

卒業論文

連続音声認識における

韻律情報に関する研究

～ I MEM法を用いた句境界検出～

東北大学工学部 情報工学科 B4

三好 哲也

目次

1	序論	1
1.1	研究の背景	1
1.2	研究の目的	1
1.3	本論文の構成	1
2	ラグウインドウ法を用いた特徴パタンの抽出法	2
2.1	原理	2
2.2	ラグウインドウ法の問題点	3
3	2段AR過程の非定常スペクトル推定法を用いた特徴パタンの抽出法	4
3.1	音声生成のモデル	4
3.2	原理	5
3.3	瞬時化最大エントロピー法の原理	6
4	実験方法	7
4.1	学習	7
4.2	認識	7
4.3	システム構成	9
4.4	使用データ	9
4.5	結果の評価	10
5	実験結果	11
6	結論	16
6.1	まとめ	16
6.2	今後の課題	16
7	謝辞	17
8	参考文献	17

1 序論

1.1 研究の背景

近年、電気機器や、OA機器などの機能の充実に伴ない操作法が複雑化してきており、それらの機械とコミュニケーションをするのが困難になってきている。

このため、人と機械が、自然にコミュニケーションをとる手段としては、音声によるものが最も望ましいと思われる。しかし、会話文の認識は、明確に発声された単語を認識する場合と比べて非常に困難であることは、容易に推測できる。このため、単語認識に用いられる局所的な音韻情報のほかに音声のもつ、イントネーション、アクセントといった韻律情報について検討する必要がある。

1.2 研究の目的

連続音声を認識する上で 韻律情報である、アクセントパターンなどにより句境界ごとに、分割することにより単語認識並みの認識率を得られることを考えると韻律情報を用いた句境界検出は、重要であるといえる。

また、韻律情報の最も基礎となるものは、ピッチ(基本周波数)であり、このピッチは、音声の生成過程では、声帯の振動数に相当する。

そこで、本研究では、励振波スペクトルを句境界情報を含む特徴パターンとして、瞬時化最大エントピー法(IMEM法)を用いて抽出し、これをパターンとして、One stage DP法を用いて、句境界を検出することをその目的とする。

1.3 本論文の構成

1章では序論として、研究の背景及び目的、2章では、従来法である、ラグウインドウ法による、特徴パタンの抽出法、3章では、2段AR過程の非定常過程のスペクトル推定法による特徴パタンの抽出法、4章では、3章の方法によるパターンでの実験方法、5章では、実験の結果及び、考察、6章では、結論として、まとめ及び今後の課題を述べ、7、8章では謝辞及び参考文献を挙げる。

2 ラグウインドウ法を用いた特徴パタンの抽出法

2.1 原理

ラグウインドウ法は、フーリエ変換を用いて、パワースペクトルから、ピッチ構造を含んだスペクトルを分離する手法である。このラグウインドウ法によって、得られた、ピッチスペクトルを20次元の特徴パターンとして、One stage DP法により句境界検出を行なう。(以後、PSPパターンと呼ぶことにする。) ラグウインドウ法による分析系を図1に表わす。

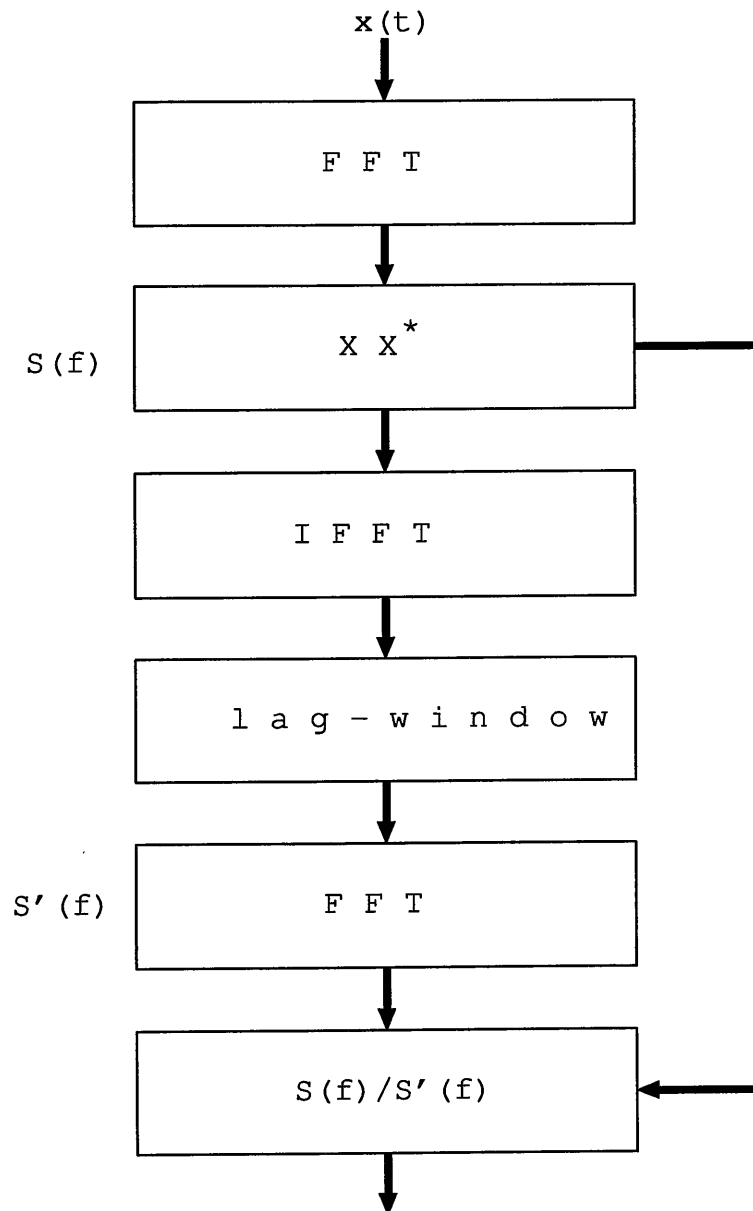


図1: ラグウインドウ法による分析系

2.2 ラグウインドウ法の問題点

ラグウインドウ法はFFT法に基づく分析法である。FFT法は、本来長いフレーム長が、必要とされるため、時間ごとに、変動する非定常過程のスペクトルを推定するには、フレーム間ごとに不規則な擾乱が生じ(時変性への追従能力が劣る。)、分析効率が、低い。このため、音声のような、非定常過程であるとみなせるような場合には、現象が時刻と共に顕著に変動し、時間依存性の情報が重要になる。そこで、スペクトルが、時刻に対して顕著に変動する非定常過程の正確なスペクトル推定法が必要となる。

3 2段AR過程の非定常スペクトル推定法を用いた特徴パタンの抽出法

3.1 音声生成のモデル

音声の生成の過程を2段階の過程で起こると仮定し、まず、1段階目で白色雑音から励振波が生成され、2段階目で音声波が生成される。この際の各段の過程をAR過程とする。この過程を図2に表わす。

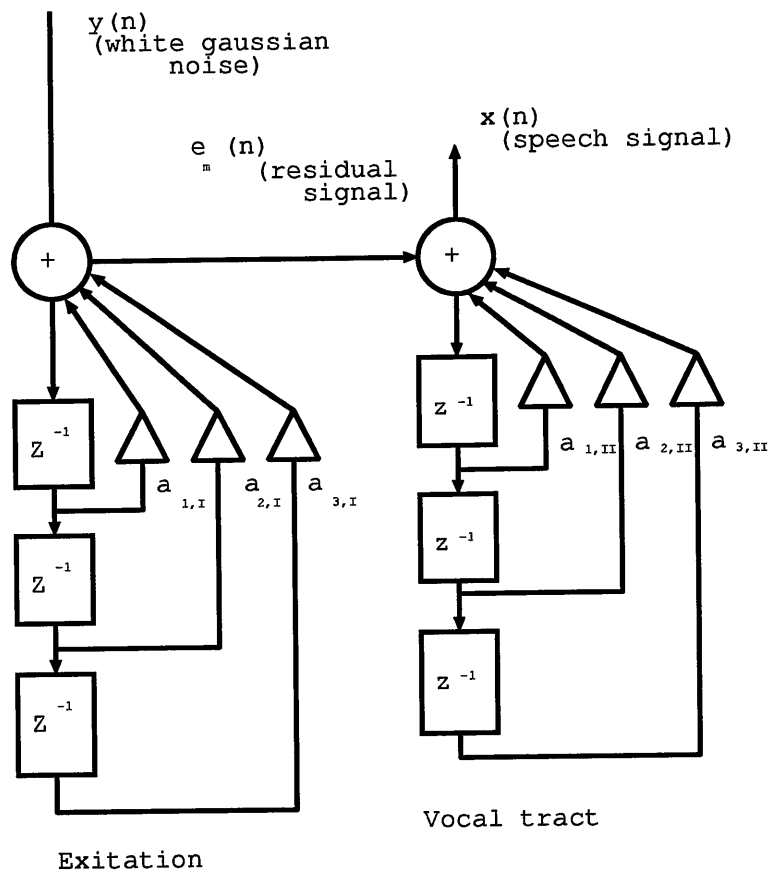


図 2: 音声生成モデル

3.2 原理

音声信号をまず、1段目の分析により、声道スペクトルを推定する。このとき、残差信号も同時に推定し得ることができる。この1段目で得られた、残差信号を低周波成分を効果的に推定するためにまず、ローパスフィルタを通過させ、次に間引きを行ない(本実験では、間引きを行なっていない。)、2段目の分析を行なうことにより、励振波スペクトルを推定することができる。この分析系の構成を図3に表わす。

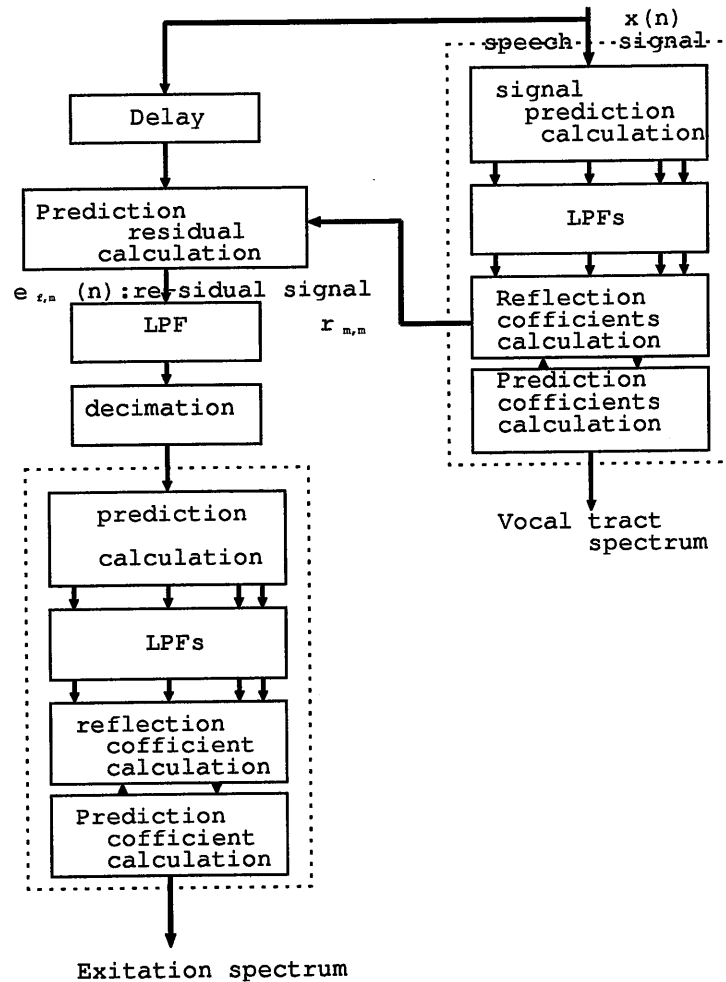


図 3: 励振波の分析系

この得られた励振波スペクトルをPSPパターンと同様に20次元のベクトルを抽出しこれを特徴パターンとする。この分析の際、非定常スペクトル推定法(瞬時化最大エントロピー法)を用いているため、フレームごとに不規則な変動が除かれた時間追従性のよい、スペクトルが得られる。そのために、これを特徴パターンとして使用することにより不規則な変動のない、句境界に関する情報が得られるのではないかと考えられる。

3.3 瞬時化最大エントロピー法の原理

瞬時化最大エントロピー法は、スペクトルが、時刻に対して顕著に変動する非定常過程のスペクトル推定の方法であり、時間依存性を損なわずに不規則な変動が除かれた、スペクトルが得られる特徴を持っている。

ここで、瞬時化最大エントロピー法の原理について述べる。

スペクトルのフーリエ変換が相関関数であるというWiener-Khitchineの関係の制約のもとで、エントロピーを増加させないように未知部分を推定することであり、AR過程($x_n = a_1x_{n-1} + a_2x_{n-2} + \dots + a_px_{n-p} + r_n$)で表わされる数値フィルターの係数を推定することにより、スペクトルを推定する方法であり、すなわち、あるスペクトルをもつランダム波を発生させる”白色雑音を入力する系”を探り出すということである。このとき、計算アルゴリズムは次のようになる。時刻nを終点とするフレームを考え、フレームの時刻をnで表わすとすると最小にすべき評価関数は次式である

$$E_m(n) = \frac{1}{2}(e_{f,m}^2(n) + e_{b,m}^2(n)) \quad (1)$$

$$e_{f,m}(n) = x(n) + \sum_{i=1}^m \gamma_{m,i}(n)x(n-i) \quad (2)$$

$$e_{b,m}(n) = x(n-m) + \sum_{i=1}^m \gamma_{m,i}(n)x(n-m+i) \quad (3)$$

ここで、 $e_{f,m}(n), e_{b,m}(n)$ は、前方および後方予測残差、 $\gamma_{m,i}(n)$ は、予測係数である。予測係数 $\gamma_{m,i}(n)$ について次式が成立する。

$$\gamma_{m,i}(n) = \gamma_{m-1,i}(n) + \gamma_{m,m}(n)\gamma_{m-1,m-i}(n) \quad (4)$$

但し $\gamma_{m,0} = 1$ である。式(1)を最小とする反射係数 $\gamma_{m,m}(n)$ は次式で表わせる。

$$\gamma_{m,m} = \frac{-2C_{m-1}(n)}{F_{m-1}(n) + B_{m-1}(n-1)} \quad (5)$$

$$C_m(n) = \sum_{i=0}^m \sum_{j=0}^m \gamma_{m,i}(n)\gamma_{m,j}(n)\phi(i, m+1-j; n)$$

$$F_m(n) = \sum_{i=0}^m \sum_{j=0}^m \gamma_{m,i}(n)\gamma_{m,j}(n)\phi(i, j; n)$$

$$B_m(n-1) = \sum_{i=0}^m \sum_{j=0}^m \gamma_{m,i}(n)\gamma_{m,j}(n)\phi(m+1-i, m+1-j; n) \quad (6)$$

$$\phi(i, j; n) = h(n) * (x(n-i)x(n-j)) \quad (7)$$

但し $h(n)$ は、ローパスフィルタのインパルス応答、 $*$ は、畳み込みである。

式(4)~(7)によって予測係数 $\gamma_{m,i}(n)$ を $m=1$ から $m=I$ (分析次数) まで再帰的に計算し、 $\gamma_{I,i}(n)$ を得る。

4 実験方法

4.1 学習

1. 入力された、音声信号を瞬時化最大エントロピー法で分析し、励振波スペクトルを推定する。これを、20次元のベクトル(46.9~492.9Hz、23.47Hzごと)に抽出し、サンプリング周期10msecごとに求めることによりベクトル列を得ることができる。
2. 前段階で得られた、励振波スペクトルのベクトル列を視察によって得られた句境界ごとに切り出しそれを500msecに線形伸縮する。
3. この句境界ごとに切り出された励振波スペクトルのベクトル列をLBG法を用いて、クラスタリングし複数の代表パターンを作成し、これをテンプレートとする。

4.2 認識

1. 入力された、連続音声信号を同様に瞬時化最大エントロピー法で分析し、励振波スペクトルのベクトル列を抽出する。
2. この前段階で得たスペクトルのベクトル列を学習により得たテンプレートとOne stage DP法により区間全体で最小2乗誤差基準による最適なテンプレート列を求める。
3. その得られたテンプレート列の境目の時刻を調べることにより句境界が、検出することができる。

この実験では、テンプレート数は、1, 2, 4, 8, 16, 32、また、各テンプレート数ごとにデルタピッチ寄与率を0.0, 0.5, 1.0で変化させて行ない計 $6 * 3 = 18$ 通りで認識を行なった。
また、今回は、間引きを行っていない。

ここで、デルタピッチとは、得られた振幅に対する重み付けされた最小2乗誤差直線の傾きであり次式で表わされる。

$$\min_{\Delta p_i, b} \sum_{k=-m/2}^{m/2} r_{i+k} w_k (p_{i+k} - (\Delta p_i k + b))^2$$

距離計算の際、パターンベクトル \hat{P}_j, \hat{P}_k の距離は、次式で表わされ

$$D(\hat{P}_j, \hat{P}_k) = \sum_{i=1}^L (1 - \alpha) |p_{ji} - p_{ki}|^2 + \alpha |\Delta p_{ji} - \Delta p_{ki}|^2$$

このとき、 α をデルタピッチ寄与率と呼ぶ。また、デルタピッチ寄与率は、パターン間の距離におけるデルタピッチ距離のしめる割合で0~1の値をとる。

また、One stage DP法のアルゴリズムについて次に述べる。

未知入力パタンのフレーム： $i = 1, \dots, N$

テンプレート番号： $k = 1, \dots, N$

テンプレート k のフレーム： $j = 1, \dots, j(k)$

(i, j, k) におけるフレーム間距離： $d(i, j, k)$

(i, j, l) における累積距離： $D(i, j, k)$

Step 1 initialize $D(1, j, k) = \sum_{n=1}^j d(1, n, k)$

Step 2

a) for $i := 2$ to N do step(b) - (e)

b) $D(i, 1, k) := d(i, 1, k) + \min(D(i-1, 1, k), D(i-1, J(k^*), k^*))$
 $(k^* = 1, \dots, K)$

c) for $j := 2$ to $J(k)$ do step(e)

e) $D(i, j, k) := d(i, j, k) + \min(D(i-1, j, k), D(i-1, j-1, k))$

Step 3 Trace back the best path

4.3 システム構成

本実験で用いたシステム構成を図4に表わす。ここで、図の点線枠の部分は、学習時に用いられ、実線枠部分は、認識時に用いられる。

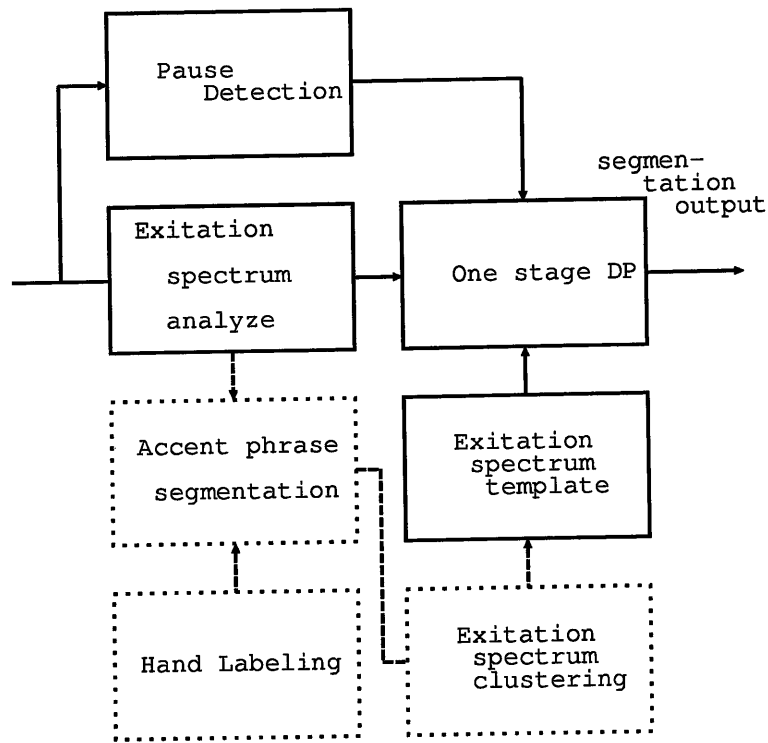


図 4: システム構成

4.4 使用データ

本実験で使用した、データベースを下の表に示す。

名称、分類	ATR 日本語音声データベース
テキスト	100 文章 内分け 2 グループ (A, B) 各 50 文章
話者	男性 1 名 (MY I)
A/D 変換	12 kHz サンプリング, 16 bit

学習用に A、認識実験用に B を用いた。

4.5 結果の評価

句境界の検出の正解の基準については、視察によって得られた、句境界の+-100 msecにあるとき、正解とした。また、評価基準として次式を用いた。

$$\text{句境界検出率} = \frac{\text{正解の句境界検出数の総数}}{\text{視察による句境界の総数}}$$

$$\text{誤認識率} = \frac{\text{誤検出した数の総数}}{\text{自動検出による句境界の総数}}$$

5 実験結果

オープン実験においては、テンプレート数8、デルタピッチ1.0のとき、認識率最高で、約76%であった。

クローズ実験においては、テンプレート数32、デルタピッチ0.0のとき、認識率最高で、約80%であった。

オープン実験とクローズ実験の認識結果を比較すると、認識率が、極端なくい違いが認められないことから、ある程度、個人についての句境界の情報を含んでいると考えられる。

しかし、従来法のラグウインドウ法による、特徴パターンと比較すると、非常に認識率が悪かった。また、誤認識をしているところは、無声音の部分や有声音のパワーの弱いところであり、パターンには、きれいにスペクトルのピークが出ていなかった。これは、本来、励振波スペクトルを推定の際には、データを間引きするのだが、本実験では、間引きを行っていないのが、原因と考えられる。間引きを行った場合と行なわない場合では、行った場合の方が、励振波スペクトルを推定する際行わない場合よりもピークが鋭敏に観察することができており、改善の余地がある。

もう1つ原因として考えられるのは、今回の実験では、パタンの抽出に、46.9 Hz ~ 492.9 Hzの20次元のベクトルであったが、本実験での分析法は、スペクトルのピークが、かなり狭い範囲で観測されるため、サンプリングの間隔の大きさによっては、ピークが抽出できずに直線的なパターンが抽出される可能性がある。しかし、サンプリングの間隔を狭めれば、当然、次元数が増加することになる。

認識率とデルタピッチ寄与率の関係のグラフを図5、図6、図7に表わす。また認識例を図8に示す。

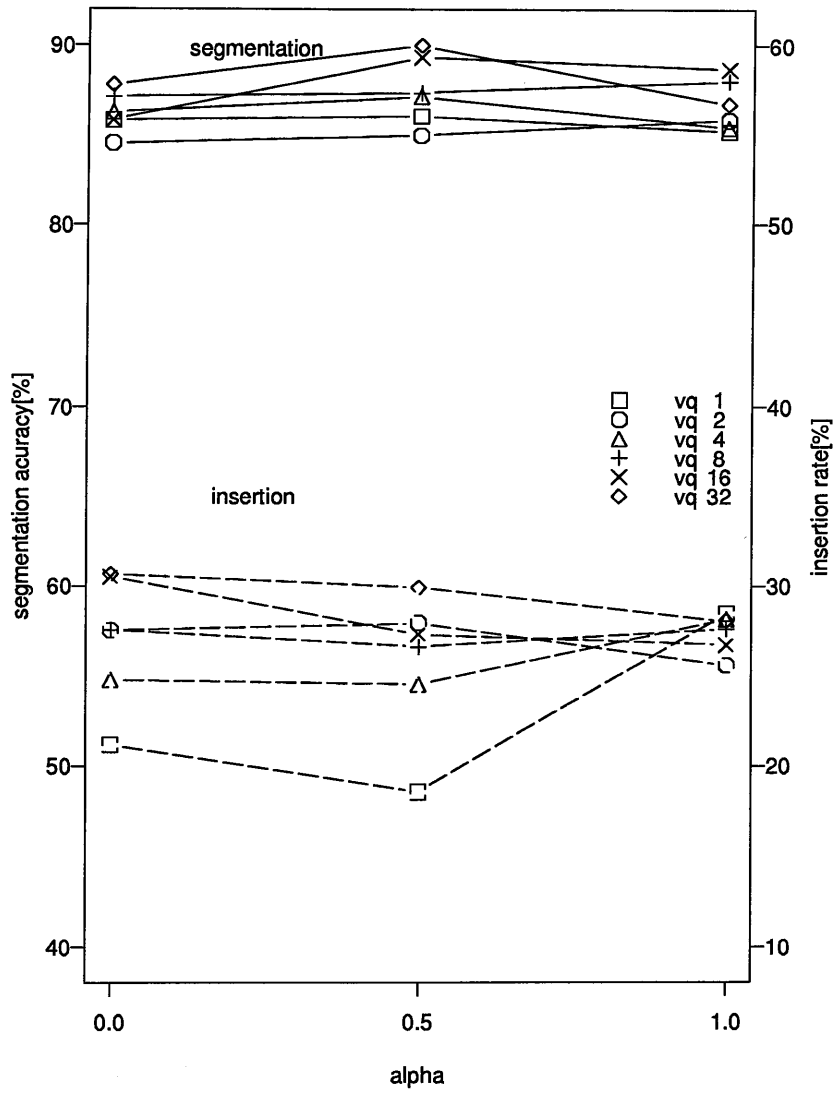


図 5: P S Pパタンの認識率

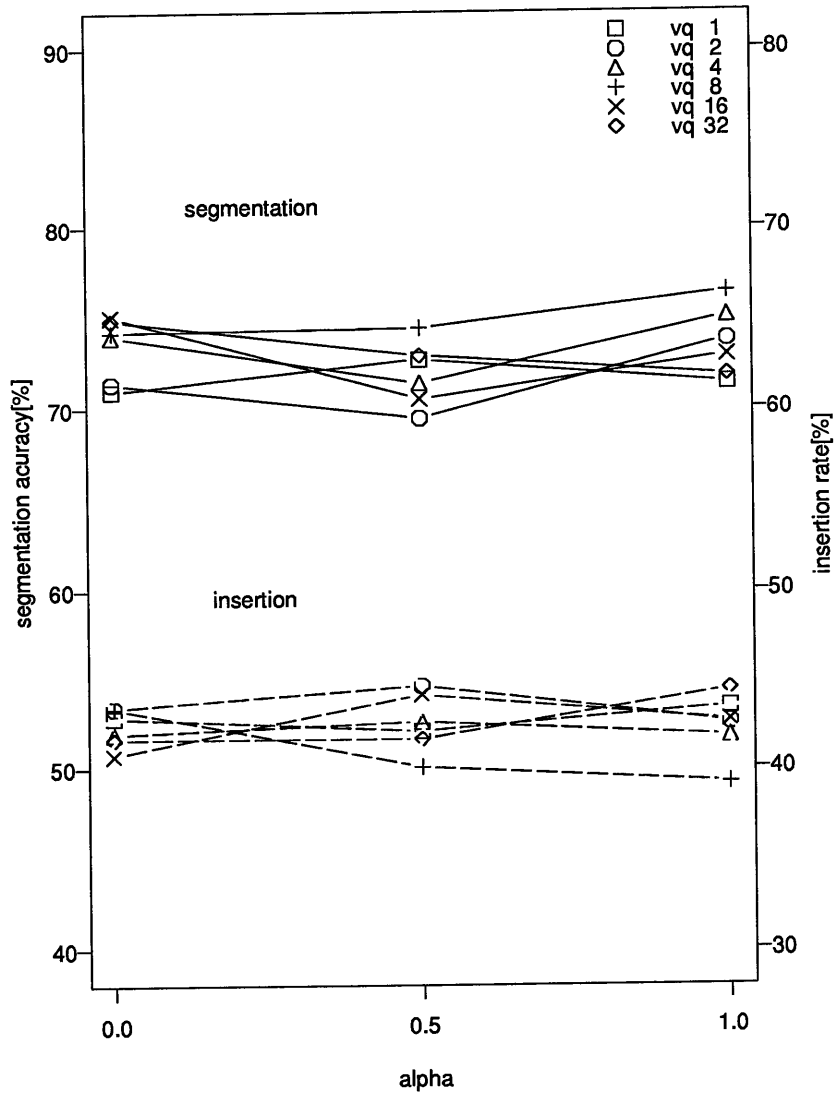


図 6: 励振波スペクトル (オープン実験)

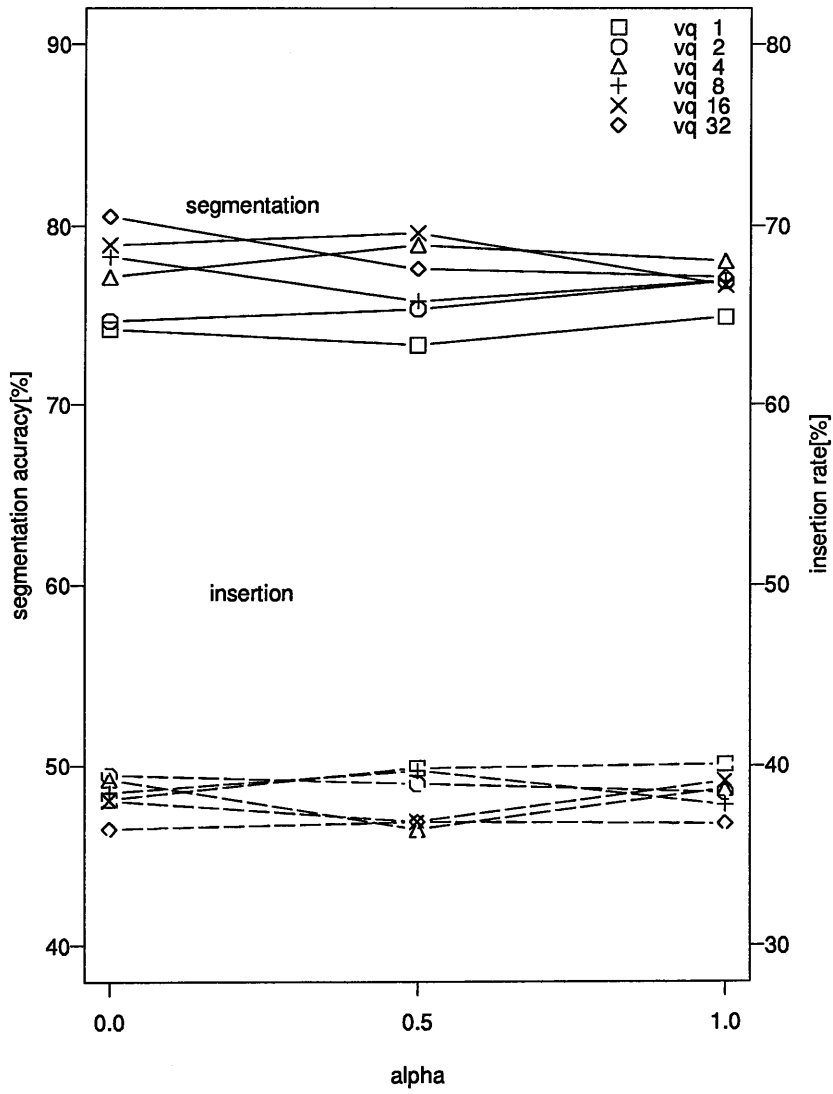


図 7: 励振波スペクトル(クローズ実験)

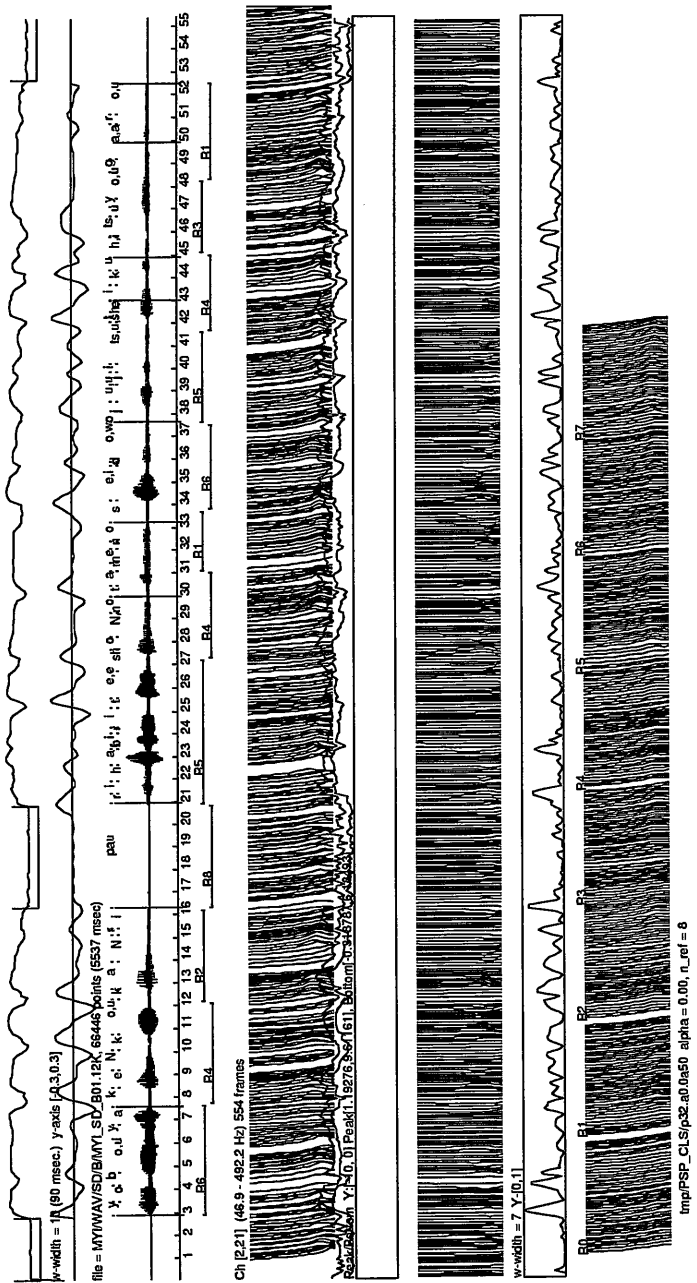


図 8: 認識例

6 結論

6.1 まとめ

従来のラグウインドウ法と異なる非定常過程の分析法(瞬時化最大エントピー法)で分析し特徴パタンの抽出を試みた。予備的な観察では、フレームごとで、不規則な擾乱の取り除かれた励振波スペクトルが、得られた。

認識実験では、オープン実験とクローズ実験から、ある程度の句境界情報を含む特徴パタンの抽出ができた。

本実験では、従来からある、ラグウインドウ法によって抽出された特徴パターンと比較すると、高精度なパタンの抽出はできなかった。

クラスタリングされた、代表パターンを調べてみると、視覚的に明らかにアクセントパターン(ピーク周波数の時間的変化)が表われているわけではなく、かなりスペクトルのピークが、平滑化されてしまっているため分かりづらいパターンとなっている。これには、本分析法の特性上、スペクトルのピークの変化が、狭い範囲で生じるためであり、このためには、今回の実験で行なった、20次元では、抽出しきれないことが分かった。

6.2 今後の課題

間引きを行った励振波スペクトル推定による特徴パタンの抽出することにより、スペクトルのピークが、きれいに観察されるため性能の改善が見込まれる。しかし、これは、その後の実験(間引きを行なったもの)でも、ピークが、あまりでておらず平滑化されてしまっているのでどうも平滑化されてしまう原因は、実験結果で述べた、次元数が少ないことによるものであるようだ。このため、次元数をもっと大きくして実験する必要があるようだ。

また、本研究では、定常性に関するスペクトル推定による特徴パタンの作成であったが、この分析法(瞬時化最大エントロピー法)は、非定常性に関するスペクトル推定も可能であるので、定常性と非定常性の2つのスペクトルをそれぞれ特徴パターンとして抽出しこの2つのパターンを用いて句境界を検出する。

本実験で本方法の有効性を示すことが、出来ていないため、他の方法を考える必要性もあるであろう。

7 謝辞

本研究を行なうにあたり、御指導、御助言を賜りました東北大学工学部阿曾弘具教授に心より感謝いたします。また、中井満氏、栗津辰功氏には、研究全般にわたり、親身な御指導を賜りました。

最後に日頃より数多くの御指導、御討論、御協力を賜りました、丸岡 研(阿曾グループ)の皆様に感謝します。

8 参考文献

- (1) 滝沢 由美、小田 啓介、渡辺 彰彦、深沢 敦司：“多段AR過程の非定常スペクトル推定”，信学技報DSP90-1, pp. 1-8
- (2) 滝沢 由美：“瞬時化最大エントロピー法に基づく非定常過程のスペクトル推定法”，電気情報通信学会論文誌 A Vol. J73-A No. 6 pp. 1083-1093 1990. 6
- (3) Hiroshi Shimodaira and Masayuki kimura：“Accent Phrase Segmentation Using Pitch Pattern Clustering”，IEEE 1992 ICASSP
- (4) 下平 博、嵯峨山 茂樹、木村 正行：“ピッチパタン連続整合による連続音声のセグメンテーション”，信学技報SP90-72, pp. 33-40