) 式を用いて、 $Q(s, \cdot)$ を最大にする状態系列 $s = (s_1, \dots, s_k)$ を Viterbi アルゴリズムにより求めて、状態系列 s を更新する。

40

【0207】

ステップ S107 では、上記ステップ S106 で更新された状態系列 s と、上記ステップ S105 で更新されたフレーズ成分、アクセント成分、ベースライン成分の条件付き期待値 x と、条件付き分散共分散行列 R とに基づいて、 $Q(s, \cdot)$ の各変数に関する偏微分を 0 と置くことにより、 $Q(s, \cdot)$ を最大にするパラメータ を求めて、各時刻における状態に応じたフレーズ指令の振幅及び各アクセント指令の振幅を表すパラメータ とを更新する。

【0208】

ステップ S108 において、予め定められた収束判定条件を満たしたか否かを判定し、収束判定条件を満たしていない場合には、上記ステップ S105 へ戻る。一方、収束判定

50

条件を満たした場合には、ステップS 109において、上記ステップS 103で学習されたパラメータ、上記ステップS 107で最終的に得られたパラメータとを、パラメータ記憶部40に格納する。また、上記ステップS 106で最終的に得られた状態系列sから、状態遷移確率を求め、パラメータ記憶部40に格納する。

【0209】

次に、予測対象のソース音声の時系列データが、基本周波数パターン予測装置100に入力されると、基本周波数パターン予測装置100において、図9に示す基本周波数パターン予測処理ルーチンが実行される。

【0210】

まず、ステップS 121において、入力された予測対象のソース音声の時系列データを読み込み、各時刻kのスペクトル特徴量ベクトル $c[k]$ を抽出する。ステップS 122において、パラメータ記憶部40に記憶されたパラメータと、上記ステップS 121で抽出された各時刻kのスペクトル特徴量ベクトル $c[k]$ とに基づいて、上記式(16)に従って、各時刻kの基本周波数 $y[k]$ を推定することにより、各時刻kの基本周波数 $y[k]$ を初期設定する。

10

【0211】

そして、ステップS 123では、パラメータ記憶部40に記憶されたパラメータと、上記ステップS 121で抽出された各時刻kのスペクトル特徴量ベクトル $c[k]$ とに基づいて、上記式(87)に従って、各時刻kのスペクトル特徴量ベクトル $c[k]$ を生成したらしい確率が最も高い正規分布のインデックス m_k を推定する。

20

【0212】

ステップS 124では、上記ステップS 122で初期設定された、または後述するステップS 125で前回更新された各時刻kの基本周波数 $y[k]$ を固定して、上記式(35)を局所最大化する状態系列sと各時刻kにおける状態に応じたフレーズ指令の振幅及び各アクセント指令の振幅を表すパラメータとを、EMアルゴリズムにより探索することにより、状態系列sを推定する。

【0213】

ステップS 125では、上記ステップS 124で推定された状態系列sと、上記ステップS 123で推定された各時刻の正規分布のインデックス m_k とに基づいて、上記式(88)に従って、各時刻kの基本周波数 $y[k]$ を更新する。

30

【0214】

ステップS 126では、予め定められた収束判定条件を満たしたか否かを判定し、収束判定条件を満たしていない場合には、ステップS 124へ戻る。一方、収束判定条件を満たした場合には、ステップS 127において、上記ステップS 125で最終的に得られた各時刻kの基本周波数 $y[k]$ を、予測対象のソース音声に対応するターゲット音声の各時刻の基本周波数の予測結果として、出力部90により出力し、基本周波数パターン予測処理ルーチンを終了する。

【0215】

以上説明したように、第1の実施の形態に係る基本周波数パターン予測装置によれば、ソース音声の各時刻のスペクトル特徴量ベクトルと、ターゲット音声の各時刻の基本周波数との間の関係をモデル化したGMMである第1確率分布 $P(q[k], c[k] | \theta_1)$ のパラメータを学習し、基本周波数パターン生成過程をモデル化した第2確率分布 $P(y, s | \theta_2)$ のパラメータを学習し、予測対象のソース音声の時系列データから抽出される各時刻のスペクトル特徴量ベクトルから、第1確率分布 $P(q[k], c[k] | \theta_1)$ と第2確率分布 $P(y, s | \theta_2)$ との積を用いて表される規準を大きくするように、予測対象のソース音声に対応するターゲット音声の各時刻の基本周波数を予測することにより、 F_0 パターンの物理的な生成過程の制約を考慮しながらスペクトル特徴量系列に対応する最適な F_0 パターンを推定することができる。

40

【0216】

[第2の実施の形態]

50

次に、本発明の第2の実施の形態に係る基本周波数パターン予測装置について説明する。なお、第1の実施の形態と同様の構成となる部分については同一符号を付して説明を省略する。

【0217】

第2の実施の形態では、第2の確率分布及びパラメータを推定する方法と、各時刻の基本周波数の予測方法とが第1の実施の形態と異なっている。

【0218】

第2の実施の形態に係る基本周波数パターン予測装置の学習部30の第2モデルパラメータ学習部38によるパラメータを学習する原理について説明する。

【0219】

まず、関連技術3の F_0 パターン生成過程モデルパラメータ推定法のパラメータ推定アルゴリズムについて説明する。

【0220】

<パラメータ推定アルゴリズム2>

観測 F_0 系列 y が与えられたもとで、モデルパラメータ θ と σ の事後確率 $P(\theta, \sigma | y)$ の局所最適解を求める反復アルゴリズムを以下に示す。状態系列 s を隠れ変数とし、事後確率 $P(\theta, \sigma | y)$ が、

【0221】

【数65】

$$P(\theta, \sigma, s | y) \propto P(y | \theta) P(\theta | s, \sigma) P(s) \quad 20$$

【0222】

を s について周辺化することで得られる点に注意すると、Q関数 $Q(\theta, \sigma, \theta', \sigma')$ は

【0223】

【数66】

$$\begin{aligned} Q(\theta, \sigma, \theta', \sigma') &= \sum_s P(s | y, \theta', \sigma') \log P(\theta, \sigma, s | y) \\ &\stackrel{c}{=} \log P(y | \theta) + \sum_s P(s | y, \theta', \sigma') \log P(\theta | s, \sigma) P(s) \end{aligned} \quad 30$$

$$\log P(y | \theta) = \log \prod_{k=1}^K \mathcal{N}(y[k]; x[k], v_n^2[k]) = - \sum_{k=1}^K \frac{C[k]}{2v_n^2[k]} \quad (64)$$

$$C[k] = \left(y[k] - \sum_{i \in \{p, a, b\}} \sum_l g_i[k-l+1] u_i[l] \right)^2 \quad (65)$$

$$\sum_s P(s | y, \theta', \sigma') \log P(\theta, \sigma | s) = - \sum_{k=1}^K \sum_t P(s_k = t | y, \theta', \sigma') \sum_{i \in \{p, a\}} \frac{D_t^{(i)}[k]}{2v_t^{(i)2}} \quad (66) \quad 40$$

$$D_t^{(i)}[k] = \left(u_i[k] - \mu_t^{(i)}[k] \right)^2 \quad (67)$$

と置ける。ただし、

【0224】

【数67】

$$\underline{\underline{c}}$$

10

20

30

40

50

【0225】

は定数項を除いて等しいことを表す。また、 $g_b[k] = \delta[k]$ (クロネッカーのデルタ) である。よって、 $P(s_k = t | y, \theta, \theta')$ を Forward-Backward アルゴリズムにより計算するステップ、 θ と θ' について $Q(\theta, \theta, \theta', \theta')$ を増加させるステップを繰り返すことで、 $P(\theta, \theta | y)$ が局所最大となる解を得ることができる。 θ はフレーズ・アクセント指令系列のペアであるため、 $Q(\theta, \theta, \theta', \theta')$ を増加させるステップにおいては、 θ の非負制約を考慮する必要がある。 θ の非負制約を満たしながら $Q(\theta, \theta, \theta', \theta')$ を増加させるような更新則は以下により導くことができる。まず、 $Q(\theta, \theta, \theta', \theta')$ の下界は Jensen の不等式より

【0226】

【数68】

$$-\left(\sum_{i \in \{p,a,b\}} \sum_l g_i[k-l+1]u_i[l]\right)^2 \geq -\sum_{i \in \{p,a,b\}} \sum_l \frac{g_i^2[k-l+1]u_i^2[l]}{\lambda_{i,k,l}}, \quad (68)$$

【0227】

のように設計することができる。また、 i, k, l は、

【0228】

【数69】

$$0 < \lambda_{i,k,l} < 1, \sum_i \sum_l \lambda_{i,k,l} = 1$$

【0229】

を満たす任意の変数である。従って Q 関数の下界は、

【0230】

【数70】

$$Q(\theta, \theta, \theta', \theta') \geq -\sum_{k=1}^K \frac{\tilde{C}[k]}{2v_n^2[k]} - \sum_{k=1}^K \sum_t P(s_k = t | y, \theta', \theta') \sum_{i \in \{p,a\}} \frac{D_t^{(i)}[k]}{2v_t^{(i)2}} \quad (69)$$

$$\tilde{C}[k] = y[k]^2 - 2y[k] \sum_{i \in \{p,a,b\}} \sum_l g_i[k-l+1]u_i[l] + \sum_{i \in \{p,a,b\}} \sum_l \frac{g_i^2[k-l+1]u_i^2[l]}{\lambda_{i,k,l}} \quad (70)$$

【0231】

と表される。この下界関数を $\lambda_{i,k,l} = 0$ に関して最大化するステップと θ に関して最大化するステップを交互に繰り返せば $Q(\theta, \theta, \theta', \theta')$ を増加させることができる。いずれのステップの更新則も解析的に求めることができ、それぞれ

【0232】

【数71】

$$\lambda_{i,k,l} = \frac{g_i[k-l+1]u_i[l]}{\sum_{i \in \{p,a,b\}} \sum_{l'} g_i[k-l'+1]u_i[l']} \quad (71)$$

$$u_i[l] = \frac{\sum_t \frac{P(s_l = t | y, \theta', \theta') \mu_t^{(i)}[l]}{v_t^{(i)2}} + \sum_k \frac{y[k]g_i[k-l+1]}{v_n^2[k]}}{\sum_t \frac{P(s_l = t | y, \theta', \theta')}{v_t^{(i)2}} + \sum_k \frac{g_i^2[k-l+1]}{v_n^2[k]\lambda_{i,k,l}}} \quad (72)$$

10

20

30

40

50

【 0 2 3 3 】

で表される。以上の反復が収束したあと、続けて を更新する。更新式は、図 1 の場合、

【 0 2 3 4 】

【 数 7 2 】

$$A^{(p)}[k] = \frac{\sum_{t \in p_0} \frac{P(s_k = t | \mathbf{y}, \mathbf{o}', \boldsymbol{\theta}') u_p[k]}{v_t^{(p)2}}}{\sum_{t \in p_0} \frac{P(s_k = t | \mathbf{y}, \mathbf{o}', \boldsymbol{\theta}')}{v_t^{(p)2}}} \quad (73)$$

10

$$A_n^{(a)} = \frac{\sum_k \sum_{t \in a_n} \frac{P(s_k = t | \mathbf{y}, \mathbf{o}', \boldsymbol{\theta}') u_a[k]}{v_t^{(a)2}}}{\sum_k \sum_{t \in a_n} \frac{P(s_k = t | \mathbf{y}, \mathbf{o}', \boldsymbol{\theta}')}{v_t^{(a)2}}} \quad (74)$$

【 0 2 3 5 】

図 2、3 の場合、

20

【 0 2 3 6 】

【 数 7 3 】

$$A_m^{(p)} = \frac{\sum_k \sum_{t \in p_m} \frac{P(s_k = t | \mathbf{y}, \mathbf{o}', \boldsymbol{\theta}') u_p[k]}{v_t^{(p)2}}}{\sum_k \sum_{t \in p_m} \frac{P(s_k = t | \mathbf{y}, \mathbf{o}', \boldsymbol{\theta}')}{v_t^{(p)2}}} \quad (75)$$

30

$$A_n^{(a)} = \frac{\sum_k \sum_{t \in a_n} \frac{P(s_k = t | \mathbf{y}, \mathbf{o}', \boldsymbol{\theta}') u_a[k]}{v_t^{(a)2}}}{\sum_k \sum_{t \in a_n} \frac{P(s_k = t | \mathbf{y}, \mathbf{o}', \boldsymbol{\theta}')}{v_t^{(a)2}}} \quad (76)$$

【 0 2 3 7 】

である。これらの更新値を \hat{o}' と $\hat{\theta}'$ に代入したのちに、 $P(s_k = t | \mathbf{y}, \hat{o}', \hat{\theta}')$ の更新を再度行い、以後同様の処理を繰り返すことで事後確率 $P(\hat{o}, \hat{\theta} | \mathbf{y})$ を単調増加させることができる。

40

【 0 2 3 8 】

以上の反復アルゴリズムが収束した後、上述したパラメータ推定アルゴリズム 1 の Viterbi アルゴリズムにより求まる最適な s を状態系列の推定結果とする。

【 0 2 3 9 】

また、推定された状態系列 s から、状態遷移確率 が求められる。

【 0 2 4 0 】

以上説明した原理に従って、第 2 モデルパラメータ学習部 3 8 は、基本周波数系列抽出部 3 4 によって抽出された各時刻 k の基本周波数 $y[k]$ に基づいて、各時刻 k の基本周波数 $y[k]$ と、各時刻 k における甲状軟骨の平行移動運動によって生じる基本周波数パターンを表すフレーズ指令 $u_p[k]$ 及び甲状軟骨の回転運動によって生じる基本周波数パターン

50

を表すアクセント指令 $u_a[k]$ のペアからなる指令関数 $o[k]$ との組み合わせの確率分布である第2確率分布のパラメータ、を学習する。

【0241】

第2の実施の形態における変換処理部50は、予測対象のソース音声の時系列データを入力として、ソース音声の時系列データから抽出される各時刻 k のスペクトル特徴量ベクトル $c[k]$ と、第1モデルパラメータ学習部36によって学習された第1確率分布のパラメータと、第2モデルパラメータ学習部38によって学習された第2確率分布のパラメータとに基づいて、第1確率分布と第2確率分布とを用いて表される規準を大きくするように、各時刻 k の基本周波数 $y[k]$ と、各時刻 k のフレーズ指令及びアクセント指令のペアからなる指令関数 $o[k]$ とを推定することにより、予測対象のソース音声に対応するターゲット音声の各時刻 k の基本周波数 $y[k]$ を予測する。

10

【0242】

ここで、予測対象のソース音声に対応するターゲット音声の各時刻 k の基本周波数 $y[k]$ を予測する原理について説明する。

【0243】

<変換処理>

学習処理で学習したパラメータ $\hat{\theta}$ 、 $\hat{\psi}$ 、 $\hat{\gamma}$ と、予測対象のソース音声の特徴量系列 $c = \{c[k]\}_{k=1}^K$ を用いて、式(9)と式(45)の積

【0244】

【数74】

20

$$\mathcal{I}_2(\mathbf{y}, \mathbf{o}) := P(\mathbf{q}|\mathbf{c}, \hat{\gamma})P(\mathbf{y}, \mathbf{o}|\hat{\theta}, \hat{\psi})^\omega \quad (82)$$

【0245】

ができるだけ大きくなるように y 、 o を推定する。ただし、 $\mathbf{q} = \mathbf{W}\mathbf{y}$ であり、

【0246】

【数75】

$$P(\mathbf{y}, \mathbf{o}|\hat{\theta}, \hat{\psi}) = P(\mathbf{y}|\mathbf{o}) \sum_{\mathbf{s}} P(\mathbf{o}|\mathbf{s}, \hat{\theta})P(\mathbf{s}|\hat{\psi}) \quad (83)$$

30

$$P(\mathbf{y}|\mathbf{o}) = \prod_k \mathcal{N}(y[k]; x[k], v_n^2[k]) \quad (84)$$

$$x[k] = g_p[k] * u_p[k] + g_a[k] * u_a[k] + \mu_b \quad (85)$$

である。

【0247】

次に、

【0248】

【数76】

40

$$\mathcal{I}_2(\mathbf{y}, \mathbf{o})$$

【0249】

を大きくするためのアルゴリズムについて述べる。上述した関連技術1のスペクトル特徴量系列からの F_0 パターン予測方法と同様に、

【0250】

【数 7 7】

$$P(\mathbf{q}[k]|\mathbf{c}[k], \hat{\gamma}) \simeq P(\hat{m}_k|\mathbf{c}[k], \hat{\gamma})P(\mathbf{q}[k]|\mathbf{c}[k], \hat{m}_k, \hat{\gamma}) \quad (99)$$

$$\hat{m}_k = \underset{m_k}{\operatorname{argmax}} P(m_k|\mathbf{c}[k], \hat{\gamma}) \quad (100)$$

【 0 2 5 1】

と近似することで以下の反復処理により y 、 o を推定することができる（ステップ 1 と 2 の実行順序は任意）。

【 0 2 5 2】

（ステップ 1）上述した関連技術 1 のスペクトル特徴量系列からの F_0 パターン予測方法の変換処理により y を初期設定する。 10

【 0 2 5 3】

（ステップ 2） c を用いて \hat{m} を式 (100) により求める。

【 0 2 5 4】

（ステップ 3） y を固定し、上述した関連技術 3 の F_0 パターン生成過程モデルパラメータ推定法のパラメータ推定アルゴリズム 2 により o を推定する。

【 0 2 5 5】

（ステップ 4） o と \hat{m} を固定して以下の式により y を更新し、ステップ 3 に戻る。

【 0 2 5 6】

【数 7 8】

$$\mathbf{y} \leftarrow (\mathbf{W}^T \mathbf{D}^{(q|c)-1} \mathbf{W} + \omega \Sigma_n^{-1})^{-1} (\mathbf{W}^T \mathbf{D}^{(q|c)-1} \mathbf{e}^{(q|c)} + \omega \Sigma_n^{-1} \mathbf{G} \mathbf{u}) \quad (101)$$

ただし、

$$\mathbf{G} = [\mathbf{G}_p \quad \mathbf{G}_a \quad \mathbf{I}] \quad (102)$$

$$\mathbf{u} = \begin{bmatrix} \mathbf{u}_p \\ \mathbf{u}_a \\ \mathbf{u}_b \end{bmatrix} \quad (103)$$

$$\mathbf{u}_p = (u_p[1], \dots, u_p[K])^T \quad (104)$$

$$\mathbf{u}_a = (u_a[1], \dots, u_a[K])^T \quad (105)$$

$$\mathbf{u}_b = (u_b[1], \dots, u_b[K])^T \quad (106)$$

【 0 2 5 7】

以上説明した原理を実現するために、第 2 の実施の形態では、図 10 に示すように、変換処理部 50 は、特徴量抽出部 52 と、基本周波数系列予測部 54 と、正規分布系列予測部 256 と、指令系列推定部 258 と、基本周波数系列更新部 260 と、収束判定部 62 とを備えている。

【 0 2 5 8】

正規分布系列予測部 256 は、第 1 モデルパラメータ学習部 36 によって学習された第 1 確率分布のパラメータ と、特徴量抽出部 52 によって抽出された各時刻のスペクトル特徴量ベクトル $\mathbf{c}[k]$ とに基づいて、上記式 (100) に従って、各時刻 k のスペクトル特徴量ベクトル $\mathbf{c}[k]$ を生成したらしい確率が最も高い正規分布のインデックス \hat{m}_k を推定する。 40

【 0 2 5 9】

指令系列推定部 258 は、基本周波数系列予測部 54 によって初期設定された、または状態系列推定部 58 によって前回更新された各時刻 k の基本周波数 $y[k]$ を固定して、上述した関連技術 3 の F_0 パターン生成過程モデルパラメータ推定法のパラメータ推定アルゴリズム 2 と同様に、事後確率 $P(o, |y)$ を局所最大化する指令系列 o を推定する。 50

【 0 2 6 0 】

基本周波数系列更新部 2 6 0 は、指令系列推定部 2 5 8 によって推定された指令系列 o と、正規分布系列予測部 5 6 によって推定された各時刻の正規分布のインデックス $\wedge mk$ とに基づいて、上記式 (1 0 1) に従って、各時刻 k の基本周波数 $y [k]$ を更新する。

【 0 2 6 1 】

収束判定部 6 2 は、予め定められた収束判定条件を満たすまで、指令系列推定部 2 5 8 及び基本周波数系列更新部 2 6 0 による各処理を繰り返させる。収束判定条件としては、例えば、予め定められた繰り返し回数に到達することである。

【 0 2 6 2 】

収束判定条件を満たしたときに、最終的に得られた各時刻 k の基本周波数 $y [k]$ を、
予測対象のソース音声に対応するターゲット音声の各時刻の基本周波数の予測結果として、出力部 9 0 により出力する。

10

【 0 2 6 3 】

< 基本周波数パターン予測装置の作用 >

【 0 2 6 4 】

次に、第 2 の実施の形態に係る基本周波数パターン予測装置 1 0 0 の作用について説明する。なお、第 1 の実施の形態と同様の処理については、同一符号を付して詳細な説明を省略する。

【 0 2 6 5 】

まず、学習サンプルのソース音声の時系列データ及びターゲット音声の時系列データからなるパラレルデータが、基本周波数パターン予測装置 1 0 0 に入力されると、基本周波数パターン予測装置 1 0 0 において、図 1 1 に示す学習処理ルーチンが実行される。

20

【 0 2 6 6 】

まず、ステップ S 1 0 1 において、入力されたソース音声の時系列データを読み込み、各時刻 k のスペクトル特徴量ベクトル $c [k]$ を抽出する。ステップ S 1 0 2 において、入力されたターゲット音声の時系列データを読み込み、ターゲット音声の各時刻 k における基本周波数 $y [k]$ を抽出し、また、各時刻 k の基本周波数 $y [k]$ とその動的成分の結合ベクトル $q [k]$ を抽出する。

【 0 2 6 7 】

そして、ステップ S 1 0 3 において、上記ステップ S 1 0 1 で抽出された各時刻 k のスペクトル特徴量ベクトル $c [k]$ と、上記ステップ S 1 0 2 で抽出された各時刻 k の基本周波数の結合ベクトル $q [k]$ とに基づいて、上記式 (1) の GMM のパラメータ を学習する。

30

【 0 2 6 8 】

ステップ S 2 0 0 では、指令系列 o と、各時刻における状態に応じたフレーズ指令の振幅及び各アクセント指令の振幅を表すパラメータ とを初期設定する。また、ターゲット音声の時系列データに基づいて、有声区間、無声区間を特定し、各時刻 k の基本周波数の不確かさの程度 $v_n^2 [k]$ を推定する。

【 0 2 6 9 】

そして、ステップ S 2 0 1 において、上記ステップ S 2 0 0 で設定された指令系列 o の初期値、または後述するステップ S 2 0 3 で前回更新された指令系列 o に基づいて、 (k, t) の全ての組み合わせについて、事後確率 $P (s_k = t | y, o , \quad)$ を更新する。

40

【 0 2 7 0 】

ステップ S 2 0 2 では、上記ステップ S 2 0 0 で設定された指令系列 o の初期値、または後述するステップ S 2 0 3 で前回更新された指令系列 o に基づいて、 (k, l) の全ての組み合わせについて、上記の式 (7 1) に従って、補助変数 $p_{, k, l}$ 、 $a_{, k, l}$ 、 $b_{, k, l}$ を算出して更新する。

【 0 2 7 1 】

次のステップ S 2 0 3 では、上記ステップ S 1 0 2 で抽出された基本周波数系列 y と、上記ステップ S 2 0 0 で算出された各時刻 k の不確かさの程度 $v_n^2 [k]$ と、上記ステップ S 2 0 1 で更新された事後確率 $P (s_k = t | y, o , \quad)$ と、上記ステップ S 2 0 2 で

50

更新された補助変数 $p_{p,k,l}$ 、 $a_{a,k,l}$ 、 $b_{b,k,l}$ とに基づいて、上記式 (72) に従って、非負値である各時刻 l のフレーズ指令 $u_p[l]$ 及びアクセント指令 $u_a[l]$ からなる指令系列 o とベース成分 u_b とを更新する。

【0272】

次のステップ S204 では、収束条件として、繰り返し回数 s が、 S に到達したか否かを判定し、繰り返し回数 s が S に到達していない場合には、収束条件を満足していないと判断して、上記ステップ S202 へ戻る。一方、繰り返し回数 s が S に到達した場合には、収束条件を満足したと判断し、ステップ S205 で、上記ステップ S203 で更新された各時刻 k のフレーズ指令 $u_p[k]$ 及びアクセント指令 $u_a[k]$ と、上記ステップ S201 で更新された事後確率 $P(s_k=t|y,o)$ とに基づいて、上記式 (73)、式 (74)、又は式 (75)、式 (76) に従って、各時刻 k のフレーズ指令の振幅 $A^{(p)}[k]$ 、及び各位置 n のアクセント指令の振幅 $A_a^{(a)}$ を更新することにより、各時刻における状態に応じたフレーズ指令の振幅及び各アクセント指令の振幅を表すパラメータ を更新する。

10

【0273】

ステップ S206 において、予め定められた収束判定条件を満たしたか否かを判定し、収束判定条件を満たしていない場合には、上記ステップ S201 へ戻る。一方、収束判定条件を満たした場合には、ステップ S207 において、

【0274】

上記ステップ S203 で最終的に更新された指令系列 o に基づいて、Viterbi アルゴリズムにより、状態系列 s を推定する。また、推定された状態系列 s から、状態遷移確率 を求める。

20

【0275】

そして、ステップ S208 において、上記ステップ S103 で学習されたパラメータ、上記ステップ S205 で最終的に得られたパラメータ と、上記ステップ S106 で得られた状態遷移確率 とを、パラメータ記憶部 40 に格納する。

【0276】

次に、予測対象のソース音声の時系列データが、基本周波数パターン予測装置 100 に入力されると、基本周波数パターン予測装置 100 において、図 12 に示す基本周波数パターン予測処理ルーチンが実行される。

30

【0277】

まず、ステップ S121 において、入力された予測対象のソース音声の時系列データを読み込み、各時刻 k のスペクトル特徴量ベクトル $c[k]$ を抽出する。ステップ S122 において、パラメータ記憶部 40 に記憶されたパラメータ と、上記ステップ S121 で抽出された各時刻のスペクトル特徴量ベクトル $c[k]$ とに基づいて、各時刻の基本周波数 $y[k]$ を初期設定する。

【0278】

そして、ステップ S123 では、パラメータ記憶部 40 に記憶されたパラメータ と、上記ステップ S121 で抽出された各時刻のスペクトル特徴量ベクトル $c[k]$ とに基づいて、各時刻 k の正規分布のインデックス \hat{m}_k を推定する。

40

【0279】

ステップ S221 では、上記ステップ S122 で初期設定された、または後述するステップ S125 で前回更新された各時刻 k の基本周波数 $y[k]$ を固定して、上記ステップ S201 ~ ステップ S206 と同様に、事後確率 $P(o, |y)$ を局所最大化する指令系列 o を推定する。

【0280】

そして、ステップ S222 において、上記ステップ S221 で推定された指令系列 o と、上記ステップ S123 で推定された各時刻の正規分布のインデックス \hat{m}_k とに基づいて、上記式 (101) に従って、各時刻 k の基本周波数 $y[k]$ を更新する。

【0281】

50

ステップS 1 2 6では、予め定められた収束判定条件を満たしたか否かを判定し、収束判定条件を満たしていない場合には、ステップS 2 2 1へ戻る。一方、収束判定条件を満たした場合には、ステップS 1 2 7において、上記ステップS 2 2 2で最終的に得られた各時刻kの基本周波数 $y[k]$ を、予測対象のソース音声に対応するターゲット音声の各時刻の基本周波数の予測結果として、出力部90により出力し、基本周波数パターン予測処理ルーチンを終了する。

【0282】

以上説明したように、第2の実施の形態に係る基本周波数パターン予測装置によれば、ソース音声の各時刻のスペクトル特徴量ベクトルと、ターゲット音声の各時刻の基本周波数との間の関係をモデル化したGMMである第1確率分布 $P(q[k], c[k] | \theta_1)$ のパラメータを学習し、基本周波数パターン生成過程をモデル化した第2確率分布 $P(y, o | \theta_2)$ のパラメータを学習し、予測対象のソース音声の時系列データから抽出される各時刻のスペクトル特徴量ベクトルから、第1確率分布 $P(q[k], c[k] | \theta_1)$ と第2確率分布 $P(y, o | \theta_2)$ との積を用いて表される規準を大きくするように、予測対象のソース音声に対応するターゲット音声の各時刻の基本周波数を予測することにより、 F_0 パターンの物理的な生成過程の制約を考慮しながらスペクトル特徴量系列に対応する最適な F_0 パターンを推定することができる。

【0283】

<実験>

【0284】

図13に示す F_0 パターンの音声データに対し、上述した従来手法である関連技術1のスペクトル特徴量系列からの F_0 パターン予測方法と、本発明の第1の実施の形態に係る手法とによりスペクトル特徴量系列から F_0 パターンの予測を行う実験を行った。図14に、両手法により予測された F_0 パターンを示す。図14では、実線が、従来手法による音声特徴量系列からの F_0 パターンの予測結果の例を示し、点線が、第1の実施の形態に係る手法による音声特徴量系列からの F_0 パターンの予測結果の例を示す。

【0285】

図13の F_0 パターンとの近さが F_0 パターンの良さの指標になる。そこで、それぞれの手法で得られた F_0 パターンと、図13の F_0 パターンとのコサイン距離(1に近いほど近いことを意味する)を測ったところ、従来手法が0.55、第1の実施の形態に係る手法が0.59であった。このことから、本発明の第1の実施の形態の手法の、従来手法に対する優位性が示された。

【0286】

なお、本発明は、上述した実施形態に限定されるものではなく、この発明の要旨を逸脱しない範囲内で様々な変形や応用が可能である。

【0287】

例えば、上述の基本周波数パターン予測装置は、内部にコンピュータシステムを有しているが、「コンピュータシステム」は、WWWシステムを利用している場合であれば、ホームページ提供環境(あるいは表示環境)も含むものとする。

【0288】

また、本願明細書中において、プログラムが予めインストールされている実施形態として説明したが、当該プログラムを、コンピュータ読み取り可能な記録媒体に格納して提供することも可能である。

【符号の説明】

【0289】

- 10 入力部
- 20 演算部
- 30 学習部
- 32 特徴量抽出部
- 34 基本周波数系列抽出部

10

20

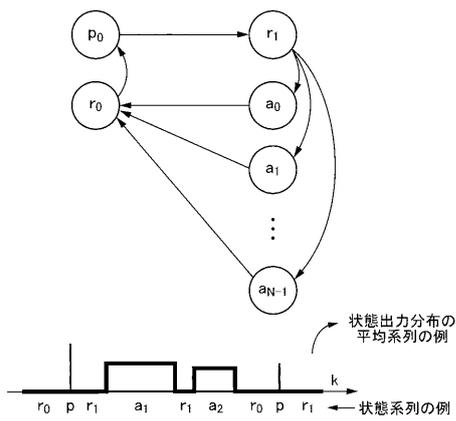
30

40

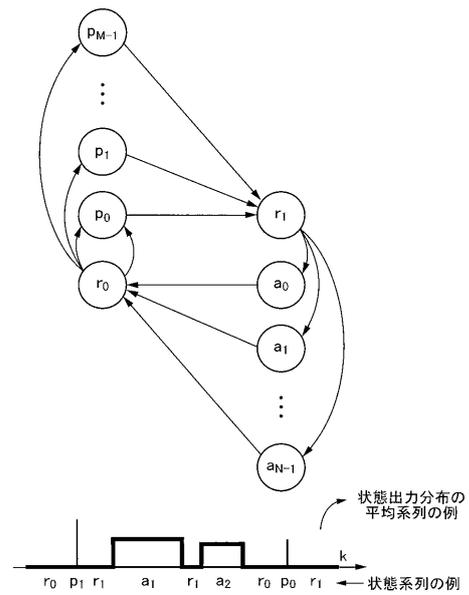
50

- 3 6 第 1 モデルパラメータ学習部
- 3 8 第 2 モデルパラメータ学習部
- 4 0 パラメータ記憶部
- 5 0 変換処理部
- 5 2 特徴量抽出部
- 5 4 基本周波数系列予測部
- 5 6 正規分布系列予測部
- 5 8 状態系列推定部
- 6 0 基本周波数系列更新部
- 6 2 収束判定部
- 9 0 出力部
- 1 0 0 基本周波数パターン予測装置
- 2 5 6 正規分布系列予測部
- 2 5 8 指令系列推定部
- 2 6 0 基本周波数系列更新部

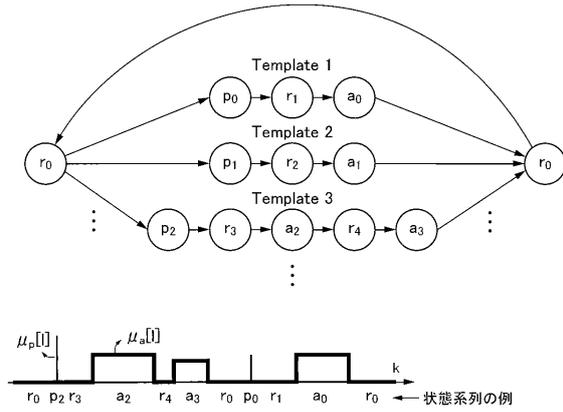
【 図 1 】



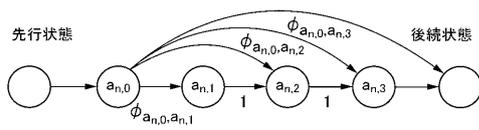
【 図 2 】



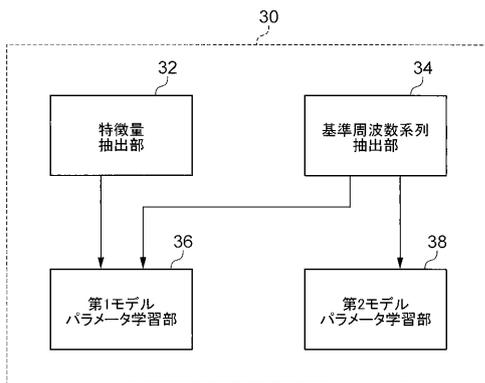
【図3】



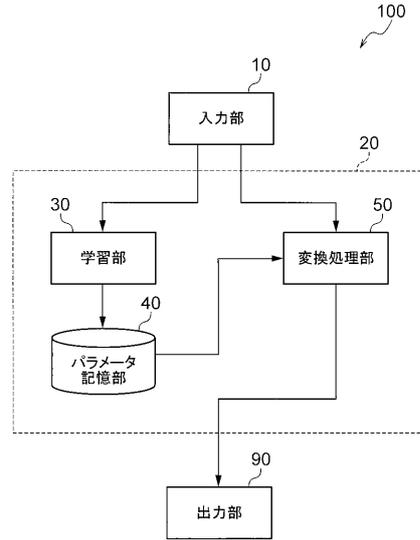
【図4】



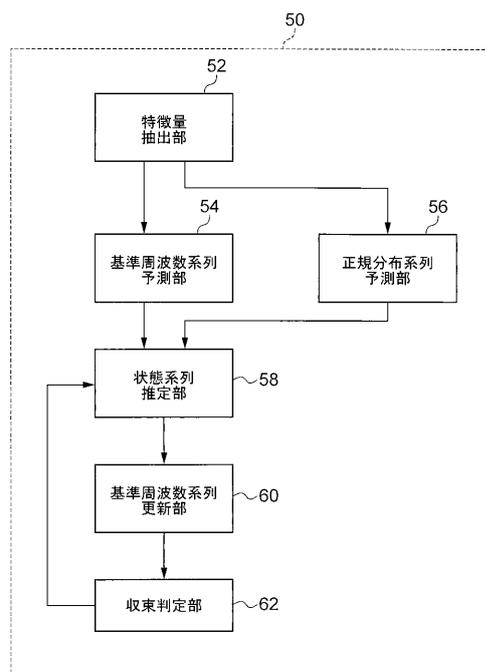
【図6】



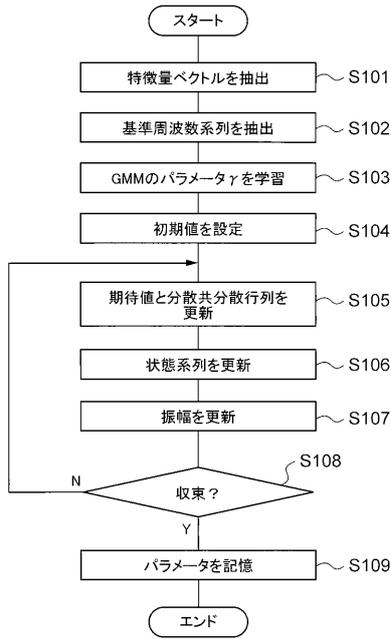
【図5】



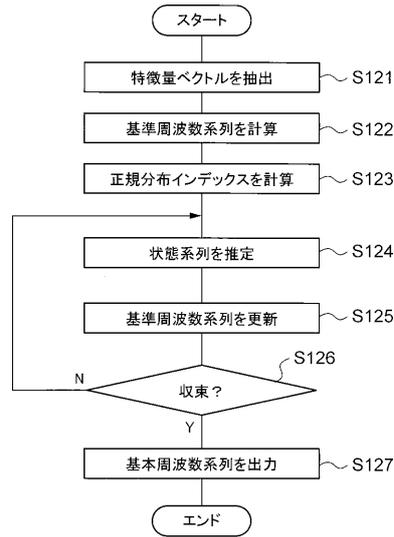
【図7】



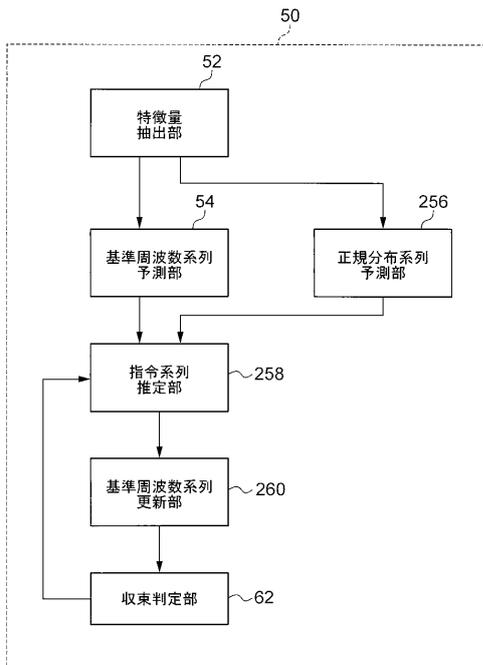
【図8】



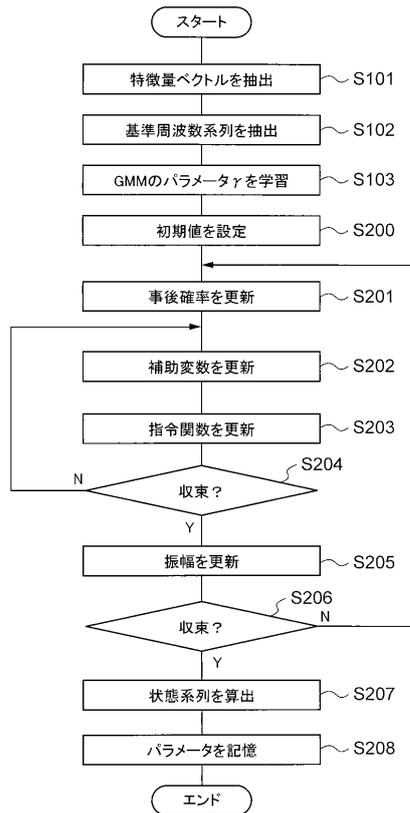
【図9】



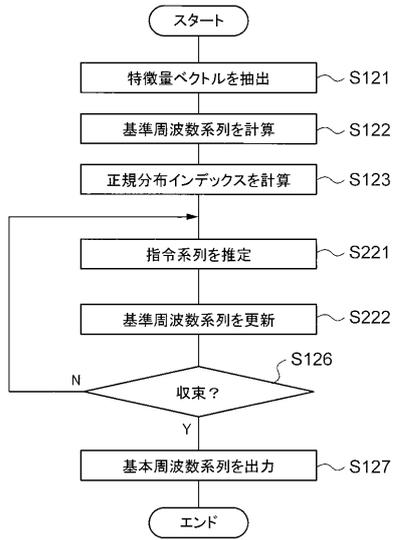
【図10】



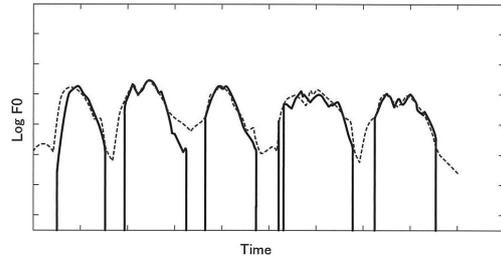
【図11】



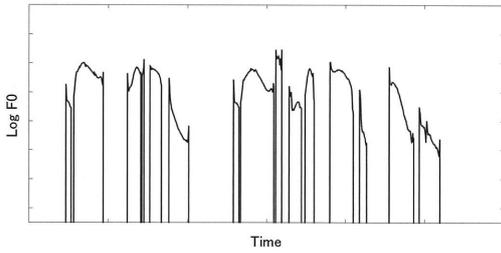
【図 1 2】



【図 1 4】



【図 1 3】



フロントページの続き

(72)発明者 戸田 智基

奈良県生駒市高山町8916-5 国立大学法人奈良先端科学技術大学院大学内

(72)発明者 中村 哲

奈良県生駒市高山町8916-5 国立大学法人奈良先端科学技術大学院大学内

審査官 上田 雄

(56)参考文献 特開2015-041081(JP,A)

国際公開第2010/137385(WO,A1)

特開2013-171196(JP,A)

特開2015-041004(JP,A)

特開2014-134730(JP,A)

国際公開第2013/132959(WO,A1)

(58)調査した分野(Int.Cl., DB名)

G10L 25/90

G10L 13/00-13/10