

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号

特許第3936266号
(P3936266)

(45) 発行日 平成19年6月27日(2007.6.27)

(24) 登録日 平成19年3月30日(2007.3.30)

(51) Int. Cl. F I
G 1 O L 15/10 (2006.01) G 1 O L 15/10 3 O O Z
G 1 O L 15/14 (2006.01) G 1 O L 15/14 2 O O A

請求項の数 10 外国語出願 (全 13 頁)

(21) 出願番号	特願2002-265510 (P2002-265510)	(73) 特許権者	393031586
(22) 出願日	平成14年9月11日(2002.9.11)		株式会社国際電気通信基礎技術研究所
(65) 公開番号	特開2004-163445 (P2004-163445A)	(74) 代理人	100099933
(43) 公開日	平成16年6月10日(2004.6.10)		弁理士 清水 敏
審査請求日	平成16年7月7日(2004.7.7)	(72) 発明者	マルコフ コンスタンチン
			京都府相楽郡精華町光台二丁目2番地2
		(72) 発明者	株式会社国際電気通信基礎技術研究所内
			中村 哲
			京都府相楽郡精華町光台二丁目2番地2
			株式会社国際電気通信基礎技術研究所内
		審査官	櫻本 剛

最終頁に続く

(54) 【発明の名称】 音声認識装置およびプログラム

(57) 【特許請求の範囲】

【請求項1】

音声認識のためのハイブリッド隠れマルコフ/ベイジアンネットワーク(HMM/BN)モデルを記憶する記憶装置を含み、隠れマルコフモデル(HMM)は時間的な音声の特徴をモデリングするのに用いられ、ベイジアンネットワーク(BN)は状態確率モデルをあらわすのに用いられ、さらに

当該記憶装置に記憶されたHMM/BNモデルを用いて入来する音声データをデコードする音声デコーダを含む、音声認識装置。

【請求項2】

BNの状態確率モデルが、変数Xと、以下で計算される条件付確率 $P(Y|Q)$ とをさらに含み、

【数1】

$$P(Y|Q) = \frac{1}{N(x)} \sum_x P(Y|X=x, Q)$$

ただし、Yは一連の観測パラメータを表わし、Qは状態変数を表わし、XはQとは独立してYの値に影響を与える所定の要素を反映する変数を表わし、xはXがとり得る値の一つであり、N(x)はXがとり得る値の数である、請求項1に記載の音声認識装置。

【請求項3】

前記変数Xが環境ノイズを表わす、請求項2に記載の音声認識装置。

【請求項 4】

前記変数 X が話者の識別情報を表わす、請求項 2 に記載の音声認識装置。

【請求項 5】

前記変数 X が話者の母語を表わす、請求項 2 に記載の音声認識装置。

【請求項 6】

状態確率モデルが変数 N と S とをさらに含み、条件付確率 $P(Y|Q)$ は以下で計算され、

$$P(Y|Q) = \frac{1}{N(n,s)} \sum_{n,s} P(Y|N=n, S=s, Q)$$

10

ただし、Y は一連の観測パラメータを表わし、Q は状態変数を表わし、N および S は、Q とは独立して Y の値に影響を与える所定の要素を反映する変数を表わし、n および s は N および S がそれぞれとり得る値の一つであり、 $N(n, s)$ は N および S がとり得る値の組合せの数である、請求項 1 に記載の音声認識装置。

【請求項 7】

前記変数 N がノイズの種類を表わす、請求項 6 に記載の音声認識装置。

【請求項 8】

前記変数 S が入来する音声データの信号対雑音比を表わす、請求項 6 または請求項 7 に記載の音声認識装置。

20

【請求項 9】

請求項 1 から請求項 8 のいずれかに記載の音声認識装置としてコンピュータを動作させるための、コンピュータで実行可能な音声認識プログラム。

【請求項 10】

請求項 1 から請求項 8 のいずれかに記載の音声認識装置としてコンピュータを動作させるための、コンピュータで実行可能な音声認識プログラムを記憶した、コンピュータ読取可能な記憶媒体。

【発明の詳細な説明】

【0001】

【発明の属する技術分野】

30

この発明は音声認識システムに関し、特に、ベイジアンネットワーク (BN) を採用した HMM (隠れマルコフモデル) に基づく音声認識システムに関する。

【0002】

【従来の技術】

多年にわたり、HMM の音声認識への導入以来、各状態 Q に対する観測の、条件付き分布 $P(y|Q)$ は確率密度関数の混合によりモデル化されてきた (離散 HMM はここでは考慮しない)。ガウス pdf (確率密度関数) およびラプラシアン pdf がこの目的でよく使用されている。後には、ハイブリッド HMM/NN (隠れマルコフモデル/ニューラルネットワーク) システムが提案されたが、ここでは所与の入力観測に対し HMM の状態の尤度を推定するのにニューラルネットワークが用いられている。

40

【0003】

多くの場合、音声スペクトルから抽出された特徴がこれらの観測値を形成する。しかしながら、音声認識に関する研究により、これらの特徴のみを用いるだけでは高いシステム性能を達成するには不十分である事が示された。このため、多くの研究者が、HMM システムに何か他の知識を表わす付加的な特徴を含めようとしてきた。

【0004】

トクダら (後掲の非特許文献 1) は、付加的なピッチ情報をモデル化するため、多元空間での確率分布を提案している。しかしほとんどの場合、付加的特徴の特性により種々の方策がとられている。この問題に対処すべき共通の、十分に柔軟性の高いフレームワークはこれまでには存在しなかった。

50

【 0 0 0 5 】

最近、HMMに対する別の選択肢として、ベイジアンネットワーク(BN)が研究者の関心を集めている。BNは人工知能の研究分野では周知でありよく研究されている。しかし、音声認識においては、これらは比較的新しい研究課題である。

【 0 0 0 6 】

ベイジアンネットワークとは有向非巡回的グラフであり、そのノードは事象を表わす。ベイジアンネットワークでは、第一のノード(ノードA)から第二のノード(ノードB)への枝は、AがBの原因である事を示す。各枝(A-B)には、確率 $P(B|A)$ を示す「リンク」行列が割当てられ、この確率は、Aの各値に対するBの各値の確率を特定する。

【 0 0 0 7 】

ベイジアンネットワークは、多くの互いに異なる(離散的または連続的な)ランダムな変数の複合的な同時確率分布を、よく構造化され容易に表現可能な仕方でモデル化できる。特に、時間的音声特徴をモデル化するのに適しているのは、ディーンらの提案するダイナミックBN(DBN)である(後掲の非特許文献2)。

【 0 0 0 8 】

音声認識におけるDBNの最初の報告のいくつかでは、ツバイクら、およびダウディらの報告でのように(後掲の非特許文献3、4)、これらは単独の語認識作業での語モデルとして用いられていた。これらの著作においては、DBNは、音声スペクトル情報に加えて、調音的特徴、サブ帯域相関、話し方のスタイル等の付加的な知識を容易に組み入れる事が可能な、HMMの一般化と考えられている。スティーブソンらはDBNのフレームワーク内で、音声学的な特徴をピッチ情報で容易に補う事ができる事を報告している。

【 0 0 0 9 】

ベイジアンネットワークの別の利点は、認識中に信頼性をもって推定するのが困難な付加的な特徴を、隠されたまま、すなわち観測不能な状態にしておける事である。

【 0 0 1 0 】

【非特許文献1】

K.トクダ、T.マスコ、N.ミヤザキ、T.コバヤシ、『ピッチパターンモデリングのための多元空間確率分布に基づく隠れマルコフモデル』、ICASSP予稿集、pp. 229-232, 1999年。(K. Tokuda, T. Masuko, N. Miyazaki, and T. Kobayashi, "Hidden Markov Models based on multi-space probability distribution for pitch pattern modeling." In Proc. ICASSP, pp. 229-232, 1999.)

【非特許文献2】

T.ディーン、K.カナザワ、「確率的な時間的推論」、AAAI, pp. 524-528, 1988年。(T. Dean and K. Kanazawa, "Probabilistic temporal reasoning," in AAAI, pp. 524-528, 1988.)

【非特許文献3】

G.ツヴァイク、S.ラッセル、『自動音声認識のためのベイジアンネットワークによる確率モデリング』、ICSLP予稿集、pp. 3010-3013, 1988年。(G. Zweig and S. Russell, "Probabilistic modeling with Bayesian Networks for automatic speech recognition," in Proc. ICSLP, pp. 3010-3013, 1988.)

【非特許文献4】

K.ダウディ、D.フォア、C.アントワーン、『確率的グラフモデルに基づく多帯域音声認識の新たなアプローチ』、ICSLP予稿集, vol. 1, pp. 329-332, 2000年。(K. Daoudi, D. Fohr, and C. Antoine, "A new approach for multi-band speech recognition based on probabilistic graphical models," in Proc. ICSLP, vol. 1, pp. 329-332, 2000.)

【非特許文献5】

T.スティーブソン、M.マシュウ、H.ボラード、『ベイジアンネットワークに基づ

10

20

30

40

50

『ASRでの補助情報のモデル化』、ユーロスピーチ予稿集、pp.] 2765 - 2768, 2001年。(T. Stephenson, M. Mathew, and H. Bourlard, "Modeling auxiliary information in Bayesian Network based ASR," in Proc. Eurospeech, pp. 2765-2768, 2001.)

【非特許文献6】

K.ダウディ、D.フォア、C.アントワヌ、「ベイジアンネットワークを用いた連続多帯域音声認識」、Proc. ASRU, 2001年。(K. Daoudi, D. Fohr, and C. Antoine, "Continuous multi-band speech recognition using Bayesian Networks," in Proc. ASRU, 2001.)

【発明が解決しようとする課題】

BNのこれらの魅力的な特性にも関わらず、音声認識への応用は依然として、小さな、単独の語認識作業に対するものに限定されている。その理由は、BNパラメータの学習および推論のための既存のアルゴリズムが、連続した音声認識(CSR)および、特に語彙の多いCSR作業にはそれほど適していないからである。連続して発話された数字の認識を可能にするDBN語モデルの拡張が、ダウディらの非特許文献6に報告されているものの、作業対象となる語彙がわずかに数百に増加するだけでも、あまりにも計算量が多くなり過ぎる。

【0011】

従って、この発明の目的の一つは、BNを採用したものであって、かつCSRに適した装置を提供する事である。

【0012】

この発明の別の目的は、BNを採用したものであって、かつ語彙量の多いCSR作業に適した装置を提供する事である。

【0013】

この発明のさらに別の目的は、BNを採用したものであって、かつBNのための煩瑣なパラメータ学習を必要としない装置を提供する事である。

【課題を解決するための手段】

この発明のある局面に従った音声認識装置は、音声認識のためのハイブリッド隠れマルコフ/ベイジアンネットワーク(HMM/BN)モデルを記憶するための記憶装置を含む。隠れマルコフモデル(HMM)は時間的な音声の特徴をモデリングするのに用いられ、ベイジアンネットワーク(BN)は状態確率モデルをあらわすのに用いられる。この装置はさらに、記憶装置に記憶されたHMM/BNモデルを用いて入来する音声データをデコードするための音声デコーダを含む。

【0014】

BNの状態確率モデルは、変数Xを含んでもよく、条件付確率 $P(Y|Q)$ は以下の式によって計算され、

【0015】

【数3】

$$P(Y|Q) = \frac{1}{N(x)} \sum_x P(Y|X=x, Q)$$

ただし、Yは一連の観測パラメータを表わし、Qは状態変数を表わし、XはQとは独立してYの値に影響を与える所定の要素を反映する変数を表わし、xはXがとり得る値の一つであり、N(x)はXがとり得る値の数である。

【0016】

好ましくは、変数Xは環境ノイズ、話者認識情報、または話者の母語を表わしてもよい。

【0017】

状態確率モデルは変数NとSとをさらに含んでもよく、条件付確率 $P(Y|Q)$ は以下で計算され、

10

20

30

40

50

【 0 0 1 8 】

【 数 4 】

$$P(Y|Q) = \frac{1}{N(n,s)} \sum_{n,s} P(Y|N=n, S=s, Q)$$

ただし、Yは一連の観測パラメータを表わし、Qは状態変数を表わし、NおよびSはQとは独立してYの値に影響を与える所定の要素を反映する変数を表わし、nおよびsはNおよびSがそれぞれとり得る値の一つであり、N(n,s)はNおよびSがとり得る値の組合せの数である。

【 0 0 1 9 】

変数Nはノイズの種類を表わし、変数Sは入来する音声データの信号対雑音比を表わしてもよい。

【 0 0 2 0 】

この発明の別の局面は、上述の音声認識装置としてコンピュータを動作させる、コンピュータで実行可能な音声認識プログラムに関する。

【 0 0 2 1 】

この発明のさらに別の局面は、上述の音声認識装置としてコンピュータを動作させる、コンピュータで実行可能な音声認識プログラムを記憶する、コンピュータ読取可能な記憶媒体に関する。

【 0 0 2 6 】

【 発明の実施の形態 】

- ハイブリッドHMM / BNモデリング -

多くの場合、音声認識でBNを使用するのは、HMMをダイナミックなBNとして表わそうとする考えに基づいている。このような表現を図1に示す。ここで、 Q_t は状態変数であり、 Y_t は時刻 $t = 1, 2, 3, 4, \dots$ での連続な観測変数である。枝は変数間の確率的依存を表わす。状態インスタンス間の枝はHMM遷移確率を表わし、状態インスタンスと観測インスタンスとの間の枝はHMM状態の条件付分布を表わす。以下の図において、四角で囲った変数は離散的であり、丸で囲った変数は連続である。ハッチングした丸 / 四角は観測可能な変数を示す。

【 0 0 2 7 】

$Y = y_1, \dots, y_T$ 、Mは本発明のモデルとしたとき、入力される観測シーケンス $P(Y|M)$ の尤度を得る必要がある場合、HMMおよびBNの表現にBN推論アルゴリズムを用いる必要がある。この作業の間に、ネットワークのサイズが入力シーケンスTのサイズに合わせて調整され、その後、ネットワーク全体から $P(Y|M)$ が推論される。

【 0 0 2 8 】

状態ノード間の枝を切断した場合を想定する。こうする事により、各時刻tに対応する、図2に示すように複数の独立したBNが得られる。時間遷移(切断された枝)が従来のHMMにより支配されるとすれば、これらのBNを適切なHMM状態に割当てる事で時間指標をなくす事ができる。BNは全て同じ共通の構造を有するので、それらを図3のように単一のBNとして表わす事ができる。図3では、変数Qは音声学的モデル内の全てのHMMの状態指標(S_{ij})の値をとり、状態確率分布 $P(Y|Q=S_{ij})$ は枝によって表わされる。

【 0 0 2 9 】

こうして、ガウスの混合ではなく、状態分布モデルとしてBNを有するように従来のHMMを修正した。HMMとBNとをこのように組合せる事で、HMM / BNモデルが階層的となる。BNは最下層にあり、HMMは最上層にある。なお、状態変数Q(図3)はBNに対しては観測可能となるが、上のHMMレベルでは、依然として隠されたままである。

【 0 0 3 0 】

状態BNは付加的な知識を表わす他のランダム変数に容易に拡張可能である。簡単な作業とは言いがたいデータからの学習によってではなく、変数間の関係に関する発明者らの知

10

20

30

40

50

識に従って、拡張BNのグラフィック構造を課する事ができる。この実施例は煩瑣なパラメータ学習とは無縁である。

【0031】

拡張状態BNの可能な構造のいくつかを図4に示す。たとえば、変数Xはこの図では環境ノイズの種類を表わす事ができ、他のWおよびZ変数は話者のID（識別情報）および話者の母語を表わす事ができる。

【0032】

このHMM/BNモデルで認識を行なう場合には、従来のHMMと同様に、各状態 $Q = q_{ij}$ について $P(y|Q)$ を計算する必要がある。ただし、「i」はHMM指標であり、「j」はi番目のHMMの状態指標である。この値をBN確率モデルから推論する事ができ、これを行うアルゴリズムとして、正確なものと近似的なものを含め多数の推論アルゴリズムがある。単純なBNでは、図4(a)に示すように、「力任せ」法さえも適用可能である。このBNに対する同時確率モデルは、以下のようにチェーンルールで表わす事ができる。

【0033】

【数5】

$$P(X, Y, Q) = P(Y|X, Q) * P(X|Q) * P(Q) \quad (1)$$

XおよびQは独立した変数なので、「 $P(X|Q) = P(X)$ 」が成立し、従って上の式は以下のように書換える事ができる。

【0034】

【数6】

$$P(X, Y, Q) = (P|X, Q) * P(X) * P(Q) \quad (2)$$

従って、求める確率 $P(Y|Q)$ はXに対するマージナライゼーションにより、以下のようになる。

【0035】

【数7】

$$\begin{aligned} P(Y|Q) &= \frac{P(Y, Q)}{P(Q)} = \frac{\sum_x P(Y, X = x, Q)}{P(Q)} \\ &= \frac{\sum_x P(Y|X = x, Q) * P(X = x) * P(Q)}{P(Q)} \\ &= \sum P(Y|X = x, Q) * P(X = x) \quad (3) \end{aligned}$$

実用上、多くの場合には全ての $X = x$ について $P(X)$ は同じであると仮定できるので、式(3)を次のように変形できる。

【0036】

【数8】

$$P(Y|Q) = \frac{1}{N(x)} \sum_x P(Y|X = x, Q) \quad (4)$$

ここで、 $N(x)$ は変数Xがとり得る値の数である。

【0037】

BNパラメータのトレーニングは、従来のHMM状態パラメータのトレーニングとほぼ同様に、各状態について独立して行なう事ができる。トレーニングの詳細は図10を参照して後述する。

【0038】

- ノイズの多い音声認識システムにおけるHMM/BNモデル -

10

20

30

40

50

音声が入力を含む場合、音声の特徴ベクトルはその分布を変え、その変化はノイズの種類とともにSNR（信号対雑音比）の値にも依存する。このため、この依存性を図5に示すような種類の状態BNで表わす事ができる。ここで、NとSとはそれぞれ、ノイズの種類とSNR値とを示す、隠れた離散変数である。この場合、状態の尤度は式(3)を導出したのと同じ方法で解析的に表現できる。

【0039】

多くの場合、事前確率P(N)およびP(S)は、ノイズの各種類および各SNR値について等しいと合理的に仮定できるので、以下のとおりとなる。

【0040】

【数9】

$$P(Y|Q) = \frac{1}{N(n,s)} \sum_{n,s} P(Y|N=n, S=s, Q) \quad (5)$$

語モデルおよびサブワードモデルも、従来のHMMの場合と同様に作成される。デコードもまた、デコーダを変更する事なく、標準的なHMMベースのシステムと同様に行なう事ができる。N(n,s)はNとSのとり得る組合せの数を示す。

【0041】

この実施例のHMM/BNモデルの構造を図6に要約して示す。

【0042】

- コンピュータでの実現例 -

図7はこの実施例の全体図である。図7を参照して、このシステムは、トレーニングデータ110とHMM/BNモデル112とに基づいてHMM/BNハイブリッドモデル114をトレーニングするためのトレーニングシステム100と、HMM/BNハイブリッドモデル114を記憶するための媒体104と、コンピュータシステムで実現され、音声データ120を媒体104に記憶されたHMM/BNハイブリッドモデル114でデコード（音声認識122）し、認識された音声124を出力するための音声認識システム（音声デコーダ）102とを含む。

【0043】

図8はこのコンピュータシステム30を概略的に示し、図9はシステム30をブロック図形式で示す。図8を参照して、このコンピュータシステム30は、FD（フレキシブルディスク）ドライブ52およびCD-ROM（コンパクトディスク読出専用メモリ）ドライブ50を有するコンピュータ40と、キーボード46と、マウス48と、モニタ42とを含む。

【0044】

図9を参照して、コンピュータ40は、FDドライブ52およびCD-ROMドライブ50に加えて、CPU（中央処理装置）56と、CPU56、FDドライブ52およびCD-ROMドライブ50に接続されたバス66と、ブートアッププログラムおよびHMMデコードプログラム等を記憶する読出専用メモリ（ROM）58と、バス66に接続され、プログラム命令、システムプログラム、およびHMM/BNハイブリッドモデルデータを記憶するランダムアクセスメモリ（RAM）60とを含む。

【0045】

ここでは示さないが、コンピュータ40はさらにローカルエリアネットワーク（LAN）への接続を提供するネットワークアダプタボードを含んでもよい。音声認識をリアルタイムで行なう場合には、コンピュータシステム30はさらにマイクロフォンとオーディオキャプチャボードとを含むことになる。

【0046】

コンピュータシステム30に音声認識を行なわせるプログラムは、CD-ROMドライブ50またはFDドライブ52に挿入されるCD-ROM62または図示しないFDに記憶され、さらにハードディスク54に転送される。またはこれに代えて、プログラムは図示しないネットワークを通じてコンピュータ40に送信されハードディスク54に記憶さ

10

20

30

40

50

れてもよい。プログラムは実行の際にRAM 60にロードされる。CD-ROM 62、FD、またはネットワークを介してRAM 60にプログラムを直接ロードしてもよい。

【0047】

プログラムは、コンピュータ40にこの実施例の音声認識を行なわせるいくつかの命令を含む。この方法を行なわせるのに必要な基本的機能のいくつかはコンピュータ40のオペレーティングシステム(OS)またはサードパーティのプログラム、もしくはコンピュータ40にインストールされるHMMツールキット等のモジュールにより提供されるので、このプログラムはこの実施例の方法を実現するのに必要な機能全てを必ずしも含まなくてよい。このプログラムは、命令のうち、所望の結果が得られるように制御されたやり方で適切な機能または「ツール」を呼出す事により音声認識プロセスを実行する命令のみを
10

【0048】

音声認識システム102は従来のHMMデコードプログラムで実現される。この実施例で新規な点は、HMM/BNハイブリッドモデル114が媒体104上に記憶されている事
である。

【0049】

- HMM/BNモデルでのトレーニングと認識 -

HMM/BNモデルのトレーニングには、HMM/NNトレーニングと同じアプローチを採用する事ができる。これはビタビトレーニングアルゴリズムに基づくものである。全体の制御フローを図10に示す。まず始めに、HMMと状態ベイジアンネットワークとのトポロジーを選択する。モデルを初期化し、特徴抽出ステップ150で入力音声140から抽出した特徴144に対し、状態分離ステップ152で、ビタビライメントが基本HMMモデル142を用いた基本認識部を使用して行なわれる。
20

【0050】

状態の分離は状態ベイジアンネットワークのためのトレーニングデータを生成するのに用いられ、つぎにこの状態ベイジアンネットワークはステップ154、156および160でトレーニングされる。このビタビトレーニング手法では、トレーニングの時間的な部分と静的な部分とが分離される。この処理では、ステップ158で終了条件が満たされるまで、ステップ154および156のBNトレーニングと、埋込まれたトレーニング処理であるステップ160の遷移確率の再推定とを交互に繰り返す。
30

【0051】

このHMM/BNモデルで認識を行なう場合、従来のHMMと同様に、通常のビタビデコードアルゴリズムが用いられる。ここで、各状態 $Q = q_{ij}$ (図6を参照)について $P(y | Q)$ を計算する必要がある。ここで i はHMM指標であり、 j は i 番目のHMMの状態指標である。この値を、標準的推論アルゴリズムを用いてBN確率モデルから推論する事ができる。

【0052】

- Aurora 2タスクでの評価 -

この実施例の音声認識システムをAurora 2タスクで評価する実験を行なった。これらの実験では、公式のAurora 2タスクで示唆されている評価用シナリオに忠実に従った。最も関心があったのは、HMM/BNシステムを複数条件でトレーニングされたHMMシステムと比較する事であった。HMM/BN状態の条件付き分布のトレーニングにあたっては、トレーニングデータをノイズの種類とSNR値とで分け、HTK(HMMツールキット、HMM処理のためのソフトウェアツールキット)を用いて各条件についてのパラメータを個別にトレーニングした。特徴ベクトル、語モデル、状態数、実験条件等の全ての他のシステムパラメータは同一にした。なお、HMM/BNシステムでは、何らかの適合やノイズに強い方法を用いたわけではない。二つのシステムでの主な機能的相違点は、HMM/BNシステムが音声の特徴およびノイズの隠れた依存性を探求する事である。
40

【 0 0 5 3 】

テストセット A (トレーニングデータと同じノイズ種類) とテストセット B (異なるノイズ) の認識結果を表 1 にまとめて示す。理解できるとおり、HMM / BN システムの性能は閉じたノイズ条件 (A セット) ではかなり高く、ずっと複雑なシステムでこの作業について得られる最新の結果に迫っている。

【 0 0 5 4 】

【表 1】

SNR	テストセット A		テストセット B	
	HMM	HMM / BN	HMM	HMM / BN
Clean	98.54	98.83	98.54	98.83
20 dB	97.52	98.12	96.96	97.26
15 dB	96.94	97.65	95.38	95.05
10 dB	94.59	96.04	92.58	90.27
5 dB	87.51	91.7	83.5	78
0 dB	59.84	76.11	58.91	48.7
-5 dB	23.46	35.79	23.86	3.18
平均*	87.29	91.92	85.46	81.85
誤差	12.71	8.08	14.54	18.15

*20dB から 0dB の値で計算

この評価は、ノイズの種類と SNR 値とを付加的なパラメータとして加え、これらの依存性を検討する事で、結果として得られるスペクトル特徴パラメータの誤差が 36.4% ($= (12.71 - 8.08) * 100 / 12.71$) 減少する事を示している。

【 0 0 5 5 】

B セットの条件についていえば、性能の劣化が見られる。これは、音声スペクトル特徴の分布のミスマッチに加えて、この HMM / BN システムには新たなノイズの依存性に関する知識が得られないという事実で説明がつく。他方で、複数条件 HMM システムでは、状態ガウス混合は多ノイズおよび SNR 条件から複雑な分布をあまりうまくモデル化できていない事が明らかである。しかしながら、このデータとモデル分布のミスマッチにはある種の平滑化の効果があり、これによって、見られないデータから一般化するというこのモデルの能力が高まる。

【 0 0 5 6 】

[ハイブリッド HMM / BN モデルの応用]

明らかに、提案されたハイブリッド HMM / BN モデルはノイズの多い音声認識システムのみでなく、観測や隠れた特徴が新たに得られる事で性能上の利益が得られる他の多くの場合に適用可能である。このアプローチは、種々の空間からの特徴を組合せ、その間の依存性を検討する事により、システムのモデル化能力を高めるといふ、より一般的なフレームワークという方がより正確である。とくに興味深いのは、BN の隠れた変数の確率が推論できる可能性である。このように、HMM / BN システムはこれらの付加的パラメータの認識に用いる事ができる。

【 0 0 5 7 】

例えば、ある付加的な隠れ変数 X が多言語システムでの言語を表わす場合、各フレームについて $P(X | Q)$ を計算し、これらの確率を、入力された発話全体にわたり累積する事ができる。その場合、 $x = \arg \max_x P(x | Q_s)$ 、 Q_s は最高の仮定状態シーケンス、となる x は、その発話がなされた言語として最も確率の高い言語を示す。従って、多言語の音声認識に加えて、このようなシステムは言語認識をも行なう事ができる。な

10

20

30

40

50

お、関数 $x = \operatorname{argmax}_x P(x | Q_s)$ は、 $P(x | Q_s)$ を最大にする x を示す。

【0058】

[この実施の形態の効果]

この実施の形態では、HMMとBNとを単一のモデル内で組合せ、HMMとBNとの両者の長所を活かしている。ハイブリッドHMM/BNモデルにより、音声認識システムに他の情報を容易に加える事が可能になり、最小限の費用でその性能を高める事ができる。さらに、HMM/BNモデルは従来のHMMと同様に、サブワードの音声ユニットを表わす事ができる。こうして、BNフレームワークを語彙数の多い連続した音声認識に用いる事が可能になる。

【図面の簡単な説明】

【図1】 HMMをDBNとして表わす模式図である。

【図2】 各時刻 t での複数BNの模式図である。

【図3】 通常のBN構造を示す、状態BNの模式図である。

【図4】 (a) は離散の変数を一つ付加した状態BNの模式図であり、(b) はより複雑な構造の状態BNの模式図である。

【図5】 ノイズおよびSNR変数を有する状態BNの模式図である。

【図6】 この発明の一実施の形態に係るHMM/BNハイブリッドモデル構造の模式図である。

【図7】 この発明の一実施の形態の音声認識システムのブロック図である。

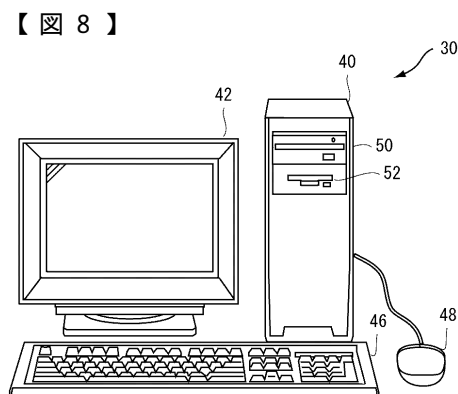
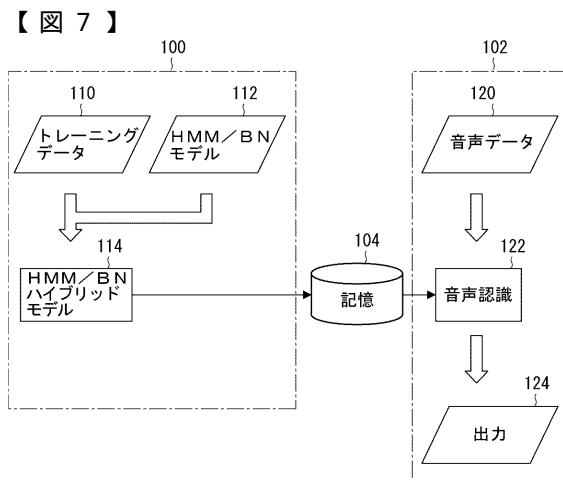
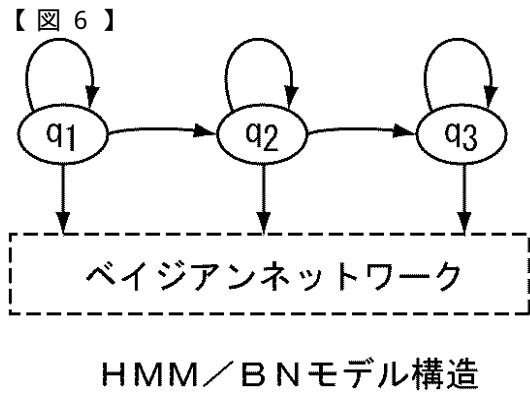
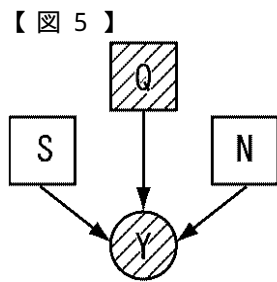
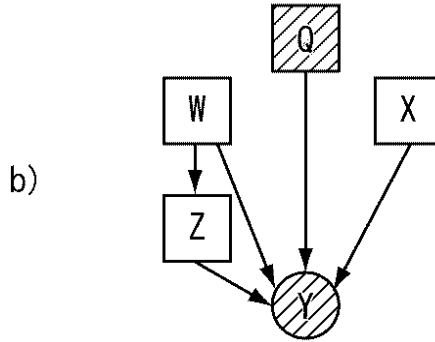
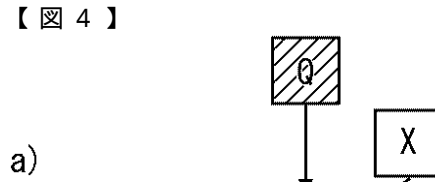
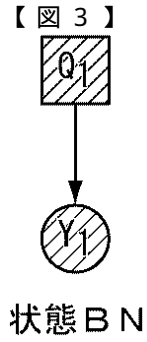
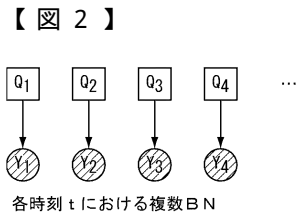
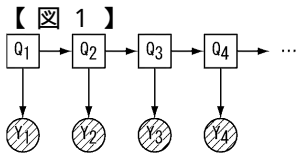
【図8】 この実施の形態の音声認識システムのプログラムを実行するコンピュータシステムの外觀図である。 20

【図9】 図8のコンピュータシステムのブロック図である。

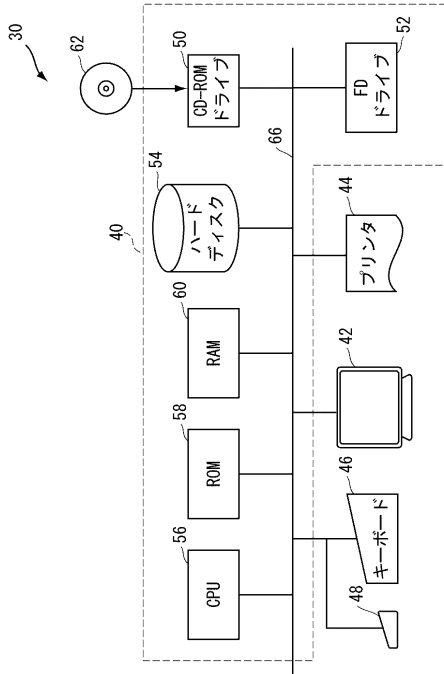
【図10】 この実施の形態のHMM/BNモデルのトレーニングプログラムの制御の流れを示すフローチャートである。

【符号の説明】

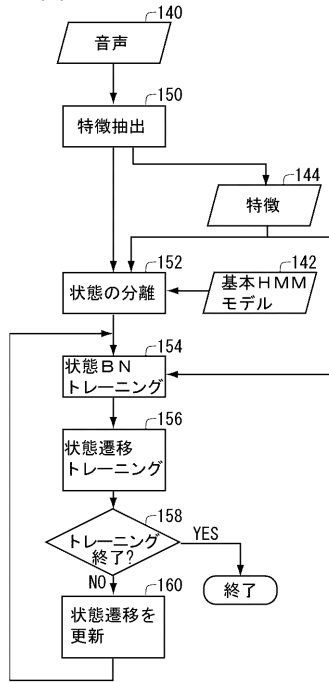
100 トレーニングシステム、102 音声認識システム、104 記憶媒体、110 トレーニングデータ、112 HMM/BNモデル、114 HMM/BNハイブリッドモデル、120 音声データ、122 音声認識、124 出力



【図9】



【図10】



フロントページの続き

(56)参考文献 特許第3508978(JP, B2)

中川, 音声認識においてHMMとトライグラムを超えるもの, 人工知能学会誌, 日本, 2002年
1月1日, Vol.17, No.12, p.35-40

本村他, ベイジアンネットワーク - 不確定性のモデリング技術 -, 人工知能学会誌, 日本, 2000年
7月, Vol.15, No.4, p.575-582

中川, 音声認識研究の動向, 電子情報通信学会論文誌 D-II, 日本, 2000年 2月25日,
Vol.J83-D-II, No.2, p.433-457

S. K. Riis et al., Hidden neural networks: a framework for HMM/NN hybrids, Proceedings
of the 1997 IEEE International Conference on Acoustics, Speech, and Signal Processing
(ICASSP '97), 米国, 1997年 4月21日, Vol.4, p.3233-3236

(58)調査した分野(Int.Cl., DB名)

G10L 15/10

G10L 15/14

JSTPlus(JDream2)

IEEE