

自譜めぐりシステムのための多重音の 音高・音源数の推定法的高速化

安部翔[†] 小田弘良[†] 松島俊明[†]

我々は自動譜めぐりシステムの開発を進めており、そのために多重音の音高・音源数の高速推定の研究を行っている。

以前、複素スペクトル内挿法により基本周波数候補を求め、それらの組み合わせによる同時発音数を評価する尺度を導入し、この評価尺度が最小なる音高の組み合わせを直接求める方法で、音源数が未知の音響信号からの高速推定法について報告したが、検出可能な同時発音数を多くするに従って計算量も増え、処理落ちが時々生じるという問題点があった。

そこで、今回、この評価尺度の最小値の変化率の推移に着目し、これを用いて多重音の音高・音源数の高速推定を試みた結果、以前よりも高速な推定が可能となったので報告する。

A fast multi-pitch estimation method for the automatic page turner

Syou Abe[†], Hiroyoshi Oda[†] and Toshiaki Matsushima[†]

We are developing an automatic page turner system, therefore are studying high-speed estimation of the number of notes and their pitches simultaneously from acoustic sound.

We already reported the following method: that is, to calculate the candidates of fundamental frequencies depending on the complex-spectrum interpolation, to introduce a evaluation function to search for the combination of the pitches which gives the minimum value of the evaluation function directly, and to presume the number of notes and their pitches from acoustic sound to a high speed.

However, by this method, computational complexity also increased as the number of detective notes increased, and there was a problem that processing omission sometimes occurred.

By considering transition of the rate change of the minimum value of the evaluation function, faster presumption of the number of notes and their pitches from acoustic sound can be attained than before.

1. はじめに

筆者らは以前より尺八譜の情報処理システムの研究を行っており、その機能の一部として自動譜めぐりの実装を試みてきた[1]。自動譜めぐり機能に要求される要素技術の1つに実時間での音高推定がある。尺八演奏だけに限れば単音を対象とすれば十分であるが、箏や三味線などとの合奏曲への応用を考慮すると、複数音への対応も必要となる。

近年、複数の楽音を含んだ音響信号から音高情報を抽出する研究が盛んに行われるようになり、例えば混合正規分布モデルのパラメータ推定[2]や混合正規分布の調波時間構造のクラスタリング[3]などの方法が提案されており、目覚ましい成果を挙げている。しかし、これらの方法では基本的にEMアルゴリズム等などによる繰り返し計算を用いているため、処理時間のコストが非常に大きくなってしまったため、自動譜めぐりなど実時間処理が必要な用途に対しては、これらの方法を直接適用することは非常に困難である。

そこで、筆者らは、複素スペクトル内挿法[4]により求めた基本周波数候補と同時発音数の組み合わせを評価する関数を用い、この評価関数が最小となる推定音高の組み合わせを直接求めることで、音源数が未知な音響信号から音高・音数の高速推定を行う方法について既に報告した[5]。しかし、この方法では評価関数に現れるペナルティ係数の決定方法の問題や、推定可能最大音源数を増やした場合に処理落ちが生じるという問題があった。

今回、評価関数の値の減少傾向と元の音響信号に実際に含まれている音数の間に関連性があることに着目し、評価関数の値の推移を観測することでより高速な音高・音源数の同時推定を試みたので報告する。本報告では、既に報告済みの高速推定方法についてその概略の説明を行った後、今回、新たに考案した推定方法について述べ、それらの方法による実験結果を示す。

2. 音高推定のための評価関数

多重音の音高推定を行うための定式化方法はいくつかあるが、本研究では嵯峨山らにより考案されたハーモニッククラスタリング[6, 7]における定式化と同様の方法により定式化を行う。音源に含まれる単音の基本周波数を Ω とすると、一般的に窓関数の影響等により観測されるスペクトルの形状は Ω を中心として左右に単調減少しかつ対象な形状となる。一般の楽音の場合は倍音構造を持つため、倍音の周波数が基本周

[†] 東邦大学理学部情報科学科
Dept. of Information Sciences, TOHO University

波数 Ω の整数倍となることを仮定すると、 ω を周波数 $f(\omega)$ を観測スペクトルとして、

$$D(\mu) = \sum_n \int_{T_n} \phi(\omega, n\mu) \cdot f(\omega) d\omega \quad (1)$$

$$\text{ただし } T_n = \left\{ \omega \mid n = \arg \min_m |\omega - m\mu| \right\} \quad (2)$$

を最小とする μ として Ω は与えられる。ここで、 $\phi(\omega, n\mu)$ は ω と $n\mu$ の距離尺度を表す関数で、 $\omega = n\mu$ で極小値を取る左右対称な関数であれば何でもよいが、文献[6]と同様に

$$\phi(\omega, n\mu) = (\omega - n\mu)^2 \quad (3)$$

に選ぶと、式(1)は

$$D(\mu) = \sum_n \int_{T_n} (\omega - n\mu)^2 \cdot f(\omega) d\omega \quad (4)$$

となる。また T_n は $(n-1)\mu$ と $n\mu$ の中点および $n\mu$ と $(n+1)\mu$ の中点を両端とする帯域を表している。

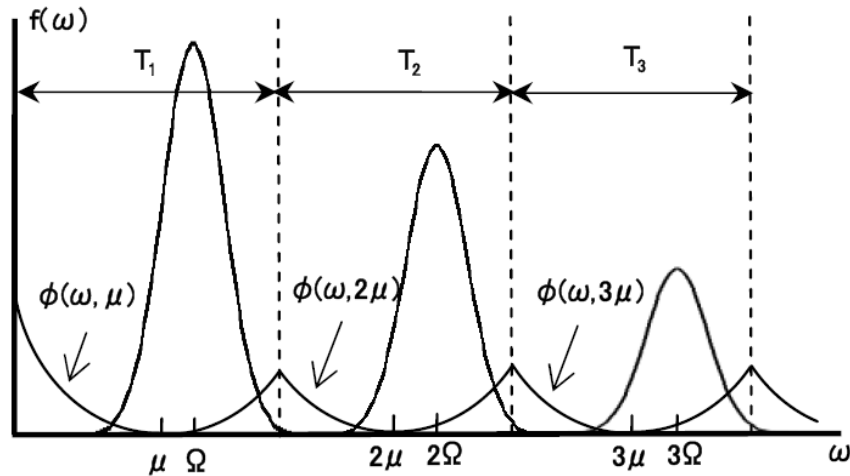


図 1 単音の場合の基本周波数の推定方法

音源数が複数の場合は、式(1)、(2)を複数音に拡張した評価関数、即ち、音源数を K 、基本周波数の推定値を $\mu = \{\mu_1, \mu_2, \dots, \mu_K\}$ とすると、評価関数

$$D(\mu) = \sum_{k=1}^K \sum_n \int_{T_n^k} \phi(\omega, n\mu_k) \cdot f(\omega) d\omega \quad (5)$$

$$\text{ただし } T_n^k = \left\{ \omega \mid (n, k) = \arg \min_{m, l} |\omega - m\mu_l| \right\} \quad (6)$$

を最小とする μ を求めれば良い。この最小値は、式(6)により与えられる周波数帯域 $T^k = T_1^k \cup T_2^k \cup \dots$ をクラスタと見なしてクラスタ重心である μ_k を

$$\mu_k = \frac{\sum_n \sum_{\omega \in T_n^k} \omega \cdot f(\omega)}{\sum_n \sum_{\omega \in T_n^k} f(\omega)} \quad (7)$$

により更新するという反復計算を行い、その収束値を求めることで得ることができるが、収束するまでの反復処理に時間がかかる、収束値が初期値に依存し、誤った値に収束することが多々あるため様々な初期値に対して試行する必要がある、等の問題がある。そのため、クラスタリングによる解法を譜めくり等リアルタイム性が要求される用途に適用するのは現状では難しい。

そこで、本研究では、周波数ピークの検出精度に優れた複素スペクトル内挿法により検出されたピーク位置の中に正しい基本周波数が含まれている確率が非常に高い点に着目し、クラスタリングにより $D(\mu)$ の収束値を求めるのではなく、様々な組み合わせの μ に対して直接 $D(\mu)$ の値を計算して μ の推定を行うことにした。

3. 音高・音源数の同時推定アルゴリズム(音源数既知)

まず、音源数(和音数) K が既知の場合の音高推定アルゴリズムを示す。短時間フーリエ変換(STFT)により得られたパワースペクトルのうち、一定以上のパワーを持つ上位 L 個(ただし $L > K$ とする)のピーク周波数 $\{\mu_1, \mu_2, \dots, \mu_L\}$ を複素スペクトル内挿法により抽出する。この L 個の基本周波数候補から選び出した K 個の基本周波数を小さい順に並べたベクトル

$$\mu = (\mu_{j(1)}, \mu_{j(2)}, \dots, \mu_{j(K)}), j(p) = 1, 2, \dots, L \quad (8)$$

に対して $D(\mu)$ を算出し、最小の $D(\mu)$ を与える μ を基本周波数の推定値とする。従って

音源数 K が既知の場合のアルゴリズムは以下のようになる。

(例) $K=2, L=4$ の場合

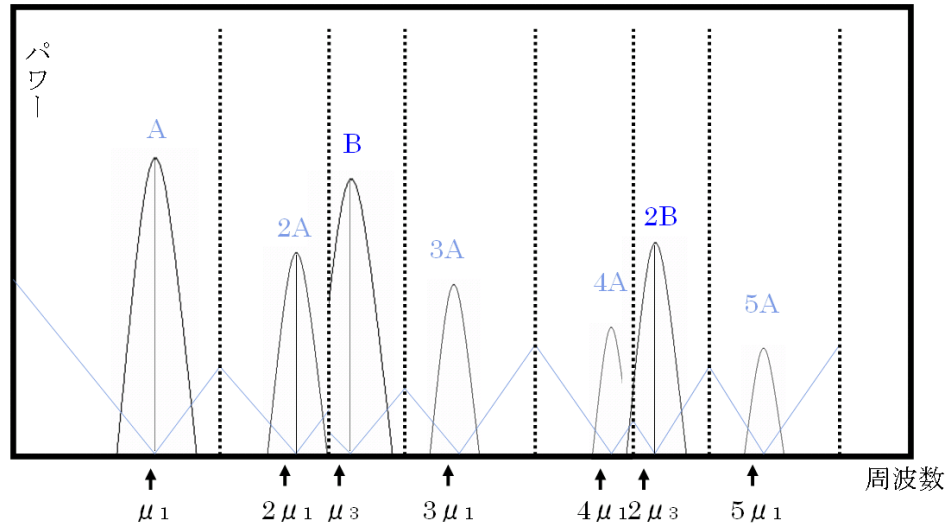


図 2. 多重音に対するハーモニッククラスタリングの使用例

アルゴリズム1：音源数 K が既知の場合

- ① 複素スペクトル内挿法によりパワーのピーク位置を検出し、パワーの大きい順に L 個の基本周波数候補 $\{\mu_1, \mu_2, \dots, \mu_L\}$ を選出する。
- ② L 個の候補から K 個を選択して得られる ${}_L C_K$ 個の基本周波数の組み合わせの集合を $M = \{\mu_1, \mu_2, \dots, \mu_{{}_L C_K}\}$ とする。
- ③ M の各要素について $D(\mu_i)$ を算出し、最小値を取る μ_i を基本周波数の推定値とする。即ち

$$\mu = \underset{\mu_i}{\operatorname{argmin}} D(\mu_i) \quad (\mu_i \in M) \quad (9)$$

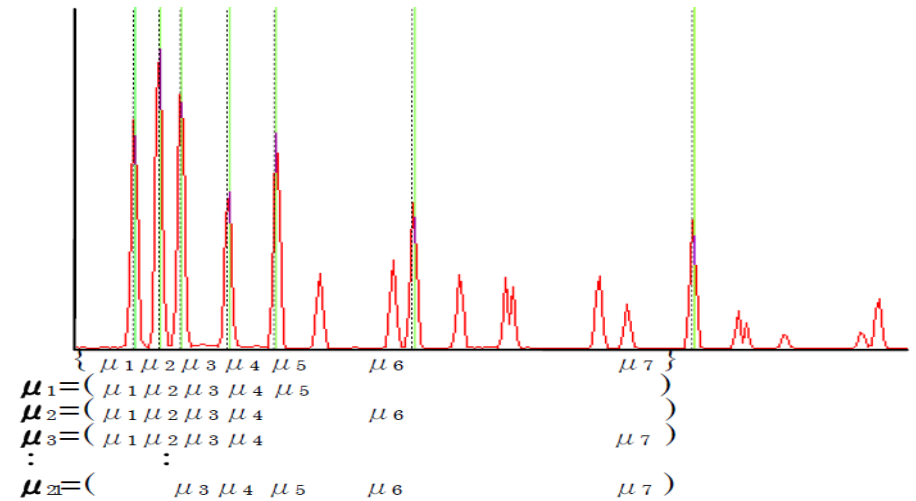


図 3. 複素スペクトル内挿法によるピーク検出と基本周波数候補の選択例

4. 音高・音源数の同時推定アルゴリズム(音源数未知)

次に、音源数が未知の場合であるが、音源数が K の時の μ を μ^K と表すと、音源数が K の場合の式(5)の最小値 $\min D(\mu^K)$ を $K=1$ から求めていき、最も小さい値を与える K および μ に決定する方法が考えられる。しかし、評価関数の値 $D(\mu^K)$ は、一般的に μ の次元 K に対して単調に減少していくので、実際の音源数よりも大きい次元の μ で最小となることが分かる。そこで、 μ の次元数の増加に対するペナルティ係数 $p(K)$ を追加し、 $p(K)D(\mu^K)$ を新たな評価関数とする。従って最大検出音源数(検出可能な最大の同時発音数)を N とすると、音源数が未知の場合の音高推定アルゴリズムは以下のようになる。

アルゴリズム2：音源数 K が未知の場合

- ① 複素スペクトル内挿法によりパワーのピーク位置を検出する。
- ② $K=1$ とする。
- ③ パワーの大きい順に L 個 ($L > K$) の基本周波数候補 $\{\mu_1, \mu_2, \dots, \mu_L\}$ を選出す

る。

④L個の候補からK個を選択して得られる ${}_L C_K$ 個の基本周波数の組み合わせの集合を $M^K = \{\mu_i^K | i=1, 2, \dots, {}_L C_K\}$ とする。

⑤ M^K の各要素について $D(\mu_i^K)$ を算出し、最小値 μ_i^K をとる音源数がKと仮定したときの基本周波数の推定値とする。即ち

$$\mu^K = \arg \min_{\mu_i^K} D(\mu_i^K) \quad (\mu_i^K \in M^K) \quad (10)$$

⑥ $K < N$ ならば $K \leftarrow K + 1$ として③へ戻る

⑦ $p(K)D(\mu^K)$ を最小とする μ^K を音源数Kと基本周波数 μ の推定値とする。即ち

$$(\mu, K) = \arg \min_{\mu^K} p(K)D(\mu^K) \quad (K = 1, 2, \dots, N) \quad (11)$$

本アルゴリズムにより、音高・音源数を同時に推定することが可能であるが、実際の音源数Kの数に関係なく、検出可能最大同時発音数Nまで $D(\mu_i^K)$ の値を計算する必要があるため、Nの値が大きくなると処理速度が低下して処理落ちが生じ易くなる。また、最適なペナルティ係数 $p(K)$ を決定方法が難しい、という問題がある。これらの問題を解決する方法として、 μ^K の値の推移に着目した方法を次に述べる。

5. 音高・音源数の同時推定アルゴリズム(音源数未知・改良型)

式(5)からも分かるように、 $D(\mu)$ の値は、 μ が基本周波数のみで構成されている場合には小さな値を取る。一方、基本周波数以外に加え、倍音の要素が加わっている場合には、クラスタ範囲が重複するだけなので、 $D(\mu)$ の値を減少させる効果は大きくないと予想できる。即ち、 μ の要素を増やしていく場合、 μ に基本周波数の要素が追加されている間は $D(\mu)$ は減少を続けるが、すべての基本周波数が μ に含まれた後、倍音の要素が追加されるようになると、 $D(\mu)$ の値の変化が小さくなることを期待できる。

この仮定が成り立っているかを実際に調べてみたところ、例えば3和音の場合は図4に示すように、推定音源数Kが実際の音源数に一致するまでは $D(\mu)$ の値は減少していくが、音源数よりもKが大きくなると $D(\mu)$ の値はほとんど変化しなくなる場合が多かった。

このことから、 $D(\mu)$ の値の変化の様子、即ち減少傾向を調べ、減少幅が小さくなっ

た時点で推定音源数を決定することができる。即ち、 $D(\mu)$ の変化率 $r(K)$ を

$$r(K) = D(\mu^{K+1}) / D(\mu^K) \quad (12)$$

と定義すると、変化率が1に近いある一定値以上の値となった場合、即ち

$$r(K) > \alpha: \text{ただし } \alpha \text{ は } 1 \text{ に近い定数} \quad (13)$$

を始めて満たした時のKは、実際の音源数より1つ大きい値を取るようになる。このことから、以下のアルゴリズムにより、音高・音源数の同時推定が可能となる。

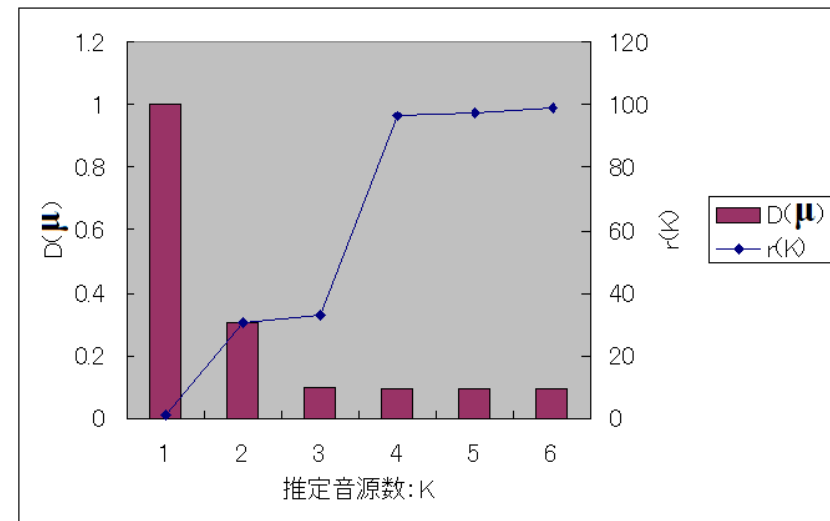


図 4. 推定音源数を増やしていった場合の $D(\mu)$ と $r(K)$ の算出結果の推移

アルゴリズム3：音源数Kが未知の場合(改良型)

- ①複素スペクトル内挿法によりパワーのピーク位置を検出する。
- ②推定音源数 $K = 1$ とする。
- ③パワーの大きい順にL個 ($L > K$)の基本周波数候補 $\{\mu_1, \mu_2, \dots, \mu_L\}$ 選出する。
- ④L個の候補からK個を選択して得られる ${}_L C_K$ 個の基本周波数の組み合わせの集合を $M^K = \{\mu_i^K | i=1, 2, \dots, {}_L C_K\}$ とする。
- ⑤ M^K の各要素について $D(\mu_i^K)$ を算出し、最小値 μ_i^K をとる音源数がKと仮定し

たときの基本周波数の推定値 μ^K とする. 即ち

$$\mu^K = \arg \min_{\mu_i^K} D(\mu_i^K) \quad (\mu_i^K \in M^K) \quad (14)$$

- ⑥ $K = 1$ ならば $K \leftarrow K + 1$ として③に戻る.
- ⑦ $2 \leq K < N$ のとき, 式(13)が成り立っていれば推定音源数を $K - 1$ に決定し, μ^{K-1} を推定音高とする. そうでなければ $K \leftarrow K + 1$ として③へもどる.
- ⑧ $K = N$ ならば推定音源数を K とし μ^K を推定音高とする.

表1にアルゴリズム2とアルゴリズム3の $D(\mu)$ の計算回数の比較を示した. この表からも分かるように, アルゴリズム3では推定音源数を実際の音源数

実際の音源数	アルゴリズム2	アルゴリズム3
単音	83	9
2和音	83	19
3和音	83	34
4和音	83	55
5和音	83	83
6和音	83	83

表1. アルゴリズム2とアルゴリズム3の $D(\mu)$ の計算回数の比較 ($L=K+2$ とした場合)

+1まで計算すれば良いため, 特に実際の音源数が少ない場合にアルゴリズム2に比べて計算回数を大幅に減らすことができることがわかる. またアルゴリズム2では必要となるペナルティ係数の設定が不要なため, アルゴリズム2の問題点を同時に解決することができる.

6. 実験結果

以下に, 提案アルゴリズムによる音高・音源数の推定実験の結果を示す. いずれもSTFTの窓長は4096, 窓関数はハミング窓を用い, 基本周波数候補数は $N=6, L=K+2$ とした. 予備実験の結果から, アルゴリズム2におけるペナルティ係数の値は $p(K)=K$ を, またアルゴリズム3における変化率の閾値 $\alpha=0.8$ を用いている.

まず, アルゴリズム2のシステムとアルゴリズム3のシステムの性能を比べるためにMIDI音源により作成した音源数が変化する音響データについて音源数・音高推定を行った結果を図5に示す(実験1).

アルゴリズム2に比べてアルゴリズム3では, CDEを含む和音のように, 比較的

基本周波数が近い音を含む和音で検出誤りが増えていることが分かる. しかし処理速度に関しては, 単音から4和音に対してはアルゴリズム2よりもアルゴリズム3の方が処理速度は向上し, 処理落ちの回数が明らかに少なくなった.

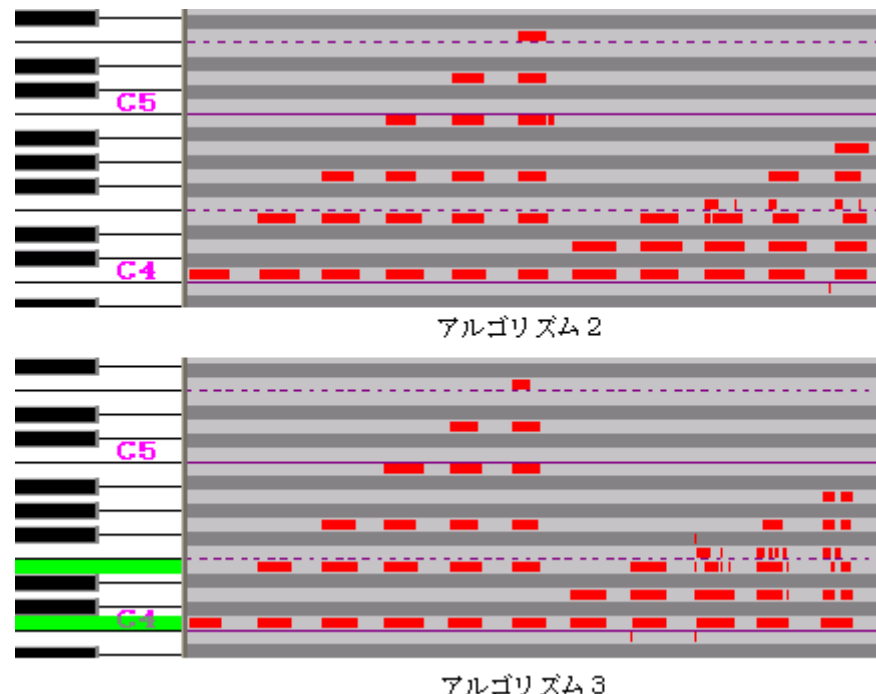


図5. 音源数が未知の場合の採譜結果
(アルゴリズム2と3の比較)

次に, RWC 研究用音楽データベース[8]に含まれている楽曲サンプル(MDB-J-2001 No.9)に対しての実験結果を図6に示す(実験2). 実験条件は実験1と同一である. どちらもやや検出もれが目立つが, 和音を含む多くの音符が検出できており, 音源数・音高の同時推定にある程度成功していることが分かる. 処理落ちに関しては, アルゴリズム2は3和音程度から顕著になってきたのに比べ, アルゴリズム3では5和音以上で処理落ちが目立つようになった. アルゴリズム2に比べて処理速度の大幅な向上を実現できたと言えるが, まだまだ改善の余地が残されている.

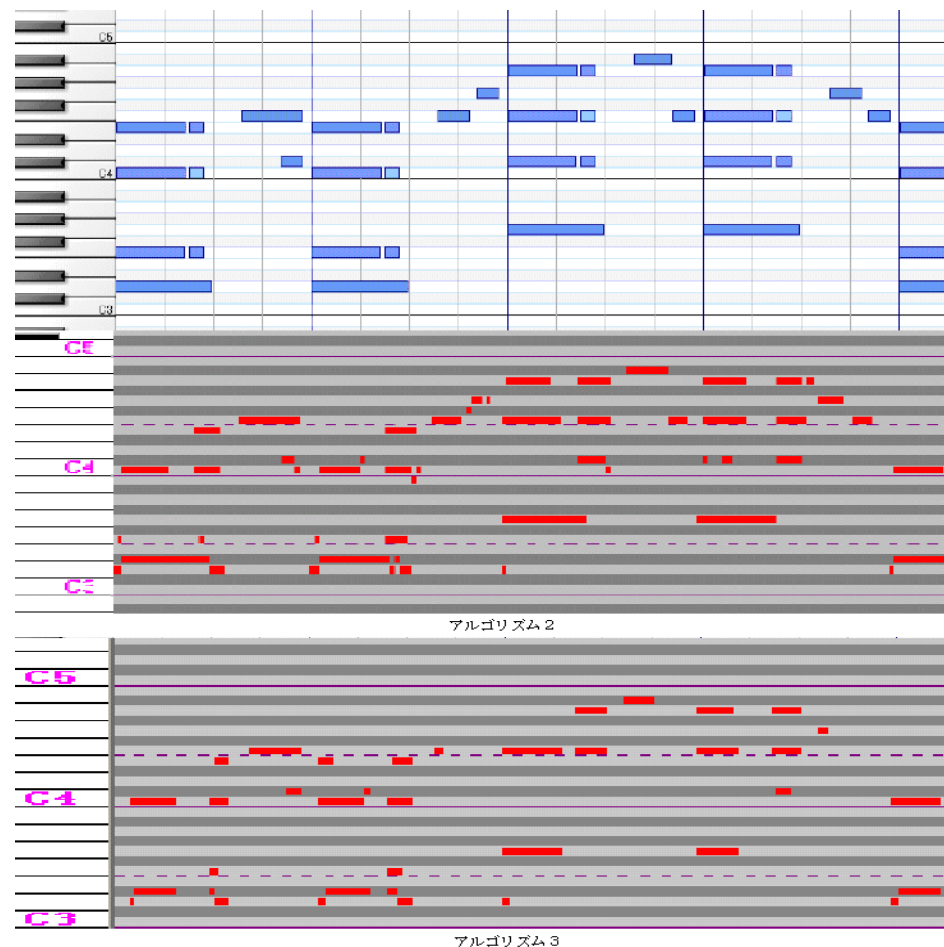


図 6. 実験 3 : 楽曲サンプルデータに対する音高推定結果
 (上段 : 正解 MIDI データ, 中段 : アルゴリズム 2 下段 : アルゴリズム 3)

7. まとめ

一般的に音高と音源数の同時推定は難しい問題であるが, 本方式では比較的簡単な処理により, 高速に音高と音源数の同時推定を行うことができた. 従来のアルゴリズム 2 では理想的なペナルティ係数の設定が難しかったが, 今回改良したアルゴリズム 3 ではそのペナルティ係数を必要とせずに音高と音源数の同時推定が可能になった. しかし, 実音を対象とした実験では, 実際の音源数よりも検出できた同時発音数が少ないことが多いため, 今後さらなる改善が必要である.

また, 音源数が多い場合(概ね 5 和音以上)については, 従来手法と処理量が変わらないため, 処理落ちの問題は完全には解決できていない. また, 低音部やオクターブ重なりなどの和音の検出精度が良くない点など, 未解決の問題点が多く残されている.

しかし, 自動譜めくりシステムでは, 通常の採譜システムとは異なり, 正解楽譜の情報をシステムが持っている状態での処理であるため, 必ずしもすべての音符情報を正確に認識できる必要はない. 4 音程度までの和音の検出が可能であれば, 既知の楽譜中からの演奏箇所の推定は可能と思われるので, 自動譜めくりシステムへの利用に十分適用可能な方法であると考えている.

参考文献

- [1] 松島俊明, "尺八譜のマルチメディア情報処理", 画像電子学会第 34 回年次大会予稿集, pp.23-30(2006)
- [2] 後藤真孝, "音楽音響信号を対象としたメロディーとベースの音高推定", 電子情報通信学会論文誌, Vol.J84-D-II, No.1, pp.12-22(2001)
- [3] 亀岡弘和, 西本卓也, 嵯峨山茂樹, "調波時間構造クラスタリング(HTC)による音楽音響特徴量の同時推定", 情報処理学会研究報告, 2005-MUS-61-12, pp.71-78(2005)
- [4] 原裕一郎, 井口征士, "複素スペクトルを用いた周波数同定", 計測自動制御学会論文集, Vol.19, No.9, pp.718-723(1983)
- [5] 宮坂広純, 小田良弘, 松島俊明, "多重音の基本周波数評価尺度の最小値選択による実時間・音源数推定", 第 7 回情報科学技術フォーラム(FIT2008), E-045, pp. 241-243(2008)
- [6] 亀岡弘和, 西本卓也, 嵯峨山茂樹, "ハーモニッククラスタリングによる多重音の基本周波数推定アルゴリズム", 情報処理学会研究報告書, 2003-MU-50-5, pp27-32(2003)
- [7] 亀岡弘和, 西本卓也, 嵯峨山茂樹 "ハーモニッククラスタリングと情報量基準による音楽の音高/音源数の推定", 情報処理学会研究報告, 2005-MUS-62-5, pp.27-32(2005)
- [8] 後藤真孝, 橋口博樹, 西村拓一, 岡隆一, "RWC 研究者用音楽データベース: クラシック音楽データベースとジャズ音楽データベース" 情報処理学会研究報告, 2002-MUS-44-5, pp.255-32(2002)