

(19)日本国特許庁 ( J P )

(12) 特 許 公 報 ( B 1 )

(11)特許番号

第2938866号

(45)発行日 平成11年(1999) 8月25日

(24)登録日 平成11年(1999) 6月11日

(51)Int.Cl. <sup>6</sup>	識別記号	F I	
G 1 0 L 3/00	5 3 7	G 1 0 L 3/00	5 3 7 J 5 3 7 D

請求項の数5 (全 15 頁)

<p>(21)出願番号 特願平10-243024</p> <p>(22)出願日 平成10年(1998) 8月28日</p> <p>審査請求日 平成10年(1998) 8月28日</p> <p>特許法第30条第1項適用申請有り 日本音響学会平成10年度春季研究発表会講演論文集(平成10年3月17日) 41-42頁に発表</p> <p>特許法第30条第1項適用申請有り 平成10年3月17日に慶応義塾大学において開催された日本音響学会平成10年度春季研究発表会において発表</p>	<p>(73)特許権者 593118597 株式会社エイ・ティ・アール音声翻訳通信研究所 京都府相楽郡精華町大字乾谷小字三平谷5番地</p> <p>(72)発明者 政瀧 浩和 京都府相楽郡精華町大字乾谷小字三平谷5番地 株式会社エイ・ティ・アール音声翻訳通信研究所内</p> <p>(72)発明者 匂坂 芳典 京都府相楽郡精華町大字乾谷小字三平谷5番地 株式会社エイ・ティ・アール音声翻訳通信研究所内</p> <p>(74)代理人 弁理士 青山 葆 (外2名)</p> <p>審査官 涌井 智則</p>
---	---

最終頁に続く

(54)【発明の名称】 統計的言語モデル生成装置及び音声認識装置

3

(57)【特許請求の範囲】

【請求項1】 複数のクラスタの統計的言語モデルを記憶する記憶手段と、  
 所定の複数の発声音声文を含む学習用テキストデータに基づいて各発声音声文に対する統計的言語モデルを生成して、各発声音声文が各クラスタに対応するように、上記生成した統計的言語モデルを上記記憶手段に記憶する初期化手段と、  
 上記学習用テキストデータの各発声音声文について、各クラスタにおける統計的言語モデルの文生成確率を計算して最大の文生成確率を有するクラスタを選択してその発声音声文を所属させるように上記記憶手段に記憶するクラスタ選択手段と、  
 発声音声文が属するクラスタが変化したときに、各クラスタ毎に、上記クラスタ選択手段によって選択された発

4

発声音声文を用いて上記記憶手段に記憶された各統計的言語モデルを更新して、各クラスタに対応した統計的言語モデルを生成するモデル変更手段と、  
 上記複数の発声音声文に属するクラスタが1文も変化しなくなるまで、上記クラスタ選択手段の処理と、上記モデル変更手段の処理を繰り返す制御手段とを備えたことを特徴とする統計的言語モデル生成装置。  
 【請求項2】 請求項1記載の統計的言語モデル生成装置において、上記統計的言語モデル生成装置はさらに、上記記憶手段に記憶された各クラスタ毎のテキストデータに基づいて、最尤推定法を用いて各クラスタ毎に単語の N - g r a m ( N は 2 以上の自然数である。 ) の遷移確率を演算する第1の演算手段と、  
 上記第1の演算手段によって演算された各クラスタ毎の単語の N - g r a m の遷移確率の出現分布を事前知識の

10

所定の確率分布と仮定し、各クラス毎の確率分布の加重平均及び加重分散を演算した後、演算された加重平均と加重分散に基づいて事前知識の確率分布のパラメータを演算する第2の演算手段と、

上記第2の演算手段によって演算された事前知識の確率分布のパラメータと、上記学習用テキストデータうちの特定クラス毎のテキストデータの事後知識における処理対象の単語列の直前の単語列の出現回数と、処理対象の単語列の出現回数とに基づいて、各クラス毎の単語の N - g r a m の遷移確率を計算することにより、各クラス毎の単語の N - g r a m の遷移確率を含む統計的言語モデルを生成する第3の演算手段とを備えたことを特徴とする統計的言語モデル生成装置。

【請求項3】 請求項2記載の統計的言語モデル生成装置において、上記統計的言語モデル生成装置はさらに、上記第3の演算手段によって演算された各クラス毎の単語の N - g r a m の遷移確率に基づいて、所定の平滑化処理を実行し、処理後の各クラス毎の単語の N - g r a m の遷移確率を含む統計的言語モデルを生成する第1の生成手段を備えたことを特徴とする統計的言語モデル生成装置。

【請求項4】 請求項3記載の統計的言語モデル生成装置において、上記統計的言語モデル生成装置はさらに、上記学習用テキストデータに基づいて、最尤推定法を用いて単語の N - g r a m ( N は 2 以上の自然数である。)の遷移確率を演算して、上記単語の N - g r a m の遷移確率を含む別の統計的言語モデルを生成する第2の生成手段を備えたことを特徴とする統計的言語モデル生成装置。

【請求項5】 入力される発声音声文の音声信号に基づいて、所定の統計的言語モデルを用いて音声認識する音声認識装置において、

請求項4記載の統計的言語モデル生成装置と、上記第2の生成手段によって生成された別の統計的言語モデルを用いて、入力される発聲音声文の音声信号を音声認識して第1の認識仮説を出力する第1の音声認識手段と、

上記第1の音声認識手段から出力される第1の認識仮説に回答して、上記第1の生成手段によって生成された各クラス毎の統計的言語モデルを用いて、入力される発聲音声文の音声信号を音声認識して、文生成確率が最大のクラスの統計的言語モデル生成装置を選択するモデル選択手段と、

上記モデル選択手段によって選択されたクラスの統計的言語モデルを用いて、上記第1の音声認識手段から出力される第1の認識仮説に対して絞込処理を行って第2の認識仮説を生成して認識結果として出力する第2の音声認識手段とを備えたことを特徴とする音声認識装置。

【発明の詳細な説明】

【0001】

【発明の属する技術分野】本発明は、学習用テキストデータに基づいて統計的言語モデルを生成する統計的言語モデル生成装置、及び上記統計的言語モデルを用いて、入力される発聲音声文の音声信号を音声認識する音声認識装置に関する。

【0002】

【従来の技術】近年、連続音声認識装置において、その性能を高めるために言語モデルを用いる方法が研究されている。これは、言語モデルを用いて、次単語を予測し探索空間を削減することにより、認識率の向上及び計算時間の削減の効果を狙ったものである。最近盛んに用いられている言語モデルとして N - g r a m ( N - グラム;ここで、Nは2以上の自然数である。)がある。これは、大規模なテキストデータを学習し、直前の N - 1 個の単語から次の単語への遷移確率を統計的に与えるものである。複数 L 個の単語列  $w_1^L = w_1, w_2, \dots, w_L$  の生成確率  $P(w_1^L)$  は次式で表される。

【0003】

【数1】

$$P(w_1^L) = \prod_{t=1}^L P(w_t | w_{1:t-1})$$

【0004】ここで、 $w_t$  は単語列  $w_1^L$  のうち t 番目の 1 つの単語を表し、 $w_i^j$  は i 番目から j 番目の単語列を表す。上記数1において、確率  $P(w_t | w_{1:t-1})$  は、N個の単語からなる単語列  $w_{1:t-1}$  が発声された後に単語  $w_t$  が発声される確率であり、以下同様に、確率  $P(A | B)$  は単語又は単語列 B が発声された後に単語 A が発声される確率を意味する。また、数1における「 $\prod$ 」は t = 1 から L までの確率  $P(w_t | w_{1:t-1})$  の積を意味し、以下同様である。

【0005】ところで、近年、上記統計的言語モデル N - g r a m を用いて連続音声認識の性能を向上させる手法が盛んに提案されている(例えば、従来技術文献1「L.R.Bahl et al., "A Maximum Likelihood Approach to Continuous Speech Recognition", IEEE Transactions on Pattern Analysis and Machine Intelligence, pp. 179 - 190, 1983年」及び従来技術文献2「清水ほか, "単語グラフを用いた自由発話音声認識", 電子情報通信学会技術報告, SP95-88, pp. 49-54, 平成7年」参照。)

【0006】しかしながら、N - g r a m はパラメータ数が多く、それぞれの値を正確に求めるためには、莫大な量のテキストデータが必要とされる。この問題を解決する方法として、学習用テキストデータに出現しない単語遷移に対しても遷移確率を与える平滑化の手法(例えば、従来技術文献3「F. Jelinek et a

l., "Interpolated estimation of Markov Source Parameters from Sparse Data", Proceedings of Workshop Pattern Recognition in Practice, pp. 381 - 387, 1980年」、従来技術文献4「S.M.Katz et al., "Estimation of Probabilities from Sparse Data for the Language model Component of a Speech Recognizer", IEEE Transactions on Acoustics, Speech, and Signal Processing, pp. 400 - 401, 1987年」及び従来技術文献5「川端ほか, "二項事後分布に基づくN-gram言語モデルのBack-off平滑化", 電子情報通信学会技術報告, SP95-93, pp. 1-6, 平成7年」参照。)や、クラス分類、可変長N-gram等パラメータの数を減少させる手法(例えば、従来技術文献6「P.F.Brown et al., "Class-Based n-gram models of natural language", Computational Linguistics, Vol. 18, No. 4, pp. 467 - 479, 1992年」、従来技術文献7「T.R.Niesler et al., "A Variable-Length Category-Based N-gram Language Model", Proceedings of ICASSP '96, Vol. 1, pp. 164 - 167, 1996年」及び従来技術文献8「政瀧ほか, "連続音声認識のための可変長連鎖統計言語モデル", 電子情報通信学会技術報告, SP95-73, pp. 1-6, 平成7年」参照。)等が数多く提案されている。しかしながら、これらの手法を用いても、精度の良い言語モデルを構築するためには、相当量のデータを用いる必要があると考えられる。

【0007】現在、実用化に向けて研究が行われている連続音声認識システムは、ホテル予約・スケジュール等、タスクを限定し、システムのパフォーマンスの向上させている物が多い。しかしながら、タスク毎に大量の言語データを集めるのは困難である。特に、日本語の場合は、英語等のように単語の区切りが明確ではなく、通常人間が手作業で単語の切り出し・形態素解析の作業を行うため、大量のデータを集めるのはさらに困難である。しかしながら、データ量を増やすために、他のタスクのデータを用いた場合、言語的特徴はタスク毎に異なるため、単純にデータを混合しても目的のタスク特有の言語特徴を効果的に表現することはできないと考えられる。

【0008】これらの問題を解決する手段として、言語

モデルのタスク適応を考えられている。すなわち、目的のタスク以外のデータも含めた大量のデータを学習することによりデータ量の問題を解決し、得られたモデルの言語特徴を目的のタスクに適応させる方法である。タスク適応の手法として、従来提案されているものには大量のデータで作成したN-gramと、目的タスクのデータで作成したN-gramとを重みづけにより混合する方法(例えば、従来技術文献9「伊藤ほか, "小量テキストによる言語モデルのタスク適応の検討", 日本音響学会講演論文集, 1-3-19, pp. 37-38, 平成8年9月」及び従来技術文献10「坂本ほか, "音声認識における統計的言語モデルの選択の効果", 日本音響学会講演論文集, 1-Q-24, pp. 157-158, 平成6年3月」参照。)がある。

【0009】例えば、従来技術文献9の手法を用いた第1の従来例のタスク適応化処理は、小量テキストに重みをかけて大量テキストと混合するものであり、次の手順によって言語モデルを作成する。

(a) 小量テキストを、重み付きで大量テキストに加える。重み係数を $w$ とすると、小量テキスト中で $m$ 回出現した単語は、大量テキスト中で $(w \cdot m)$ 回出現した単語と同等に扱われることになる。

(b) このようにしてできたテキストの中で、出現頻度が一定回数未満の単語を、未知語を表す記号に置き換える。すなわち、語彙の制限を行う。

(c) テキストから統計を取り、N-gramモデルを構築する。

【0010】しかしながら、第1の従来例のタスク適応化処理においては、重み係数 $w$ は1つのみしか使用していないので、言語モデルにおける遷移確率の予測精度はいまだ比較的低く、当該言語モデルを用いて音声認識をしたときの音声認識率は比較的低いという問題点があった。

【0011】この第1の従来例の問題点を解決するために、本発明者は、特開平10-198395号の特許出願(以下、第2の従来例という。)において、特定のタスクのN-gram言語モデルの精度を向上させるため、MAP推定(最大事後確率推定)によるタスク適応の手法を提案している。

【0012】

【発明が解決しようとする課題】しかしながら、第2の従来例の手法は単語列毎にタスク間のN-gram遷移確率の平均、及び分散を用いてパラメータ推定を行うため、テキスト全体があらかじめ複数のタスクに分割されている必要があり、単一のタスクのデータには適用できないという問題点があった。また、データ量が比較的多いタスクにおいては、タスク適応の効果が小さいという問題があった。

【0013】本発明の目的は以上の問題点を解決し、各タスクに対する適応効果が大きく、従来例に比較して遷

移確率の予測精度及び信頼性を改善することができる統計的言語モデルを生成することができる統計的言語モデル生成装置、及び、当該統計的言語モデルを用いて従来例に比較して高い音声認識率で音声認識することができる音声認識装置を提供することにある。

【0014】

【課題を解決するための手段】本発明に係る請求項1記載の統計的言語モデル生成装置は、複数のクラスタの統計的言語モデルを記憶する記憶手段と、所定の複数の発声音声文を含む学習用テキストデータに基づいて各発声音声文に対する統計的言語モデルを生成して、各発声音声文が各クラスタに対応するように、上記生成した統計的言語モデルを上記記憶手段に記憶する初期化手段と、上記学習用テキストデータの各発声音声文について、各クラスタにおける統計的言語モデルの文生成確率を計算して最大の文生成確率を有するクラスタを選択してその発声音声文を所属させるように上記記憶手段に記憶するクラスタ選択手段と、発声音声文が属するクラスタが変化したときに、各クラスタ毎に、上記クラスタ選択手段によって選択された発声音声文を用いて上記記憶手段に記憶された各統計的言語モデルを更新して、各クラスタに対応した統計的言語モデルを生成するモデル変更手段と、上記複数の発声音声文に属するクラスタが1文も変化しなくなるまで、上記クラスタ選択手段の処理と、上記モデル変更手段の処理を繰り返す制御手段とを備えたことを特徴とする。

【0015】また、請求項2記載の統計的言語モデル生成装置は、請求項1記載の統計的言語モデル生成装置において、さらに、上記記憶手段に記憶された各クラスタ毎のテキストデータに基づいて、最尤推定法を用いて各クラスタ毎に単語のN - g r a m ( N は 2 以上の自然数である。 ) の遷移確率を演算する第1の演算手段と、上記第1の演算手段によって演算された各クラスタ毎の単語のN - g r a m の遷移確率の出現分布を事前知識の所定の確率分布と仮定し、各クラスタ毎の確率分布の加重平均及び加重分散を演算した後、演算された加重平均と加重分散に基づいて事前知識の確率分布のパラメータを演算する第2の演算手段と、上記第2の演算手段によって演算された事前知識の確率分布のパラメータと、上記学習用テキストデータうちの特定クラスタのテキストデータの事後知識における処理対象の単語列の直前の単語列の出現回数と、処理対象の単語列の出現回数とに基づいて、各クラスタ毎の単語のN - g r a m の遷移確率を計算することにより、各クラスタ毎の単語のN - g r a m の遷移確率を含む統計的言語モデルを生成する第3の演算手段とを備えたことを特徴とする。

【0016】さらに、請求項3記載の統計的言語モデル生成装置は、請求項2記載の統計的言語モデル生成装置において、さらに、上記第3の演算手段によって演算された各クラスタ毎の単語のN - g r a m の遷移確率に基

づいて、所定の平滑化処理を実行し、処理後の各クラスタ毎の単語のN - g r a m の遷移確率を含む統計的言語モデルを生成する第1の生成手段を備えたことを特徴とする。

【0017】またさらに、請求項4記載の統計的言語モデル生成装置は、請求項3記載の統計的言語モデル生成装置において、さらに、上記学習用テキストデータに基づいて、最尤推定法を用いて単語のN - g r a m ( N は 2 以上の自然数である。 ) の遷移確率を演算して、上記単語のN - g r a m の遷移確率を含む別の統計的言語モデルを生成する第2の生成手段を備えたことを特徴とする。

【0018】本発明に係る請求項5記載の音声認識装置は、入力される発聲音声文の音声信号に基づいて、所定の統計的言語モデルを用いて音声認識する音声認識装置において、請求項4記載の統計的言語モデル生成装置と、上記第2の生成手段によって生成された別の統計的言語モデルを用いて、入力される発聲音声文の音声信号を音声認識して第1の認識仮説を出力する第1の音声認識手段と、上記第1の音声認識手段から出力される第1の認識仮説に回答して、上記第1の生成手段によって生成された各クラスタ毎の統計的言語モデルを用いて、入力される発聲音声文の音声信号を音声認識して、文生成確率が最大のクラスタの統計的言語モデル生成装置を選択するモデル選択手段と、上記モデル選択手段によって選択されたクラスタの統計的言語モデルを用いて、上記第1の音声認識手段から出力される第1の認識仮説に対して絞込処理を行って第2の認識仮説を生成して認識結果として出力する第2の音声認識手段とを備えたことを特徴とする。

【0019】

【発明の実施の形態】以下、図面を参照して本発明に係る実施形態について説明する。

【0020】図1に本発明に係る一実施形態の連続音声認識装置のブロック図を示す。本実施形態の連続音声認識装置は、図1において、特に、クラスタリング処理部40と、2つの言語モデル生成部41、42と、2つの単語仮説絞込部6a、6bを備えたことを特徴としている。本実施形態では、上述した第2の従来例の問題点を解決するためには、同一タスクの文でも、さまざまな内容の文が存在し、タスクという尺度よりも、文毎の内容で分類した方が言語的特徴がより明確になる考えられる。本実施形態では、これらの問題を解決し、さらに精度の高いN - g r a m 言語モデルを得るため、学習用テキストデータメモリ20内のテキストデータのコーパス全体をクラスタリング処理部40によって自動的にクラスタリングし、言語モデル生成部42において、MAP推定を用いてそれぞれのクラスタ毎にN - g r a m 言語モデルを構築する。また、精度を向上させるため、従来の単語N - g r a m に代り、可変長N - g r a m ( 品詞

と可変長単語列の複合 N - g r a m ) を用いる。

【 0 0 2 1 】すなわち、本実施形態の連続音声認識装置は、言語モデルの学習に用いるテキストコーパスをクラスタリングし、クラスタ毎の言語特徴を明確にさせ、言語モデルの精度を向上させる。しかしながら、入力された音声の発話文が属するクラスタをあらかじめ知ることは不可能である。このため、まず最初に、入力音声をコーパス全体で作成した言語モデルで認識を行い、次に、認識結果からクラスタ別の言語モデルを 1 つのみ選択し再度認識を行う、という 2 段階の認識を行う。

【 0 0 2 2 】中間認識結果 W からクラスタ別の言語モデル  $L M c$  の選択は、次式のように確率が最も高いものを選択することにより行う。

【 0 0 2 3 】

【数 2】

$$L M c = \underset{L M k}{\operatorname{argmax}} P ( L M k | W )$$

【 0 0 2 4 】上記式は、Bayes 則を用いると、次式のように表される。

【 0 0 2 5 】

【数 3】

$$L M c = \underset{L M k}{\operatorname{argmax}} P ( W | L M k ) P ( L M k )$$

【 0 0 2 6 】ここで、 $P ( L M k )$  は、言語モデル  $L M k$  の事前確率である。この確率は前発話の内容等より求めることができるが、本実施形態ではこの確率に関しては考慮しない。従って、次の式によりクラスタ言語モデルを選択する。

【 0 0 2 7 】

【数 4】

$$L M c = \underset{L M k}{\operatorname{argmax}} P ( W | L M k )$$

【 0 0 2 8 】すなわち、K 個のクラスタのそれぞれのモデル  $L M k$  で認識結果文 W に対する生成確率  $P ( W | L M k ) ( 1 \leq k \leq K )$  を求め、確率の最も高いクラスタモデル  $L M c$  を選択する。この選択処理は、言語モデル選択部 8 により行われる。

【 0 0 2 9 】次いで、クラスタリング処理部 4 0 によるコーパスのクラスタリングについて説明する。コーパスを自動クラスタリングするために、公知の K - m e a n s 法に類似した方法を用いた。K - m e a n s 法は、サンプルを距離が最も近いクラスタ中心に所属させる形でクラスタリングを行う手法である。この手法を文のクラスタリングに適用するため、次の 2 点で通常の方法と異なる。

( 1 ) クラスタ中心ベクトルをそのクラスタに属する文で生成される言語モデルとする。

( 2 ) 距離尺度に文の生成確率  $P ( W | L M k )$  を用いる。

【 0 0 3 0 】これらは、前述した認識結果からクラスタ

モデルの選択で用いる手法と同一であり、妥当な方法であると考えられる。以下に、クラスタリングの手順を示す。

<ステップ S S 1> クラスタモデルの初期化：クラスタ数を K とし、適当な手法によりコーパスから K 文を選択して全クラスタに 1 文ずつ配置し、クラスタ毎の言語モデル  $L M 1$ 、 $L M 2$ 、...、 $L M K$  を作成する。

<ステップ S S 2> クラスタの選択：コーパスの全文について、各クラスタにおける言語モデルの文生成確率を求め、最も確率の高いクラスタを選択し、その文を所属させる。

<ステップ S S 3> クラスタモデルの変更：各クラスタ毎に、ステップ S S 2 で選択した文を用いて言語モデル  $L M 1$ 、 $L M 2$ 、...、 $L M K$  を更新する。

<ステップ S S 4> 終了条件：文の属するクラスタが 1 文も変化しない場合、クラスタリングを終了する。それ以外の場合は、ステップ S S 2 及び S S 3 の処理を繰り返す。ただし、ある程度の回数を繰り返してもクラスタリングが収束しない場合は強制終了させる。

【 0 0 3 1 】次いで、MAP 推定による N - g r a m の適応について説明する。クラスタリングを行うことにより、クラスタ毎の言語的特徴は明確になるものの、クラスタ毎のデータ量は減少するため、N - g r a m のパラメータ推定の精度が低下することが考えられる。このため、第 2 の従来例で開示された MAP 推定を用いた適応の手法を用い、パラメータ推定の精度を向上させる。

【 0 0 3 2 】まず、MAP 推定法による遷移確率の算出について述べる。通常、N - g r a m の遷移確率は、ML (Maximum Likelihood; 最尤) 推定法により算出される。最尤推定法は、観測したサンプル値 (すなわち、テキストデータの単語)  $x$  に対して、遷移確率  $p$  が決まったときに単語  $x$  となる尤度関数  $f ( x | p )$  を最大にさせる値として、N - g r a m 遷移確率  $p_{ML}$  が次式で定められる。

【 0 0 3 3 】

【数 5】

$$p_{ML} = \underset{p}{\operatorname{argmax}} f ( x | p )$$

【 0 0 3 4 】ここで、関数  $\operatorname{argmax}$  は、 $p$  を変化したときに関数の引数が最大となるときの  $p$  の値を示す関数である。議論を簡単にするため、単語  $w_k$  から  $w_l$  への bigram の遷移確率  $p ( w_l | w_k )$  について考え、また、以下のような変数の定義を行う。

( a ) N : 学習用テキストデータ中の単語  $w_k$  の出現回数。

( b ) n : 学習用テキストデータ中の単語列  $w_k w_l$  の出現回数。

( c ) p : 単語  $w_k$  から  $w_l$  への遷移確率。

【 0 0 3 5 】このとき尤度関数  $f ( x | p )$  は、単語  $w_k$  が N 回観測され、次に単語  $w_l$  に続く回数が n 回で、それ以外の単語に続く回数が  $( N - n )$  回であるから、次

式を得ることができる。

【0036】

$$\text{【数6】 } f(p) = p^n (1-p)^{N-n}$$

【0037】  $f(p)$  の最大化条件  $d \log f(p) / dp = 0$  を解くことにより、N - gram の遷移確率は次式のように計算される。

【0038】

$$\text{【数7】 } p_{ML} = n / N$$

【0039】 従って、もし単語列  $w_k w_l$  が観測データ上で出現しない場合、 $n = 0$  であるから、遷移確率は 0 と推定されてしまう。これに対して、MAP (Maximum a posteriori Probability; 最大事後確率) 推定においては、最尤推定法を用いて、観測したサンプル値  $x$  に対して、遷移確率  $p$  が決定される事後確率関数  $h(p | x)$  を最大化する値として、N - gram の遷移確率が次式で求められる。

【0040】

$$\text{【数8】 } p_{MAP} = \underset{p}{\operatorname{argmax}} h(p | x)$$

【0041】 ここで、Bayes 則を用いると、上記数 8 は次式のように変形される。

【0042】

$$\text{【数9】 } p_{MAP} = \underset{p}{\operatorname{argmax}} f(x | p) g(p)$$

【0043】 ここで、 $g(p)$  は、各クラスタのテキストデータに基づいて予め決定される、N - gram の遷移確率  $p$  の事前分布である。すなわち、MAP 推定法を用いると、N - gram の遷移確率はある事前知識より

$$\begin{aligned} p_{MAP} &= \underset{p}{\operatorname{argmax}} \{ p^n (1-p)^{N-n} a p^{\alpha-1} (1-p)^{\beta-1} \} \\ &\equiv \underset{p}{\operatorname{argmax}} L(p) \end{aligned}$$

【0047】 ここで、関数  $L(p)$  が最大となるための条件  $d \log L(p) / dp = 0$  を  $p$  について解くと、単語の bigram の遷移確率  $p_{MAP}$  は次式のように求まる。

【0048】

$$\text{【数11】 } p_{MAP} = (n + \alpha - 1) / (N + \alpha + \beta - 2)$$

【0049】 ここで、パラメータ  $\alpha$  及び  $\beta$  は、事前分布であるベータ分布のパラメータであるが、これらは、次式のように求めることができる。なお、ベータ分布の平均  $\mu$  及び分散  $\sigma^2$  は以下の式となることが知られている (例えば、従来技術文献 5 参照。 )。

【0050】

$$\begin{aligned} \text{【数12】 } \mu &= \alpha / (\alpha + \beta) \\ \sigma^2 &= (\alpha \beta) / \{ (\alpha + \beta)^2 (\alpha + \beta + 1) \} \end{aligned}$$

\* 得られる分布  $g(p)$  に従う変数とし、この事前分布と実際に観測されたサンプル値とを用いて、実際の遷移確率が推定される。このため、観測データで出現しない単語遷移に対しても、事前知識により 0 でない遷移確率を与えることができる。

【0044】 次いで、bigram を例にとり、MAP 推定法により N - gram の遷移確率を求める方法を示す。ただし、変数の定義は上述と同じものを用いる。まず、遷移確率  $p$  の事前分布としてベータ分布  $(a p^{\alpha-1} (1-p)^{\beta-1})$  , ここで、 $\alpha$  及び  $\beta$  はベータ分布の正のパラメータであり、 $a$  は正規化のための正の定数である。)を用いる。なお、 $0 < p < 1$  である。ベータ分布を用いる理由は次の 2 点である。

(a) ベータ分布は 2 項分布の自然共役事前分布で、MAP 推定によるパラメータの解が求まりやすい。

(b) ベータ分布のパラメータ  $\alpha, \beta$  を変化させることにより、様々な形状の分布を表すことができる。

ここで、ベータ分布は、連続変数の確率分布の一種であり、ガンマ関数をもとにして構成されるベータ関数が表示に含まれる。なお、本実施形態においては、ベータ分布を用いるが、本発明はこれに限らず、ベータ分布に代えて、ディリクレ分布を用いてもよい。ディリクレ分布は、ベータ分布を多変量分布に拡張したものであり、多項分布の自然共役事前分布である。

【0045】 上記数 9 の MAP 推定法の定義に従うと、遷移確率  $p_{MAP}$  は、尤度関数  $f(p)$  と事前分布  $g(p)$  とを用いて次式のように求められる。

【0046】

【数10】

【0051】 これらの式を  $\mu, \sigma^2$  について解くと、次式が得られる。

【0052】

$$\begin{aligned} \text{【数13】 } \mu &= \{ \mu^2 (1-\mu) \} / \{ \mu^2 - \mu \} \\ \sigma^2 &= \{ \mu (1-\mu) \} / \{ \mu^2 - 1 \} \end{aligned}$$

【0053】 以上より、観測テキストデータから頻度を計算することにより得られるパラメータ  $N, n$ 、及び事前分布の平均  $\mu$  及び分散  $\sigma^2$  により、上記数 11 及び数 13 を用いて、単語の bigram の遷移確率を求めることができる。

【0054】 これまでの議論は、単語の bigram についてのみの議論であったが、一般に、MAP 推定法による N - gram の遷移確率  $p(w_n | w_1^{n-1})$  は、直前の単語  $w_k$  を直前の単語列  $w_1^{n-1}$  と置き換え、パラメータ  $N$  及び  $n$  を次のように定義すれば、同じ議論が通用

10

20

40

50

することは明らかである。

( a ) N : 学習用テキストデータ中の単語列  $w_1^{n-1}$  の出現回数 (  $c ( w_1^{n-1} )$  )、すなわち、処理対象の単語列の直前の単語列の出現回数である。

( b ) n : 学習用テキストデータ中の単語列  $w_1^n$  の出現回数 (  $c ( w_1^n )$  )、すなわち、処理対象の単語列の出現回数である。

【 0 0 5 5 】次いで、MAP 推定法を用いたクラスタ適応化処理について述べる。上述のMAP 推定法による N - g r a m をクラスタ適応化に応用するために、図 3 に示すように、複数のクラスタより構成される大量のテキストデータから構成される不特定のタスクのテキストデータに基づく N - g r a m を事前知識とし、目的の特定クラスタ i のテキストデータ  $2 1 - i$  を事後知識とみな\*

$$\sigma^2 = \sum_i c_i ( w_1^{n-1} ) p_i ( w_n | w_1^{n-1} )^2 / \sum_i c_i ( w_1^{n-1} ) - \mu^2$$

【 0 0 5 7 】ここで、 $c_i ( w_1^{n-1} )$  はクラスタ i において単語列  $w_1^{n-1}$  の出現頻度であり、 $p_i ( w_n | w_1^{n-1} )$  はクラスタ i における単語列  $w_1^{n-1}$  から  $w_n$  への遷移確率である。また、事後知識を目的のクラスタのテキストデータとすると、前述のパラメータ N 及び n は次のように表される。

( a ) N : 目的の特定クラスタ i のテキストデータ  $2 1 - i$  中の単語列  $w_1^{n-1}$  の出現頻度、すなわち、処理対象の単語列の直前の単語列の出現回数である。

( b ) n : 目的の特定クラスタ i のテキストデータ  $2 1 - i$  中の単語列  $w_1^n$  の出現頻度、すなわち、処理対象の単語列の出現回数である。

以上の加重平均  $\mu$ 、加重分散  $\sigma^2$ 、パラメータ n 及び N を上述の数 1 0 及び数 1 2 に代入することにより、MAP 推定法によるタスク適応後の N - g r a m 遷移確率が得られる。

【 0 0 5 8 】さらに、Back - o f f 平滑化法による遷移確率の平滑化について述べる。上記でMAP 推定法によるタスク適応の基本原理解を述べたが、実際に言語モデルとして使用するには、2 つの問題がある。1 つは、平滑化の問題である。不特定タスクの大量のテキストデータを用いても、出現しない単語列が存在し、MAP 推定法を用いても、N - g r a m の遷移確率が 0 となってしまう。従って、平滑化処理によりテキストに出現しな

$$\begin{aligned} & P s ( w_n | w_1^{n-1} ) \\ & = P h ( w_n | w_1^{n-1} ) , c_i ( w_1^{n-1} ) > 0 \text{ のとき} \\ & = ( w_1^{n-1} ) P s ( w_n | w_2^{n-1} ) , c_i ( w_1^{n-1} ) = 0 , c_i ( w_2^{n-1} ) > 0 \text{ のとき} \\ & = P s ( w_n | w_2^{n-1} ) , c_i ( w_1^{n-1} ) = 0 , c_i ( w_2^{n-1} ) = 0 \text{ のとき} \end{aligned}$$

【 0 0 6 1 】上記の数 1 6 において、Ph はクラスタ適応化により得られる確率に軽減係数をかけたものであり、次式で与えられる。

【 0 0 6 2 】

$$\text{【数 1 7】 } P h ( w_n | w_1^{n-1} ) = \{ c_i ( w_1^n ) + 1 \} / \{ c_i ( w_1^n ) \} \times \{ n c_i ( w_1^n ) + 1 \} / \{ n$$

\* す。不特定のクラスタの N - g r a m を事前知識とみなしたとき、その事前分布は、各クラスタにおける N - g r a m 遷移確率の分布と考えることができる。ただし、各クラスタにおける N - g r a m 遷移確率は最尤推定法により求められる。この事前分布をベータ分布と仮定してMAP 推定法の前分布として用いる。このとき、前分布の加重平均  $\mu$ 、及び加重分散  $\sigma^2$  は次式で求められる。

【 0 0 5 6 】

$$\text{【数 1 4】 } \mu = \sum_i c_i ( w_n ) p_i ( w_n | w_1^{n-1} ) / \sum_i c_i ( w_1^{n-1} )$$

【数 1 5】

い単語組に対しても、0 でない遷移確率を与える必要がある。もう 1 つの問題は、本発明に係るタスク適応化処理は、全ての遷移確率を独立に求める手法であるため、遷移確率の和が 1 になるとは限らない。連続音声認識等に適用する際は、問題とはならないが、パープレキシティで評価する際は、1 に正規化されていないと、正しい評価ができない。従って、近年盛んに用いられている B a c k - O f f 平滑化法 (例えば、従来技術文献 4 参照。) を拡張して、これらの問題を解決する方法を述べる。

【 0 0 5 9 】単語列  $w_1^n$  が不特定のクラスタのテキストデータ  $2 1 - k$  に含まれる場合は、上記のタスク適応化処理により、遷移確率  $p_{MAP} ( w_n | w_1^{n-1} )$  を求め、チューリング ( T u r i n g ) 推定法により、確率  $p_{MAP} ( w_n | w_1^{n-1} )$  を軽減する。ただし、軽減係数は不特定のクラスタのテキストデータ  $3 1$  の頻度 (  $c_i ( w_1^n )$  ) を用いて計算する。当該軽減により生じた確率の余剰分を  $w_1^n$  が不特定のクラスタのテキストデータ  $3 1$  に含まれない単語連鎖に対して、( n - 1 ) - g r a m の遷移確率に比例して配分する。以上をまとめると、クラスタ適応化された N - g r a m の平滑化後の遷移確率  $P s ( w_n | w_1^{n-1} )$  は次式で表される。

【 0 0 6 0 】

$$\text{【数 1 6】 } P s ( w_n | w_1^{n-1} ) = \{ c_i ( w_1^n ) \} \cdot p_{MAP} ( w_1^n )$$

【 0 0 6 3 】ここで、 $n_c$  は、不特定のクラスタのテキストデータ  $3 1$  中に c 回出現する単語列の種類数 (異なり) であり、また、数 1 6 で、 $( w_1^{n-1} )$  は正規化のための係数であり、次のように求められる。

【 0 0 6 4 】

17

【数18】  $(w_1^{n-1}) = Aa / Ab$ 

ここで、

$$Aa \equiv 1 - \sum_{w_n : c_1(w_1^n) > 0} Ph(w_n | w_1^{n-1})$$

$$w_n : c_1(w_1^n) > 0$$

$$Ab \equiv 1 - \sum_{w_n : c_1(w_2^n) > 0} Ph(w_n | w_1^{n-1})$$

$$w_n : c_1(w_2^n) > 0$$

【0065】以上のBack-off平滑化法を応用した手法を用いることにより、学習データ上に出現しない単語連鎖に対しても確率値を与えることができる。また、遷移確率 $p_{MAP}$ が正規化されていなくても、上記数

【0066】従って、本実施形態で用いるMAP推定法によるN-gramの適応手法について要約すると、以下の通りとなる。MAP推定法による単語列 $h$ から次単語 $w$ への単語N-gramの遷移確率 $P(w|h)$ は次式により与えられる。

【0067】

$$\text{【数19】 } P(w|h) = \{N(h, w) + \dots - 1\} / (N(h) + \dots - 2)$$

【0068】ここで、 $N(\#)$ はそのクラスタでの単語(列) $\#$ の出現頻度である。また、 $\dots$ 及び $\dots$ は事前分布として用いるベータ分布 $(\alpha p^{-1} (1-p)^{-1}, \alpha)$ のパラメータであり、次式により求められる。

【0069】

$$\text{【数20】 } \dots = \{ \mu^2 (1 - \mu) / \dots^2 \} - \mu$$

$$\text{【数21】 } \dots + \dots = \{ \mu (1 - \mu) / \dots^2 \} - 1$$

【0070】上式の $\mu$ 及び $\dots^2$ は、クラスタ毎の遷移確率 $P(w|h)$ の分布の平均及び分散である。また、本実施形態で用いる可変長N-gramは、クラスN-gramを基本としたモデルであり、遷移確率は $P(ws|c(ws)) \cdot P(c(ws)|c(h))$ として与えられる。ただし、 $ws$ は可変長の単語列で、 $c(\#)$ は単語(列) $\#$ の属するクラスである。 $P(c(ws)|c(h))$ はクラス間の遷移確率であり、上記数19と同様に与えることができる。また、 $P(ws|c(ws))$ はその単語の属するクラスから単語の出現確率であり、MAP推定により次式で与えられる。

【0071】

$$\text{【数22】 } P(ws|c(ws)) = \{N(ws) + \dots - 1\} / \{N(c(ws)) + \dots - 2\}$$

【0072】また、公知のBack-off平滑化法を用い、コーパス上に出現しなかった単語遷移に対して確率を与えるとともに、遷移確率の和が1になるよう確率の正規化を行う。

【0073】図4は、図1のクラスタリング処理部40によって実行されるクラスタリング処理を示すフローチャートである。図4において、まず、ステップS1において学習用テキストデータメモリ20からK個の発声音

18

選択した発声音声を、学習用テキストデータメモリ21の各メモリ21-1乃至21-Kに、クラスタ1からクラスタMへの順番に1文ずつ書き込む。次いで、ステップS3において学習用テキストデータメモリ21の各クラスタのテキストデータを読み出し、ステップS4において上記読み出した各クラスタのテキストデータから、各クラスタ毎に統計的言語モデルを生成する。ここで、生成された統計的言語モデルはクラスタリング処理部40の内部メモリ又は統計的言語モデルメモリ32に記憶される。

【0074】さらに、ステップS5において学習用テキストデータメモリ20から1文ずつ読み出し、ステップS6においてステップS4で生成した各統計的言語モデルに対して、ステップS5で読み出した文の生成確率を計算し、確率の最も高いクラスタCを選択し、ステップS7においてステップS5で読み出した文を、学習用テキストデータメモリ21のメモリ21-CにクラスタCとして書き込む。そして、ステップS8においてステップS5で読み出した文は最後の文か否かが判断され、NOであるときは、次の文を処理するために、ステップS5に戻る。一方、ステップS8でYESのときは、ステップS9において、ステップS6で選択されたクラスタCが1文でも変化したか否かが判断され、YESのときは再度クラスタリング処理を実行するために、ステップS3に戻る。ステップS9でNOであるときは、当該クラスタリング処理を終了する。

【0075】図5は、図1の言語モデル生成部41によって実行される言語モデル生成処理を示すフローチャートである。図5において、まず、ステップS11において学習用テキストデータメモリ20からコーパスのテキストデータを読み出し、ステップS12において読み出したテキストデータに基づいて最尤推定法を用いて単語bigramの遷移確率を数6を用いて計算する。次いで、ステップS13において計算された単語bigramの遷移確率を含む統計的言語モデルを生成して、統計的言語モデルメモリ31に記憶して当該言語モデル生成処理を終了する。

【0076】図6は、図1の言語モデル生成部42によって実行される言語モデル生成処理を示すフローチャートである。図6において、まず、ステップS21において学習用テキストデータメモリ20から各クラスタkのテキストデータ21-k( $k=1, 2, \dots, K$ )を読み出す。次いで、ステップS22において、読み出した各クラスタkのテキストデータ21-kに基づいて最尤推定法を用いて各クラスタk毎に単語bigramの遷移確率を数6を用いて計算し、ステップS23において各クラスタkの単語bigramの遷移確率の出現頻度分布をベータ分布と仮定し、ベータ分布の加重平均 $\mu$ 及び加重分散 $\dots^2$ を数13及び数14を用いて計算した後これらに基づいて数12を用いてパラメータ $\dots$ 及び $\dots$ を計



算する。さらに、ステップ S 2 4 において事前知識のパラメータ及び  $t$ 、特定クラスタのテキストデータ  $21-i$  の事後知識のパラメータ  $N$  及び  $n$  とに基づいて数 10 を用いて各クラスタ  $k$  毎の単語  $bigram$  の遷移確率  $p$  を計算する。ここで、上記ステップ S 2 1 から S 2 4 までの処理は、すべてのクラスタ  $k = 1, 2, \dots, K$  について実行される。さらに、ステップ S 2 5 において各クラスタ  $k$  毎の単語  $bigram$  の遷移確率  $p$  に基づいて Back-off 平滑化処理を実行し、処理後の各クラスタ  $k$  毎の単語  $bigram$  の遷移確率を含む統計的言語モデルを生成して、クラスタ適応化された統計的言語モデルメモリ 3 2 に記憶して、当該言語モデル生成処理を終了する。

【0077】次いで、図 1 に示す連続音声認識装置の構成及び動作について説明する。図 1 において、単語照合部 4 に接続された音素隠れマルコフモデル（以下、隠れマルコフモデルを HMM という。）メモリ 1 1 内の音素 HMM は、各状態を含んで表され、各状態はそれぞれ以下の情報を有する。

( a ) 状態番号、( b ) 受理可能なコンテキストクラスタ、( c ) 先行状態、及び後続状態のリスト、( d ) 出力確率密度分布のパラメータ、及び ( e ) 自己遷移確率及び後続状態への遷移確率。

なお、本実施形態において用いる音素 HMM は、各分布がどの話者に由来するかを特定する必要があるため、所定の話者混合 HMM を変換して生成する。ここで、出力確率密度関数は 3 4 次元の対角共分散行列をもつ混合ガウス分布である。また、単語照合部 4 に接続された単語辞書メモリ 1 2 内の単語辞書は、音素 HMM メモリ 1 1 内の音素 HMM の各単語毎にシンボルで表した読みを示すシンボル列を格納する。

【0078】図 1 において、話者の発声音声はマイクロホン 1 に入力されて音声信号に変換された後、特徴抽出部 2 に入力される。特徴抽出部 2 は、入力された音声信号を A/D 変換した後、例えば LPC 分析を実行し、対数パワー、16 次ケプストラム係数、対数パワー及び 16 次ケプストラム係数を含む 3 4 次元の特徴パラメータを抽出する。抽出された特徴パラメータの時系列はバッファメモリ 3 を介して単語照合部 4 に入力される。

【0079】単語照合部 4 は、ワン・パス・ピタビ復号化法を用いて、バッファメモリ 3 を介して入力される特徴パラメータのデータに基づいて、音素 HMM 1 1 と単語辞書 1 2 とを用いて単語仮説を検出し尤度を計算して出力する。ここで、単語照合部 4 は、各時刻の各 HMM の状態毎に、単語内の尤度と発声開始からの尤度を計算する。尤度は、単語の識別番号、単語の開始時刻、先行単語の違い毎に個別にもつ。また、計算処理量の削減のために、音素 HMM 1 1 及び単語辞書 1 2 とに基づいて計算される総尤度のうちの低い尤度のグリッド仮説を削減する。単語照合部 4 は、その結果の単語仮説と尤度の

情報を発声開始時刻からの時間情報（具体的には、例えばフレーム番号）とともにバッファメモリ 5 を介して単語仮説絞込部 6 a に出力する。

【0080】単語仮説絞込部 6 a は、単語照合部 4 からバッファメモリ 5 を介して出力される単語仮説に基づいて、統計的言語モデルメモリ 3 2 内の統計的言語モデルを参照して、終了時刻が等しく開始時刻が異なる同一の単語の単語仮説に対して、当該単語の先頭音素環境毎に、発声開始時刻から当該単語の終了時刻に至る計算された総尤度のうちの最も高い尤度を有する 1 つの単語仮説で代表させるように単語仮説の絞り込みを行った後、絞り込み後のすべての単語仮説の単語列のうち、最大の総尤度を有する仮説の単語列を認識結果としてバッファメモリ 7 を介して言語モデル選択部 8 に出力する。本実施形態においては、好ましくは、処理すべき当該単語の先頭音素環境とは、当該単語より先行する単語仮説の最終音素と、当該単語の単語仮説の最初の 2 つの音素とを含む 3 つの音素並びをいう。

【0081】単語仮説絞込部 6 a の処理においては、例えば、図 2 に示すように、(  $i-1$  ) 番目の単語  $W_{i-1}$  の次に、音素列  $a_1, a_2, \dots, a_n$  からなる  $i$  番目の単語  $W_i$  がくるときに、単語  $W_{i-1}$  の単語仮説として 6 つの仮説  $W_a, W_b, W_c, W_d, W_e, W_f$  が存在している。ここで、前者 3 つの単語仮説  $W_a, W_b, W_c$  の最終音素は  $/x/$  であるとし、後者 3 つの単語仮説  $W_d, W_e, W_f$  の最終音素は  $/y/$  であるとする。終了時刻  $t$  と先頭音素環境が等しい仮説（図 2 では先頭音素環境が “  $x/a_1/a_2$  ” である上から 3 つの単語仮説）のうち総尤度が最も高い仮説（例えば、図 2 において 1 番上の仮説）以外を削除する。なお、上から 4 番めの仮説は先頭音素環境が違うため、すなわち、先行する単語仮説の最終音素が  $x$  ではなく  $y$  であるので、上から 4 番めの仮説を削除しない。すなわち、先行する単語仮説の最終音素毎に 1 つのみ仮説を残す。図 2 の例では、最終音素  $/x/$  に対して 1 つの仮説を残し、最終音素  $/y/$  に対して 1 つの仮説を残す。

【0082】次いで、言語モデル選択部 8 は、上述のように、数 4 に従ってクラスタの統計的言語モデルを統計的言語モデルメモリ 3 2 から選択し、すなわち、 $K$  個のクラスタのそれぞれのモデル  $LM_k$  で認識結果文  $W$  に対する生成確率  $P(W | LM_k)$  (  $1 \leq k \leq K$  ) を求め、確率の最も高いクラスタモデル  $LM_c$  を選択して、その選択情報を単語仮説絞込部 6 b に出力する。これに回答して、単語仮説絞込部 6 b は、単語仮説絞込部 6 a によって絞り込まれた単語仮説に対して、再度、統計的言語モデルメモリ 3 2 で選択された統計的言語モデルを用いて、単語仮説絞込部 6 a と同様の処理を実行して、単語仮説の絞り込み処理を実行して、処理後の例えば最尤の絞り込んだ単語仮説（ここで、 $n-best$  でもよい。）を認識結果として出力する。

【0083】以上の実施形態においては、当該単語の先頭音素環境とは、当該単語より先行する単語仮説の最終音素と、当該単語の単語仮説の最初の2つの音素とを含む3つの音素並びとして定義されているが、本発明はこれに限らず、先行する単語仮説の最終音素と、最終音素と連続する先行する単語仮説の少なくとも1つの音素とを含む先行単語仮説の音素列と、当該単語の単語仮説の最初の音素を含む音素列とを含む音素並びとしてもよい。

【0084】以上の実施形態において、特徴抽出部2と、単語照合部4と、単語仮説絞込部6a、6bと、クラスタリング処理部40と、言語モデル生成部41、42とは、例えば、デジタル電子計算機などのコンピュータで構成され、バッファメモリ3、5と、音素HMMメモリ11と、単語辞書メモリ12と、学習用テキストデータメモリ20、21と、統計的言語モデルメモリ31、32とは、例えばハードディスクメモリなどの記憶装置で構成される。

【0085】以上実施形態においては、単語照合部4と単語仮説絞込部6a、6bとを用いて音声認識を行っているが、本発明はこれに限らず、例えば、音素HMM11を参照する音素照合部と、例えばOne Pass DPアルゴリズムを用いて統計的言語モデルを参照して単語の音声認識を行う音声認識部とで構成してもよい。ただし、本実施形態の場合、統計的言語モデルメモリ31を参照して音声認識する第1の音声認識部と、統計的言語モデルメモリ32内で言語モデル選択部8によって\*

パープレキシティによる比較

全体モデル	クラスタモデル(クラスタ数)				
	4	8	16	32	64
14.21	13.00	12.33	11.44	10.44	9.72

【0089】ここで、パープレキシティとは以下のように定義される。例えば、複数n個の単語からなる長い単語列  $w_1^n = w_1 w_2 \dots w_n$  があるときのエントロピー  $H(n)$  は次式で表される。

【0090】

【数23】

$$H(n) = -(1/n) \cdot \log_2 P(w_1^n)$$

【0091】ここで、 $P(w_1^n)$  は単語列  $w_1^n$  の生成確率であり、パープレキシティ  $PP(n)$  は次式で表される。

【0092】

$$【数24】 PP(n) = 2^{H(n)}$$

【0093】上記表1より、クラスタ数に比例してパープレキシティが減少しており、クラスタ毎の言語的特徴がよりできたと考えられる。クラスタ数が64の時は、全体モデルよりもパープレキシティが約32%減少し

\* 選択された1つのクラスタの統計的言語モデルを参照して音声認識する第2の音声認識部とを備えることになる。

【0086】

【実施例】本発明者は、本実施形態で用いるタスク適応化された統計的言語モデルの性能を確認するため、評価実験を行った。実験で用いたデータは、本特許出願人が所有する自然発話データベース(例えば、従来技術文献11「T. Morimoto et al., "A Speech and Language Database for Speech Translation Research", ICSLP, pp. 1791-1794, 1994年」参照。)であり、本データベースのサイズは、1,332対話、32,074文、597,626単語で、語彙は7,221語である。このうち評価用として「ホテルの部屋の予約」タスクから40対話、1166文、18,381単語を選択し、残りのデータを言語モデルの学習に使用した。

【0087】最初にテストセットパープレキシティにより評価を行った。可変長N-gramは活用形及び活用型を含む158品詞による初期クラスから、500クラス分離を行ったモデルを使用した。クラスタ数4、8、16、32、64の時のクラスタモデルと、データベース全体で作成したモデル(クラスタ数1)とのパープレキシティの比較を表1に示す。

【0088】

【表1】

た。また、評価に用いた「ホテルの部屋の予約」タスクのデータは、データ量が多いために第2の従来例では、タスク適応の効果は、単語bigramで5%程度と小さかったが、本実施形態に係る装置では、文の内容毎に適応モデルを作成するため、大きな精度向上が得られたと考えられる。計算量の都合のため、クラスタ数は最大64としたが、さらにクラスタ数を増加させることにより、パープレキシティは減少すると考えられる。ただし、クラスタ数を多くしすぎると各クラスタのデータ量が少なくなりすぎ、パラメータ推定が困難になるため、限界はあると考えられる。

【0094】次に、連続音声認識に適用した際の認識率によって評価を行った。音響モデルにはML-SSS法(従来技術文献12「M. Ostendorf et al., "HMM topology design using maximum likelihoods

10

20

40

50

uccessive state splitting", Computer Speech and Language, No. 11, pp. 17 - 41, 1997年」参照。)によるHMM網(801状態5混合分布)の不特定話者モデルを用い、単語グラフサーチ法(従来技術文献13「清水ほか, "単語グラフを用いた自由発話音声認識", 電子情報通信学会研究報告, SP95-88, pp. 49-54, 1995年12月」参\*

連続音声認識における性能比較

認識率の種類	全体モデル	クラスタモデル(クラスタ数)		
		4	16	64
単語認識率	77.66	78.69	79.06	78.54
文認識率	33.43	35.82	36.12	37.31

【0096】上記表2より、単語認識率はクラスタ数16の時に全体モデルより約1.4%向上(改善率約6%)し、文認識率はクラスタ数64の時に最大約3.9%向上(改善率約6%)し、連続音声認識における有効性を確認した。クラスタ数64の時の単語認識率はクラスタ数4、16の時よりも低下しているが、これは、誤認識が生じた際にクラスタモデルの選択が正しく行われないことが原因と考えられる。

【0097】以上説明したように、本実施形態によれば、コーパスの各文をクラスタリングし、それぞれのクラスタ毎にMAP推定によるN-gram型の言語モデルを作成することにより言語特徴をより効果的に表現できる手法を開示している。実験の結果、パープレキシティは最大約32%減少し、また、連続音声認識に適用した際、単語認識率及び文認識率共に最大約6%改善し、本手法の有効性を確認した。すなわち、本実施形態によれば、少量のテキストデータを用いて、従来例に比較して、より高い遷移確率の予測精度及び信頼性を有する統計的言語モデルを生成することができるとともに、タスク選択を自動的に行うことができ、選択された統計的言語モデルを用いて音声認識することにより、従来例に比較してより高い音声認識率で連続的に音声認識することができる。

【0098】以上の実施形態において、統計的言語モデルは、N-gramの言語モデルを含むが、ここで、Nは2及び3に限らず、4以上の自然数であってもよい。

【0099】

【発明の効果】以上詳述したように本発明に係る請求項1記載の統計的言語モデル生成装置によれば、複数のクラスタの統計的言語モデルを記憶する記憶手段と、所定の複数の発声音声文を含む学習用テキストデータに基づいて各発声音声文に対する統計的言語モデルを生成して、各発声音声文が各クラスタに対応するように、上記

\*照。)により認識解の探索を行った。言語モデルは、コーパス全体で作成したモデルとクラスタ数4、16、64のクラスタモデルとを比較した。表2に単語認識率(Accuracy)(%)及び文認識率(%)を示す。

【0095】

【表2】

20

30

40

50

生成した統計的言語モデルを上記記憶手段に記憶する初期化手段と、上記学習用テキストデータの各発声音声文について、各クラスタにおける統計的言語モデルの文生成確率を計算して最大の文生成確率を有するクラスタを選択してその発声音声文を所属させるように上記記憶手段に記憶するクラスタ選択手段と、発声音声文が属するクラスタが変化したときに、各クラスタ毎に、上記クラスタ選択手段によって選択された発声音声文を用いて上記記憶手段に記憶された各統計的言語モデルを更新して、各クラスタに対応した統計的言語モデルを生成するモデル変更手段と、上記複数の発声音声文に属するクラスタが1文も変化しなくなるまで、上記クラスタ選択手段の処理と、上記モデル変更手段の処理を繰り返す制御手段とを備える。従って、少量のテキストデータを用いて、従来例に比較して、より高い遷移確率の予測精度及び信頼性を有する統計的言語モデルを生成することができる。

【0100】また、請求項2記載の統計的言語モデル生成装置によれば、請求項1記載の統計的言語モデル生成装置において、さらに、上記記憶手段に記憶された各クラスタ毎のテキストデータに基づいて、最尤推定法を用いて各クラスタ毎に単語のN-gram(Nは2以上の自然数である。)の遷移確率を演算する第1の演算手段と、上記第1の演算手段によって演算された各クラスタ毎の単語のN-gramの遷移確率の出現分布を事前知識の所定の確率分布と仮定し、各クラスタ毎の確率分布の加重平均及び加重分散を演算した後、演算された加重平均と加重分散に基づいて事前知識の確率分布のパラメータを演算する第2の演算手段と、上記第2の演算手段によって演算された事前知識の確率分布のパラメータと、上記学習用テキストデータうちの特定クラスタのテキストデータの事後知識における処理対象の単語列の直前の単語列の出現回数と、処理対象の単語列の出現回数

とに基づいて、各クラスタ毎の単語の N - g r a m の遷移確率を計算することにより、各クラスタ毎の単語の N - g r a m の遷移確率を含む統計的言語モデルを生成する第 3 の演算手段とを備える。従って、少量のテキストデータを用いて、従来例に比較して、より高い遷移確率の予測精度及び信頼性を有する統計的言語モデルを生成することができる。

【0101】さらに、請求項 3 記載の統計的言語モデル生成装置によれば、請求項 2 記載の統計的言語モデル生成装置において、さらに、上記第 3 の演算手段によって演算された各クラスタ毎の単語の N - g r a m の遷移確率に基づいて、所定の平滑化処理を実行し、処理後の各クラスタ毎の単語の N - g r a m の遷移確率を含む統計的言語モデルを生成する第 1 の生成手段を備える。従って、少量のテキストデータを用いて、従来例に比較して、より高い遷移確率の予測精度及び信頼性を有する統計的言語モデルを生成することができる。

【0102】またさらに、請求項 4 記載の統計的言語モデル生成装置によれば、請求項 3 記載の統計的言語モデル生成装置において、さらに、上記学習用テキストデータに基づいて、最尤推定法を用いて単語の N - g r a m ( N は 2 以上の自然数である。 ) の遷移確率を演算して、上記単語の N - g r a m の遷移確率を含む別の統計的言語モデルを生成する第 2 の生成手段を備える。従って、少量のテキストデータを用いて、従来例に比較して、より高い遷移確率の予測精度及び信頼性を有する統計的言語モデルを生成することができる。

【0103】本発明に係る請求項 5 記載の音声認識装置によれば、入力される発声音声文の音声信号に基づいて、所定の統計的言語モデルを用いて音声認識する音声認識装置において、請求項 4 記載の統計的言語モデル生成装置と、上記第 2 の生成手段によって生成された別の統計的言語モデルを用いて、入力される発聲音声文の音声信号を音声認識して第 1 の認識仮説を出力する第 1 の音声認識手段と、上記第 1 の音声認識手段から出力される第 1 の認識仮説に回答して、上記第 1 の生成手段によって生成された各クラスタ毎の統計的言語モデルを用いて、入力される発聲音声文の音声信号を音声認識して、文生成確率が最大のクラスタの統計的言語モデル生成装置を選択するモデル選択手段と、上記モデル選択手段によって選択されたクラスタの統計的言語モデルを用いて、上記第 1 の音声認識手段から出力される第 1 の認識仮説に対して絞込処理を行って第 2 の認識仮説を生成して認識結果として出力する第 2 の音声認識手段とを備える。従って、少量のテキストデータを用いて、従来例に比較して、より高い遷移確率の予測精度及び信頼性を有する統計的言語モデルを生成できるとともに、タスク選択を自動的に行うことができ、選択された統計的言語モデルを用いて音声認識することにより、従

来例に比較してより高い音声認識率で連続的に音声認識することができる。

【図面の簡単な説明】

【図 1】 本発明に係る一実施形態である連続音声認識装置のブロック図である。

【図 2】 図 1 の連続音声認識装置における単語仮説絞込部 6 a 及び 6 b の処理を示すタイミングチャートである。

【図 3】 図 1 の言語モデル生成部 4 2 の処理を示すブロック図である。

【図 4】 図 1 のクラスタリング処理部 4 0 によって実行されるクラスタリング処理を示すフローチャートである。

【図 5】 図 1 の言語モデル生成部 4 1 によって実行される言語モデル生成処理を示すフローチャートである。

【図 6】 図 1 の言語モデル生成部 4 2 によって実行される言語モデル生成処理を示すフローチャートである。

【符号の説明】

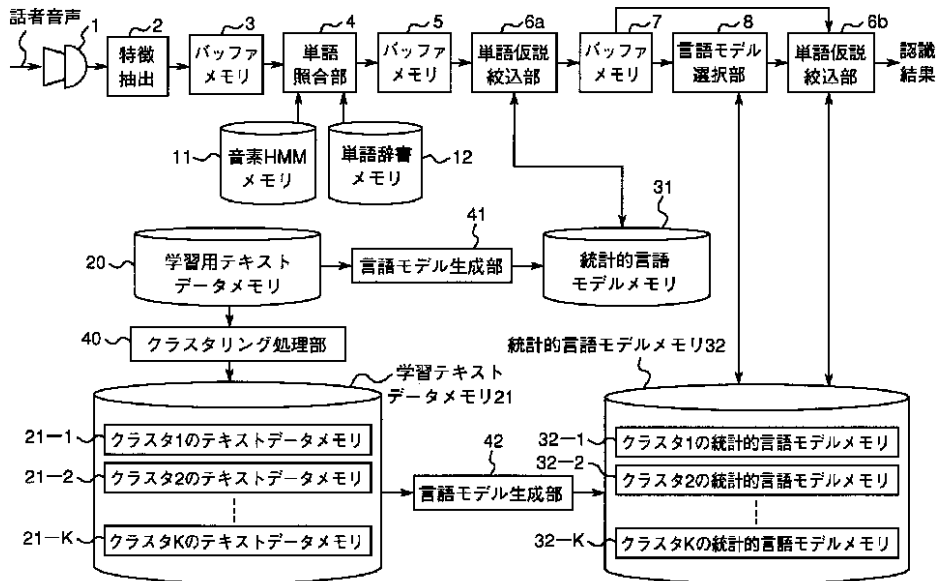
- 1...マイクロホン、
- 2...特徴抽出部、
- 3, 5, 7...バッファメモリ、
- 4...単語照合部、
- 6 a, 6 b...単語仮説絞込部、
- 8...言語モデル選択部、
- 11...音素 HMM メモリ、
- 12...単語辞書メモリ、
- 20, 21...学習用テキストデータメモリ、
- 21 - k...各クラスタのテキストデータメモリ、
- 31, 32...統計的言語モデルメモリ、
- 32 - k...各クラスタの統計的言語モデルメモリ、
- 41, 42...言語モデル生成部。

【要約】

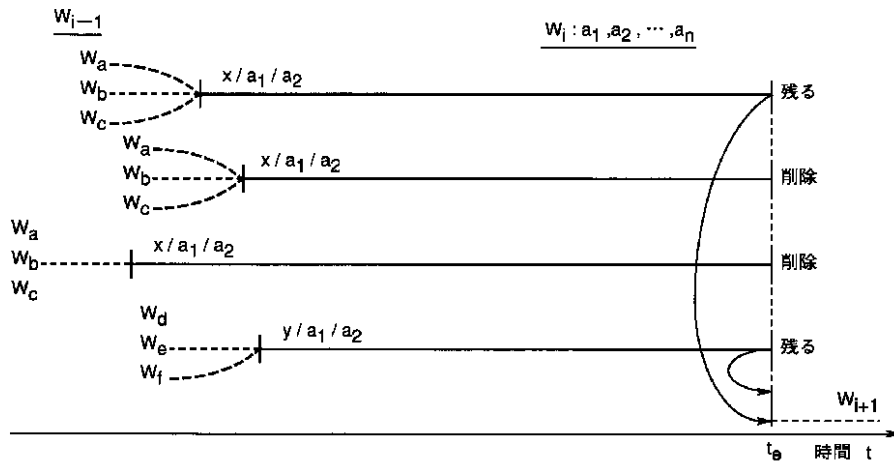
【課題】 遷移確率の予測精度及び信頼性を改善可能な統計的言語モデルを生成し、統計的言語モデルを用いてより高い音声認識率で音声認識する。

【解決手段】 学習用テキストデータ全体をクラスタリング処理部 4 0 によって自動的にクラスタリングしてクラスタ毎のテキストデータをメモリ 2 1 に記憶し、言語モデル生成部 4 2 により M A P 推定法を用いて各クラスタ毎の統計的言語モデルを生成してメモリ 3 2 に記憶する。一方、学習用テキストデータ全体に対して統計的言語モデルを生成してメモリ 3 1 に記憶する。単語照合部 4 による単語仮説の生成の後、単語仮説絞込部 6 a はメモリ 3 1 内の統計的言語モデルを用いて単語仮説の絞込処理を実行した後、言語モデル選択部 8 はメモリ 3 2 内の各クラスタの統計的言語モデルのうちで文生成確率が最大のモデルを選択して、単語仮説絞込部 6 b は選択されたモデルを用いて再度の絞込処理を行って認識結果を出力する。

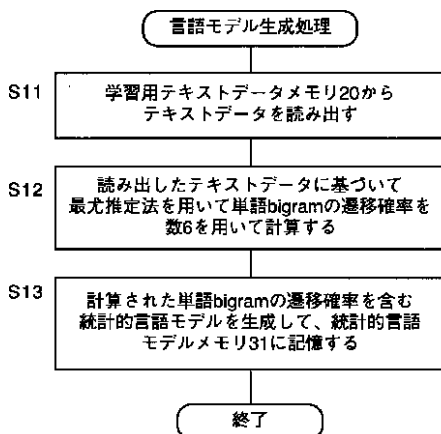
【図1】



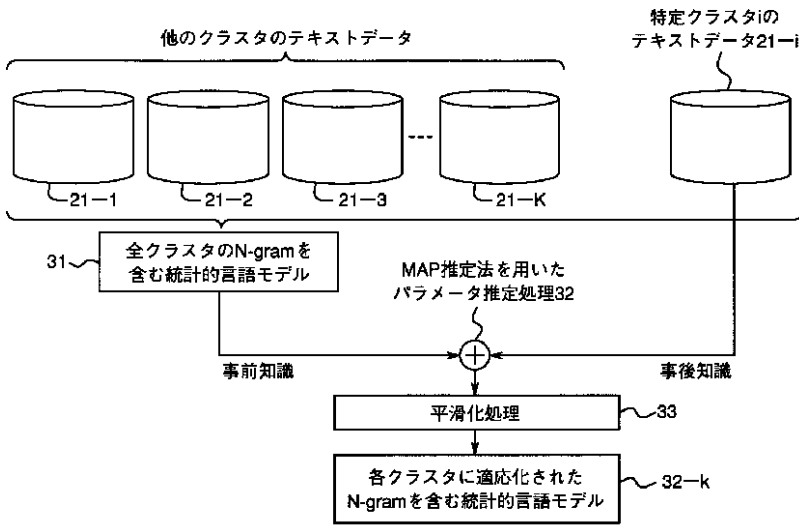
【図2】



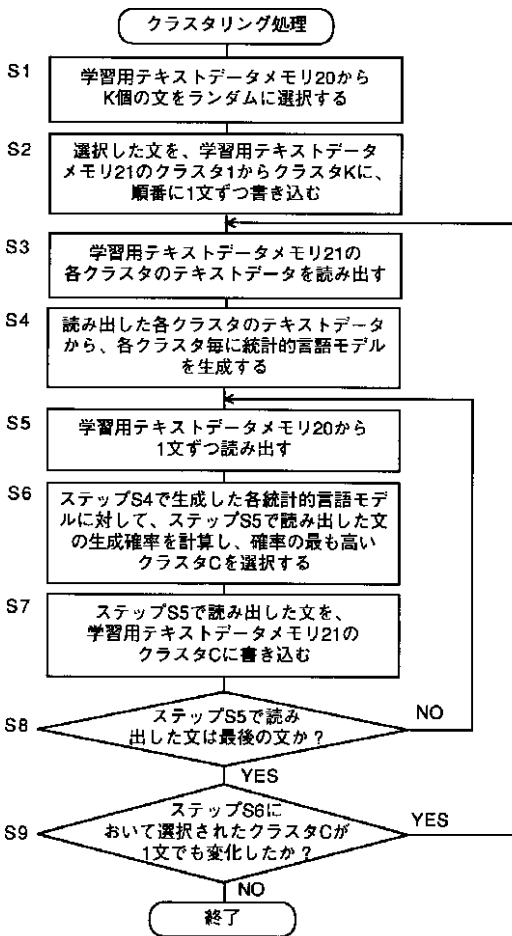
【図5】



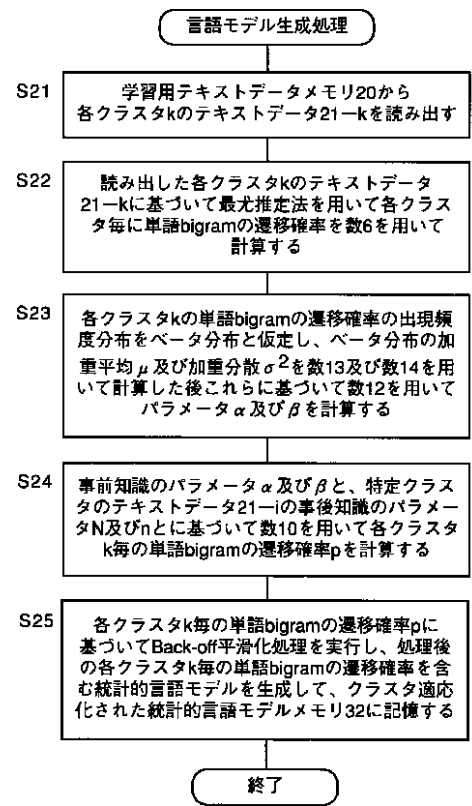
【図 3】



【図 4】



【図 6】



フロントページの続き

(56)参考文献 特開 平 6 - 27985 ( J P , A )

(58)調査した分野(Int.Cl.<sup>6</sup>, DB名)

G10L 3/00 - 9/20

J I C S T ファイル ( J O I S )