

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号
特許第6371516号
(P6371516)

(45) 発行日 平成30年8月8日(2018.8.8)

(24) 登録日 平成30年7月20日(2018.7.20)

(51) Int.Cl.

F I

G 1 O L 21/0308 (2013.01)

G 1 O L 25/18 (2013.01)

G 1 O L 21/0308 Z

G 1 O L 25/18

請求項の数 15 (全 17 頁)

(21) 出願番号	特願2013-237353 (P2013-237353)	(73) 特許権者	000001007
(22) 出願日	平成25年11月15日 (2013.11.15)		キヤノン株式会社
(65) 公開番号	特開2015-96921 (P2015-96921A)		東京都大田区下丸子3丁目30番2号
(43) 公開日	平成27年5月21日 (2015.5.21)	(74) 代理人	100076428
審査請求日	平成28年11月15日 (2016.11.15)		弁理士 大塚 康德
		(74) 代理人	100112508
			弁理士 高柳 司郎
		(74) 代理人	100115071
			弁理士 大塚 康弘
		(74) 代理人	100116894
			弁理士 木村 秀二
		(74) 代理人	100130409
			弁理士 下山 治
		(74) 代理人	100134175
			弁理士 永川 行光

最終頁に続く

(54) 【発明の名称】 音響信号処理装置および方法

(57) 【特許請求の範囲】

【請求項 1】

音響信号に非負値行列因子分解を適用することで、複数の基底スペクトルから構成される基底行列を取得する取得手段と、

前記取得手段により取得される前記基底行列を構成する基底スペクトルに対応する基底ケプストラムの、所定範囲のケフレンシーにおける最大値に基づいて、該基底スペクトルを第一の基底スペクトル群と第二の基底スペクトル群の何れかに分類する分類手段と、を備えることを特徴とする音響信号処理装置。

【請求項 2】

前記取得手段により取得された前記複数の基底スペクトルを前記分類手段により分類した結果に基づいて、前記音響信号に含まれる雑音成分が抑制された目的音信号を生成する生成手段を有することを特徴とする請求項 1 に記載の音響信号処理装置。

【請求項 3】

前記基底行列を構成する複数の基底スペクトルのうち、前記分類手段により前記第一の基底スペクトル群に分類される基底スペクトルの重心周波数は、前記分類手段により前記第二の基底スペクトル群に分類される基底スペクトルの重心周波数よりも高いことを特徴とする請求項 1 又は 2 に記載の音響信号処理装置。

【請求項 4】

前記分類手段は、ケフレンシーが所定値以下である低ケフレンシー部分における該基底ケプストラムの最大値に基づいて、該基底スペクトルを前記第一の基底スペクトル群と前

10

20

記第二の基底スペクトル群の何れかに分類することを特徴とする請求項 1 又は 2 に記載の音響信号処理装置。

【請求項 5】

前記分類手段は、前記低ケフレンシー部分における該基底ケプストラムの最大値と閾値との比較により、該基底スペクトルを前記第一の基底スペクトル群と前記第二の基底スペクトル群の何れかに分類することを特徴とする請求項 4 に記載の音響信号処理装置。

【請求項 6】

音響信号に非負値行列因子分解を適用することで、複数の基底スペクトルから構成される基底行列を取得する取得手段と、

前記取得手段により取得される前記基底行列を構成する基底スペクトルに対応する基底ケプストラムの、所定の基本周波数範囲に対応する部分の値に基づいて、該基底スペクトルを、第一の基底スペクトル群と、前記第一の基底スペクトル群に分類される基底スペクトルよりも調波成分が大きい基底スペクトルが分類される第二の基底スペクトル群との何れかに分類する分類手段と、を備えることを特徴とする音響信号処理装置。

10

【請求項 7】

前記分類手段は、該基底ケプストラムの前記所定の基本周波数範囲に対応する部分の最大値と閾値との比較により、該基底スペクトルを前記第一の基底スペクトル群と前記第二の基底スペクトル群の何れかに分類することを特徴とする請求項 6 に記載の音響信号処理装置。

【請求項 8】

20

前記第一の基底スペクトル群へ分類される基底スペクトルの数である第一の基底数と、前記第二の基底スペクトル群へ分類される基底スペクトルの数である第二の基底数の少なくとも一方を調整する調整手段をさらに備える請求項 1 乃至 7 の何れか 1 項に記載の音響信号処理装置。

【請求項 9】

前記分類手段により前記第一の基底スペクトル群に分類される基底スペクトルと、前記非負値行列因子分解の適用により取得されるアクティビティ行列を構成する複数のアクティビティベクトルのうちの前記第一の基底スペクトル群に対応するアクティビティベクトルとを用いて、第一の音響復元信号を生成する生成手段を更に備えることを特徴とする請求項 1 乃至 8 の何れか 1 項に記載の音響信号処理装置。

30

【請求項 10】

前記分類手段により前記第二の基底スペクトル群に分類される基底スペクトルと、前記非負値行列因子分解の適用により取得されるアクティビティ行列を構成する複数のアクティビティベクトルのうち前記第二の基底スペクトル群に対応するアクティビティベクトルとを用いて、第二の音響復元信号を生成する生成手段を更に備えることを特徴とする請求項 1 乃至 8 の何れか 1 項に記載の音響信号処理装置。

【請求項 11】

前記分類手段により前記第一の基底スペクトル群に分類される基底スペクトルと、前記非負値行列因子分解の適用により取得されるアクティビティ行列を構成する複数のアクティビティベクトルのうちの前記第一の基底スペクトル群に対応するアクティビティベクトルとを用いて、第一の音響復元信号を生成する第一生成手段と、

40

前記分類手段により前記第二の基底スペクトル群に分類される基底スペクトルと、前記アクティビティ行列を構成する複数のアクティビティベクトルのうち前記第二の基底スペクトル群に対応するアクティビティベクトルとを用いて、第二の音響復元信号を生成する第二生成手段をさらに備え、

前記第一生成手段により生成される前記第一の音響復元信号と前記第二生成手段により生成される前記第二の音響復元信号の少なくとも一方を用いて、前記音響信号から非目的音を除去することを特徴とする請求項 1 乃至 8 のいずれか 1 項に記載の音響信号処理装置。

【請求項 12】

50

前記取得手段は、前記音響信号に対して時間周波数変換を行うことで得られる行列を、前記非負値行列因子分解により基底行列とアクティビティ行列に分解することで、前記基底行列を取得することを特徴とする請求項 1 乃至 11 のいずれか 1 項に記載の音響信号処理装置。

【請求項 13】

音響信号処理装置の制御方法であって、

音響信号に非負値行列因子分解を適用することで、複数の基底スペクトルから構成される基底行列を取得する取得工程と、

前記取得工程において取得される前記基底行列を構成する基底スペクトルに対応する基底ケプストラムの、所定範囲のケフレンシーにおける最大値に基づいて、該基底スペクトルを第一の基底スペクトル群と第二の基底スペクトル群の何れかに分類する分類工程と、を有することを特徴とする音響信号処理装置の制御方法。

10

【請求項 14】

音響信号処理装置の制御方法であって、

音響信号に非負値行列因子分解を適用することで、複数の基底スペクトルから構成される基底行列を取得する取得工程と、

前記取得工程において取得される前記基底行列を構成する基底スペクトルに対応する基底ケプストラムの、所定の基本周波数範囲に対応する部分の値に基づいて、該基底スペクトルを、第一の基底スペクトル群と、前記第一の基底スペクトル群に分類される基底スペクトルよりも調波成分が大きい基底スペクトルが分類される第二の基底スペクトル群との何れかに分類する分類工程と、を有することを特徴とする音響信号処理装置の制御方法。

20

【請求項 15】

コンピュータを、請求項 1 乃至 12 のいずれか 1 項に記載の音響信号処理装置の各手段として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、音響信号をたとえば目的音と雑音といったように複数の音響信号に分離する音響信号処理装置および方法に関する。

【背景技術】

30

【0002】

音響信号から非目的音である雑音を除去する技術は、音響信号に含まれる目的音に対する聴感を改善し、また、音声認識においては認識率を高めるために重要な技術である。

【0003】

音響信号から雑音を除去する技術として、非負値行列因子分解を用いたものがある。これは、音響信号を時間周波数変換した行列を非負値行列因子分解によって基底行列とアクティビティ行列に分解するもので、これらの行列が目的音に係る部分行列と雑音に係る部分行列に分けられるという仮定に基づく。そして、目的音に係る部分基底行列である目的音基底行列と、目的音に係る部分アクティビティ行列である目的音アクティビティ行列を用いて、雑音が除去された目的音復元信号を生成する。

40

【0004】

特許文献 1 では、雑音除去対象である音響信号とは別に目的音と雑音をそれぞれ用意し、それらを事前学習することで目的音と雑音それぞれの教師基底行列および教師アクティビティ行列を得る。そして教師基底行列および、教師アクティビティ行列の統計量情報を用い、音響信号を時間周波数変換した行列を分解して目的音復元信号を得る。

【0005】

特許文献 2 では、2ch の音響信号それぞれを時間周波数変換した 2 つの行列を非負値行列因子分解する。そして、各 ch の基底行列の各列を構成する基底スペクトルについて、ch 間の相関が高いものを雑音基底スペクトル、それ以外のものを目的音基底スペクトルとする。そして、目的音基底スペクトルで構成される目的音基底行列と、それに対応す

50

る目的音アクティビティ行列を用いて目的音復元信号を生成する。

【先行技術文献】

【特許文献】

【0006】

【特許文献1】特開2009-128906号公報

【特許文献2】特開2012-022120号公報

【発明の概要】

【発明が解決しようとする課題】

【0007】

特許文献1の方法では、別に用意した音から基底行列を事前学習し、これを用いて復元信号を生成することになる。これは、調波楽器の音階のように、基底スペクトルの形状（調波構造）が大体決まっている場合の楽器音分離（例えば自動採譜に用いる）などには好適であると考えられる。しかし、それ以外の場合は、分離対象の音響信号に含まれる音とは異なる基底スペクトルを用いて復元信号を生成する可能性があるため、音質劣化につながり得る。

【0008】

特許文献2の方法は、雑音除去対象の音響信号から基底行列を得るため、目的音基底行列と雑音基底行列にうまく分けられれば、実際の目的音に対応する基底スペクトルを用いて目的音復元信号を生成できると期待される。しかしながら、目的音基底スペクトルと雑音基底スペクトルへの分類はc h間の相関に基づいて行うため、複数c hの音響信号を必要としている。

【0009】

また、相関は2つの基底スペクトルの組み合わせに対して算出される量であるが、基底スペクトル間のユークリッド距離や内積を用いるとしている。しかし、このような単純な相関指標は物理的な意味が不明瞭であり、基底スペクトルの分類において必ずしも好適ではない。

【0010】

本発明は上述した問題を解決するためになされたものであり、音響信号の各々の基底スペクトルを高精度に分類することが可能な音響信号処理装置および制御方法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】

【0011】

上記目的を達成するための本発明の一態様による音響信号処理装置は以下の構成を備える。すなわち、

音響信号に非負値行列因子分解を適用することで、複数の基底スペクトルから構成される基底行列を取得する取得手段と、

前記取得手段により取得される前記基底行列を構成する基底スペクトルに対応する基底ケプストラムの、所定範囲のケフレンシーにおける最大値に基づいて、該基底スペクトルを第一の基底スペクトル群と第二の基底スペクトル群の何れかに分類する分類手段と、を備える。

【発明の効果】

【0012】

本発明によれば、音響信号の各々の基底スペクトルから算出された指標を用いて基底スペクトルを分類するので、基底スペクトルを高精度に分類することができる。

【図面の簡単な説明】

【0013】

【図1】実施形態に係る音源分離装置のブロック図。

【図2】音響信号および振幅スペクトログラムを説明するための図。

【図3】実施形態に係る音響分離処理のフローチャート。

【図4】第一実施形態に係る基底スペクトルの評価指標を説明する図。

10

20

30

40

50

【図 5】第一実施形態に係る基底番号のソートおよび目的音復元信号の S N R を説明する図。

【図 6】実施形態に係る目的音復元信号を説明する図。

【図 7】実施形態に係る基底スペクトルの分類を説明する図。

【図 8】第二実施形態に係る基底スペクトルの評価指標を説明する図。

【図 9】第二実施形態に係る基底番号のソートおよび目的音復元信号の S N R を説明する図。

【図 10】第三実施形態に係る基底スペクトルの評価指標を説明する図。

【図 11】第三実施形態に係る基底番号のソートおよび目的音復元信号の S N R を説明する図。

10

【発明を実施するための形態】

【0014】

以下、添付の図面を参照して、本発明をその好適な実施形態に基づいて詳細に説明する。なお、以下の実施形態において示す構成は一例に過ぎず、本発明は図示され、以下に説明される構成に限定されるものではない。

【0015】

< 第一実施形態 >

図 1 は、第一実施形態による音響信号処理装置としての音源分離装置の構成例を示すブロック図である。図 1 に示す音源分離装置は、主たるシステムコントローラ 100 の中に、全構成要素の制御を行うシステム制御部 101、各種データを記憶しておく記憶部 102、信号の解析処理を行う信号解析処理部 103 を備える。また、システムコントローラ 100 は、音響信号を入出力するための音響信号入出力部 104 を備える。

20

【0016】

図 1 に示す音源分離装置は、音響信号入出力部 104 を介して、例えば外部の記憶媒体やネットワークから音響信号が入力され、記憶部 102 へ記録される。ここで、音響信号とは、目的音に除去対象となる雑音が混合した混合音を指すものとする。なお、たとえば、不図示のマイクによって収音されたマイク信号に、増幅および A/D 変換が施されたものが音響信号として入力され、記憶部 102 へ逐次記録されるようにしてもよい。

【0017】

第一実施形態においては、目的音を図 2 (a) に示すような音声、雑音を図 2 (b) に示すような風雑音とし、これらが混合した図 2 (c) に示すような混合音を雑音除去の対象である音響信号とする。以下、信号解析処理部 103 が中心となり、図 3 のフローチャートに沿って雑音除去処理が行われる。なお、図 3 のフローチャートは、音響信号の所定の時間ブロック長ごとの処理を表すものとし、第一実施形態では時間ブロック長を 3 秒としている。

30

【0018】

S301 において、信号解析処理部 103 は、音響信号入出力部 104 より入力され記憶部 102 に記憶されている音響信号から時間ブロック長の音響信号を得て、これを時間周波数変換して音響信号の複素スペクトログラム Y を取得する。具体的には、時間ブロック長より短い所定の時間フレーム長ずつ音響信号を切り出して行き、フーリエ変換することで複素フーリエ係数を得る。このとき、切り出される時間フレームは時間フレーム長の半分ずつシフトして行くものとし、時間フレーム長によって時間周波数変換における時間解像度および周波数解像度が決まる。時間ブロックに含まれる時間フレーム数を T、ナイキスト周波数までの周波数分割数を F とすると、上記の処理によって $F \times T$ 個の複素フーリエ係数が得られる。複素スペクトログラム Y は、これらのフーリエ係数を要素とするサイズ $\{F \times T\}$ の複素行列である。なお、フーリエ変換の前に時間信号に対して窓掛けを行うのが好適であり、窓掛けは逆フーリエ変換によって再び時間信号に戻した後にも行う。このため、50% ずつオーバーラップする時間フレームに対し、2 回の窓掛けにおける再構成条件を考慮して、窓関数にはサイン窓などを用いる。

40

【0019】

50

次にS302において、信号解析処理部103は、音響信号の振幅スペクトログラム $|Y|$ を、基底行列Hとアクティビティ行列Uに非負値行列因子分解する。ここで、振幅スペクトログラム $|Y|$ とは、S301で取得した複素スペクトログラムYについて、要素ごとに複素数の絶対値を取った非負値行列である。

【0020】

図2(d)は、図2(c)の音響信号の振幅スペクトログラムを表している。ただし、表示においては振幅値を二値化しており、白が大きい方を、黒が小さい方を表す。図2(d)より、低域では風雑音が卓越しており、また中高域では音声の調波成分による縞模様が見えていることが分かる。

【0021】

非負値行列因子分解に際して指定する基底数をKとすると、サイズ $\{F \times T\}$ の振幅スペクトログラム $|Y|$ を、式(1)のようにサイズ $\{F \times K\}$ の基底行列Hとサイズ $\{K \times T\}$ のアクティビティ行列Uに分解できる。ここで、「 $*$ 」は行列(ベクトル、スカラ含む)の積を表すものとする。

$$|Y| = H * U \quad \dots \quad (1)$$

なお、式(1)の収束計算におけるHとUの更新式は、 $(H * U)$ の $|Y|$ からの乖離度を表す規準に応じたものを用いればよい。行列の乖離度を表す規準としては、ユークリッド距離(二乗誤差)、一般化Kullback-Leiblerダイバージェンス、板倉斎藤距離などが挙げられる。

【0022】

図2(d)の振幅スペクトログラム $|Y|$ を例えば基底数 $K=20$ で非負値行列因子分解すると、サイズ $\{F \times 20\}$ の基底行列Hとサイズ $\{20 \times T\}$ のアクティビティ行列Uが得られる。図4(b)は、基底行列Hの各列を構成する、サイズ $\{F \times 1\}$ の20個の基底スペクトルを正規化して表示したものであり、図の縦軸の番号は基底番号を表している。なお、人の聴感に合わせて周波数軸は対数で、振幅はデシベルで表示している。また、図4(c)は、アクティビティ行列Uの各行を構成する、サイズ $\{1 \times T\}$ の20個のアクティビティベクトルを正規化して表示したものである。

【0023】

ここで、同じ基底番号の基底スペクトルとアクティビティベクトルを掛け合わせることで、サイズ $\{F \times T\}$ の基底別振幅スペクトログラム $|Y_i|$ が式(2)のように得られる。ただし、 $(:, i)$ は行列の第i列を取り出す操作を、 $(i, :)$ は行列の第i行を取り出す操作を表すものとする。

$$|Y_i| = H(:, i) * U(i, :) \quad [i = 1 \sim K] \quad \dots \quad (2)$$

【0024】

さらに、基底別振幅スペクトログラム $|Y_i|$ に複素スペクトログラムYの位相成分を掛け合わせることで、サイズ $\{F \times T\}$ の基底別複素スペクトログラム Y_i が式(3)のように得られる。ここで、「 $\cdot *$ 」は行列の要素ごとの積を、 j は虚数単位を表すものとする。また、 $\arg(Y)$ はYの要素ごとに複素数の偏角を取った行列を表すものとする。

$$Y_i = |Y_i| \cdot \exp(j * \arg(Y)) \quad [i = 1 \sim K] \quad \dots \quad (3)$$

【0025】

そして、基底別複素スペクトログラム Y_i を逆時間周波数変換することで、基底別復元信号 y_i [$i=1 \sim K$]を生成することができる。具体的には、複素スペクトログラムをサンプリング周波数まで対称復元した後、列ごとに逆フーリエ変換することで各時間フレームの復元信号が得られるため、これを窓掛けしてオーバーラップ加算すればよい。

【0026】

図4(d)は、20個の基底別復元信号を正規化して表示したものである。図4(d)の一番上に示す目的音の波形と見比べてみると、例えば基底番号4番や16番は、目的音と雑音のうち目的音に係るものであると予想できる。このように、例えば基底別復元信号を人が見たり聴いたりすれば、各基底を目的音と雑音のものに高精度に分類できそうだが

10

20

30

40

50

、本実施形態では以下のように、各々の基底スペクトルから物理的意味が明確な指標を算出することで自動的に行う。

【 0 0 2 7 】

S 3 0 3 において、信号解析処理部 1 0 3 は、各々の基底スペクトルから算出する評価指標に基づいて、基底番号のソートを行う。具体的には、風雑音の基底スペクトルなら低域が卓越しており、音声の基底スペクトルなら中高域含む広い周波数範囲に分布しているであろうという考え方に基づく。そこで、このような基底スペクトルの周波数軸上での分布状態を数値化するため、第一実施形態では各々の基底スペクトルの評価指標として、各基底スペクトルの周波数分布における重心周波数を算出する。

【 0 0 2 8 】

まず、図 4 (b) に示すような基底スペクトルの表現を得るため、基底スペクトルをデシベル表現する。ただし、微小な値がデシベル表現によって負の大きな値となると不都合であるため、基底行列の最大値 - 6 0 d B 未満の値は、最大値 - 6 0 d B に丸めるなどする。そして、例えば 0 から 1 の範囲内に正規化した後、対数周波数軸上で等間隔に振幅値を得るためオクターブスムージングを行う。

【 0 0 2 9 】

このような表現を行った基底スペクトルを h 、対数周波数軸上の対象周波数範囲（例えば 5 0 ~ 3 k H z ）における等間隔なサンプル点の番号を s (= 0 ~) とすると、重心周波数に対応するサンプル点番号 s_g を式 (4) のように算出することができる。ただし、 $h (s)$ は基底スペクトル h のサンプル点番号 s における値を表し、 \sum は s に関して和を取る操作を表すものとする。

$$s_g = \left(\sum (s * h (s)) \right) / \left(\sum (h (s)) \right) \quad \dots (4)$$

式 (4) で算出される s_g は一般に小数值であり、これを対数周波数軸上に対応付けた値が重心周波数となる。

【 0 0 3 0 】

図 4 (b) の各々の基底スペクトルについて、上記のようにして求めた重心周波数を模式的に示したのが図 4 (b) の黒丸であり、また重心周波数の値を棒グラフで示したのが図 4 (a) である。これらの図より、目的音に係ると予想した基底番号 4 番や 1 6 番の基底スペクトルは広い周波数範囲に分布しており、それゆえ重心周波数も他と比べて高くなっていることが分かる。

【 0 0 3 1 】

図 5 (a) は、図 4 (a) の重心周波数に従って基底スペクトルを重心周波数の昇順にソートしたものであり、横軸がソートされた基底番号を表している。ここで、左の基底番号ほど、基底スペクトルの重心周波数が低いため、低域が卓越した風雑音である可能性が高い。また、右の基底番号ほど重心周波数が高いため、基底スペクトルが広い周波数範囲に分布した音声である可能性が高いと考えられる。

【 0 0 3 2 】

S 3 0 4 において、信号解析処理部 1 0 3 は、基底スペクトルを第一の基底スペクトルとしての目的音基底スペクトルと、第二の基底スペクトルとしての非目的音基底スペクトル（雑音基底スペクトルともいう）とに分類する。そして、分類された基底スペクトルを用いて音響復元信号を生成する。はじめに、信号解析処理部 1 0 3 は、S 3 0 3 でソートされた基底番号に基づいて、基底行列 H の各列（基底スペクトル）を並べ替える。すなわち、図 5 (a) に示されるソート結果に従って元の基底行列の第 1 5 列を並べ替え後の基底行列の第 1 列とし、元の基底行列の第 1 2 列を並べ替え後の基底行列の第 2 列とする、といった具合に 2 0 個の基底スペクトルを並べ替える。また、アクティビティ行列 U の各行（アクティビティベクトル）も同様に並べ替える。

【 0 0 3 3 】

このように、ソートされた基底番号に基づいて基底行列を並べ替えれば、後は目的音基底数または雑音基底数を定めることで、基底スペクトルを目的音基底スペクトルと雑音基底スペクトルに分類することができる。すなわち、雑音基底数を K_n とすれば、並べ替え

10

20

30

40

50

られた基底行列 H の第 1 列から第 K_n 列までの基底スペクトルが雑音基底スペクトルとして、第 $K_n + 1$ 列から第 K 列までの基底スペクトルが目的音基底スペクトルとして分類される。そして、雑音基底スペクトルで構成される雑音基底行列 H_n と、目的音基底スペクトルで構成される目的音基底行列 H_s が、それぞれ式 (5) と式 (6) のように得られる。ただし、 $(: , 1 : K_n)$ は行列の第 1 列から第 K_n 列を取り出す操作を、 $(: , K_n + 1 : K)$ は行列の第 $K_n + 1$ 列から第 K 列を取り出す操作を表すものとする。

$$H_n = H (: , 1 : K_n) \quad \dots \quad (5)$$

$$H_s = H (: , K_n + 1 : K) \quad \dots \quad (6)$$

【0034】

また、基底スペクトルと同様に、アクティビティベクトルも第一のアクティビティベクトルとしての目的音アクティビティベクトルと第二のアクティビティベクトルとしての雑音アクティビティベクトル（非目的音アクティビティベクトル）に分類できる。雑音アクティビティベクトルで構成される雑音アクティビティ行列 U_n と、目的音アクティビティベクトルで構成される目的音アクティビティ行列 U_s が、それぞれ式 (7) と式 (8) のように得られる。ただし、 $(1 : K_n , :)$ は行列の第 1 行から第 K_n 行を取り出す操作を、 $(K_n + 1 : K , :)$ は行列の第 $K_n + 1$ 行から第 K 行を取り出す操作を表すものとする。

$$U_n = U (1 : K_n , :) \quad \dots \quad (7)$$

$$U_s = U (K_n + 1 : K , :) \quad \dots \quad (8)$$

【0035】

目的音基底数を $K_s (= K - K_n)$ とすると、サイズ $\{F \times K_s\}$ の目的音基底行列 H_s と、サイズ $\{K_s \times T\}$ の目的音アクティビティ行列 U_s を掛け合わせることで、サイズ $\{F \times T\}$ の目的音振幅スペクトログラム $|Y_s|$ が式 (9) のように得られる。

$$|Y_s| = H_s * U_s \quad \dots \quad (9)$$

【0036】

さらに、複素スペクトログラム Y の位相成分を掛け合わせることで、サイズ $\{F \times T\}$ の目的音複素スペクトログラム Y_s が式 (10) のように得られる。

$$Y_s = |Y_s| \cdot \exp(j * \arg(Y)) \quad \dots \quad (10)$$

【0037】

そして、目的音複素スペクトログラム Y_s を逆時間周波数変換することで、音響復元信号として目的音復元信号 y_s を生成することができる。なお、サイズ $\{F \times K_n\}$ の雑音基底行列 H_n と、サイズ $\{K_n \times T\}$ の雑音アクティビティ行列 U_n を用いて、音響復元信号としての雑音復元信号 y_n も同様に生成することができる。

【0038】

図 6 (a) は、雑音基底数を 0 から 20 まで増やして行ったときの、それぞれの目的音復元信号を示したものである。図 6 (a) より、雑音基底数 K_n を多くするにつれて、言い換えれば目的音基底数 K_s を絞って行くにつれて、風雑音が除去されて音声が増えられて行く様子が分かる。また図 5 (b) は、雑音基底数と目的音復元信号の SNR の関係を示したグラフであり、雑音基底数 $K_n = 17$ (目的音基底数 $K_s = 3$) のとき、SNR が最大の 2.21 dB となる。このとき図 5 (a) より、基底番号 16 番、4 番、7 番が用いられていることが分かる。しかしながら図 5 (b) より、さらに雑音基底数を増やして目的音基底数を絞って行くと、SNR が低下してしまう。このため、目的音基底数または雑音基底数を適切に定めることが大切と考えられる。

【0039】

目的音基底数の決め方としては、重心周波数は $[Hz]$ の単位を持つ物理的意味が明確な指標であるため、例えば閾値を 200 Hz とし、重心周波数が閾値以上の基底スペクトルの数を目的音基底数としてもよい。図 7 (c) は、図 4 (a) をヒストグラム化して重心周波数の分布を表したものであり、閾値を 200 Hz とすると図 7 (c) の実線で分けることになるため、目的音基底数は 3 となる。また、図 7 (c) のヒストグラムを混合正規分布と見なして、EM アルゴリズムにより 2 群に分類することで目的音基底数を決

10

20

30

40

50

定してもよい。

【 0 0 4 0 】

他にも、音響信号とは別に音声と風雑音を用意し、それぞれから求めた重心周波数のヒストグラムを利用してもよい。例えば、図 7 (b) に示す風雑音のみの音響信号から得られたヒストグラムの範囲を考慮して、図 7 の実線で図 7 (c) を分けると目的音基底数は 3 となる。または、図 7 (a) に示す音声のみの音響信号から得られたヒストグラムの範囲を考慮して、図 7 の点線で図 7 (c) を分けると目的音基底数は 4 となる。この方法は、特許文献 1 の事前学習に似ているようにも見えるが、別に用意する目的音または雑音は、目的音基底数または雑音基底数の決定に用いるだけであり、基底行列は音響信号から得るため特許文献 1 とは異なる。

10

【 0 0 4 1 】

なお、システム制御部 1 0 1 と相互に結ばれた不図示の入出力 G U I 部 (例えばタッチパネルで構成される) を介して、ユーザが目的音基底数または雑音基底数を調整できるようにしてもよい。

【 0 0 4 2 】

以上のようにして、目的音基底数または雑音基底数を適切に定めて生成された目的音復元信号は、記憶部 1 0 2 へ記録される。記録された目的音復元信号は、音響信号入出力部 1 0 4 を介して外部に出力されたり、 D A 変換および増幅を行ったのち、不図示のイヤホン、ヘッドホン、スピーカ等によって再生されたりする。

20

【 0 0 4 3 】

< 第二実施形態 >

第一実施形態では、 S 3 0 3 において、各々の基底スペクトルから重心周波数という評価指標を算出したが、評価指標の算出はこれに限られるものではない。第二実施形態においては、基底スペクトルを変換することで得られるケプストラム (以下、基底ケプストラムと呼ぶ) から評価指標を算出する例を説明する。

【 0 0 4 4 】

基底ケプストラムは、サンプリング周波数まで対称復元した基底スペクトルに対し、対数を取って逆フーリエ変換した結果の実部として得られる。図 8 (c) は、図 8 (b) の各々の基底スペクトルから求めた基底ケプストラムを正規化して表示したものであり、横軸は時間の次元を持つケフレンシーとなる。なお、図 8 (b)、(d) は、それぞれ図 4 (b)、(d) と同じものである。

30

【 0 0 4 5 】

図 8 (c) において、基底番号 4 番や 1 6 番の基底ケプストラムは、図 8 (b) の点線丸で示すように低ケフレンシー部分が一際大きくなっている。ここで、一般にケプストラムの低ケフレンシー部分はスペクトルの包絡成分に対応するが、スペクトルの包絡成分の大きさは、スペクトルの広がり状態を表していると考えられる。実際に、基底ケプストラムにおいて低ケフレンシー部分に大きい値を持つ基底番号 4 番や 1 6 番では、図 8 (b) の点線丸で示すように、基底スペクトルが広い周波数範囲に分布していることが確認できる。なお、低ケフレンシー部分とは、ケフレンシーが所定値以下の部分であり、本実施形態では、たとえば 2 m s 以下の部分とする。

40

【 0 0 4 6 】

そこで、第二実施形態では、 S 3 0 3 において基底スペクトルの周波数軸上での分布状態を数値化するための評価指標を、各々の基底ケプストラムから算出する。例えば、基底スペクトルの広がり状態を表す包絡成分の大きさを数値化するため、基底ケプストラムの所定のケフレンシー以下の部分について、その最大値を評価指標とする。より簡単には、図 8 (c) より基底ケプストラムの低ケフレンシー部分の最大値は、基底ケプストラム全体の最大値であるとして差し支えなさそうであるため、これを評価指標として用いてもよい。

【 0 0 4 7 】

図 8 (c) の各々の基底ケプストラムについて、基底ケプストラム全体の最大値を棒ゲ

50

ラフで示したのが図 8 (a) であり、これを昇順にソートしたものが図 9 (a) である。図 8 (a) において、棒グラフが右へ延びる基底番号ほど基底スペクトルの包絡成分が大きい。逆に、棒グラフが左側にとどまる基底番号ほど、基底スペクトルの包絡成分が小さいため、基底スペクトルが狭い周波数範囲 (低域) に集中した風雑音である可能性が高いと考えられる。

【 0 0 4 8 】

ここで、第一実施形態と第二実施形態はともに、基底スペクトルの周波数軸上での分布状態を数値化するという考えに基づいている。このため、図 5 (a) と図 9 (a) の棒グラフは全体的に傾向が似ており、ソートされた基底番号の並び順も類似している。特に、上位 4 つ (1 6 番、4 番、7 番、2 番) は同じである。そのため、第二実施形態における雑音基底数と目的音復元信号の S N R の関係を示した図 9 (b) は、図 5 (b) と同じく雑音基底数が 1 7 (目的音基底数が 3) のときに S N R が最大の 2 . 2 1 d B となる。

【 0 0 4 9 】

なお、S 3 0 4 の目的音基底数または雑音基底数の決定において、図 7 の実線や点線のように評価指標の値で分ける場合、音響信号の大きさに評価指標の値が依存しないことが望ましい。第一実施形態の重心周波数はその性質上、音響信号の大きさには依存しないが、一般に基底スペクトルや基底ケプストラムの大きさは音響信号の大きさに依存する。そこで、式 (1) の非負値行列因子分解において基底行列 H を正規化しておけば、基底スペクトルや基底ケプストラムが音響信号の大きさに依存しなくなり、基底ケプストラムから算出する評価指標も音響信号の大きさに依存しなくなるため好適である。

【 0 0 5 0 】

なお、図 3 のフローチャートにおいて上述した S 3 0 3 の処理以外は、第一実施形態と同様である。

【 0 0 5 1 】

< 第三実施形態 >

上述の第一実施形態、第二実施形態では、S 3 0 3 において、基底スペクトルの周波数軸上での分布状態を数値化した評価指標を用いた。第三実施形態においては、音声の基底スペクトルなら風雑音の基底スペクトルより調波成分が大きいうという考え方に基づき、このような調波成分の大きさを数値化するため、各々の基底ケプストラムから評価指標を算出する。

【 0 0 5 2 】

一般にケプストラムのピークは、スペクトルの調波成分の大きさと、その基本周波数を示している。例えば、ケプストラムがケフレンシー 5 m s の位置にピークを持てば、スペクトルはその逆数の 2 0 0 H z を基本周波数とする調波成分を持つ。

【 0 0 5 3 】

基底スペクトルの調波成分の大きさを数値化するためには、基底ケプストラムのピークの大きさを調べればよい。簡単には基底ケプストラムの最大値を算出することが考えられる。しかしながら第二実施形態で述べたように、基底ケプストラム全体の最大値は、実際には低ケフレンシー部分の最大値であるため、結局は基底スペクトルの周波数軸上での分布状態を見ていることになる。そこで第三実施形態では、基底ケプストラムの音声の基本周波数範囲に対応する部分においてその最大値を取ることで調波成分の大きさを数値化した評価指標とする。

【 0 0 5 4 】

図 1 0 (c) は、図 1 0 (b) の各々の基底スペクトルから求めた基底ケプストラムについて、音声の基本周波数範囲である 1 0 0 ~ 4 0 0 H z に対応する部分、すなわちケフレンシーで 2 . 5 ~ 1 0 m s の部分を拡大して表示したものである。なお、図 1 0 (b) 、(d) は、それぞれ図 4 (b) 、(d) と同じものである。

【 0 0 5 5 】

図 1 0 (c) の点線丸で示すように、目的音に係ると予想される基底番号 4 番や 1 6 番

10

20

30

40

50

の基底ケプストラムは、音声の基本周波数範囲に対応する部分にピークを持っており、このようなピークは基底番号 7 番や 11 番の基底ケプストラムにも見られる。

【0056】

音声の基本周波数範囲に対応する部分である、図 10 (c) の各々の基底ケプストラムについて最大値を棒グラフで示したのが図 10 (a) であり、これを昇順にソートしたものが図 11 (a) である。図 11 (a) において、右の基底番号ほど、基底スペクトルの調波成分が大きいと音声である可能性が高く、逆に左の基底番号ほど、基底スペクトルの調波成分が小さいと風雑音である可能性が高いと考えられる。なお、上位 4 つ (7 番、4 番、11 番、16 番) の組み合わせを見ると、上記第一、第二実施形態における 2 番の代わりに、11 番が含まれていることが分かる。

10

【0057】

図 11 (b) は、本実施形態における雑音基底数と目的音復元信号の S N R の関係を示したグラフである。第一、第二実施形態と少し異なり、雑音基底数が 16 (目的音基底数が 4) のとき S N R が最大の 2.98 dB となっており、上記第一、第二実施形態における 2.21 dB よりも高い。これは、先に述べた基底番号 11 番が用いられたためと考えられる。特に、図 10 (d) の基底別復元信号で示される点線丸の部分が加わることで、第三実施形態の目的音復元信号を示した図 6 (b) において、図 6 (a) と比較して点線丸の部分の音声は復元されていることが分かる。

【0058】

なお、特許文献 2 には、雑音復元信号から調波成分を抽出して目的音復元信号に合成する処理があるが、このような処理は特に、目的音基底スペクトルが雑音基底スペクトルとして分類されてしまった場合に必要になると考えられる。第三実施形態は、基底スペクトルの分類における評価指標を調波成分の大きさとするすることで、上記のような分類ミスを防止しようとするものであり、特許文献 2 とは異なる。

20

【0059】

なお、時間ブロック長は 3 秒として説明を行ってきたが、第三実施形態において音声の音素 (好ましくは母音) ごとに基底スペクトルを得るために、例えば 0.3 秒といった短い時間ブロック長を用いてもよい。この場合、非負値行列因子分解における行列サイズが縮小されるため、計算時間の短縮にもつながる。

【0060】

なお、図 3 のフローチャートにおいて上述した S 303 の処理以外は、第一実施形態と同様である。

30

【0061】

ここで、風雑音に音声は少し混ざったような基底スペクトルがある場合、各実施形態どのように分類され得るかを考える。基底スペクトルの概形としては、低域が卓越した風雑音が支配的となるため、第一実施形態では、重心周波数が低くなって雑音基底スペクトルに分類される可能性が高い。また、基底スペクトルが狭い周波数範囲 (低域) に集中しており、基底スペクトルの包絡成分が小さくなるため、第二実施形態においても雑音基底スペクトルに分類されると考えられる。一方、音声の調波成分は含まれているため、第三実施形態では目的音基底スペクトルに分類される可能性がある。以上より、第一実施形態や第二実施形態のように、基底スペクトルの周波数軸上での分布状態を見る方法は、雑音の除去を重視した方法であると解釈することができる。一方で、第三実施形態のように基底スペクトルに含まれる調波成分を見る方法は、目的音の保存を重視した方法であると解釈でき、このような重視する点の違いが、目的音復元信号の S N R の値にも表れたと考えられる。

40

【0062】

以上の各実施形態においては、目的音を音声、雑音を風雑音として説明を行ってきたが、本発明は他の音の組み合わせに対しても適用できることは言うまでもない。第一実施形態や第二実施形態の方法は、たとえ目的音と雑音で調波成分の強さが同程度であっても、基底スペクトルの周波数軸上での分布状態が異なっていればよいと、例えば目的音がセ

50

せらぎの音で、雑音が車のロードノイズであるような場合にも適用できる。また、第三実施形態の方法は、たとえ目的音と雑音で周波数帯が重なっていても、基底スペクトルの調波成分の大きさが異なっていればよい。よって、例えば音声や鳥の鳴き声のように調波成分を持つ目的音と、ざわめきのような雑音の組み合わせに対しても適用可能であるし、調波楽器と打楽器のような音の組み合わせにも用いることができる。このように、基底スペクトルの周波数軸上での分布状態を見る方法と、基底スペクトルに含まれる調波成分を見る方法によって、本発明は様々な音の組み合わせに対応することができる。

【0063】

なお、目的音復元信号と雑音復元信号の少なくとも一方を用いて、元の音響信号から雑音除去を行うようにしてもよい。例えば、雑音復元信号を音響信号からスペクトル減算することで雑音除去を行ってもよいし、目的音復元信号と雑音復元信号から生成したウィナーフィルタを音響信号に適用してもよい。

【0064】

なお、雑音をもう一つの目的音と考えれば、本発明を雑音除去ではなく音源分離に用いることもできる。さらに、ソートされた基底番号を2つではなく3つ以上に分けることで、音響信号を3つ以上の音に分離することも可能である。また、複数 c の音響信号に対しても、 c ごとに本発明の処理を適用できることは言うまでもない。

【0065】

なお、以上の実施形態においては、音響信号の振幅スペクトログラム $|Y|$ を非負値行列因子分解することで、基底行列とアクティビティ行列を得ていたがこれに限られるものではない。たとえば、複素 NMF を用いることで、音響信号の複素スペクトログラム Y を、基底行列とアクティビティ行列および、サイズ $\{F \times T\}$ の K 個の位相スペクトログラム $P_i [i = 1 \sim K]$ に分解することができる。このとき、例えば目的音複素スペクトログラム Y_s は、式 (11) のように算出することになる。ただし、 i は目的音に対応する K 個の i に関して、和を取る操作を表すものとする。

$$Y_s = (H(:, i) * U(i, :)) \cdot P_i \quad \dots \quad (11)$$

【0066】

なお、振幅スペクトログラム $|Y|$ は、複素スペクトログラム Y の要素ごとに複素数の絶対値を取ったものとしていたが、代わりに絶対値の指数乗（例えば0.5乗や2乗）を取ったものとしてもよい。また、時間周波数変換において、フーリエ変換の他にウェーブレット変換などを用いてもよく、その場合はスカログラムが振幅スペクトログラムの代わりとなる。

【0067】

以上説明したように、上記各実施形態によれば、音響信号の各々の基底スペクトルから物理的意味が明確な指標を算出し、算出された指標により基底スペクトルを目的音と雑音に分類するため、音響信号から高精度に雑音を除去することができる。また、単一の音響信号から教師基底なしで高精度に雑音除去することができる。

【0068】

< その他の実施形態 >

以上、実施形態例を詳述したが、本発明は例えば、システム、装置、方法、プログラム若しくは記録媒体（記憶媒体）等としての実施態様をとることが可能である。具体的には、複数の機器（例えば、ホストコンピュータ、インタフェース機器、撮像装置、Webアプリケーション等）から構成されるシステムに適用しても良いし、また、一つの機器からなる装置に適用しても良い。

【0069】

また、本発明の目的は、以下のようにすることによって達成されることはいうまでもない。即ち、前述した実施形態の機能を実現するソフトウェアのプログラムコード（コンピュータプログラム）を記録した記録媒体（または記憶媒体）を、システムあるいは装置に供給する。係る記憶媒体は言うまでもなく、コンピュータ読み取り可能な記憶媒体である。そして、そのシステムあるいは装置のコンピュータ（またはCPUやMPU）が記録媒

10

20

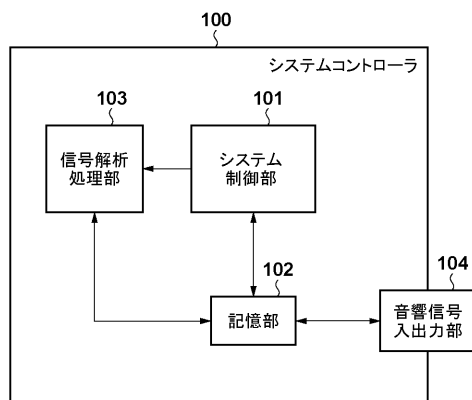
30

40

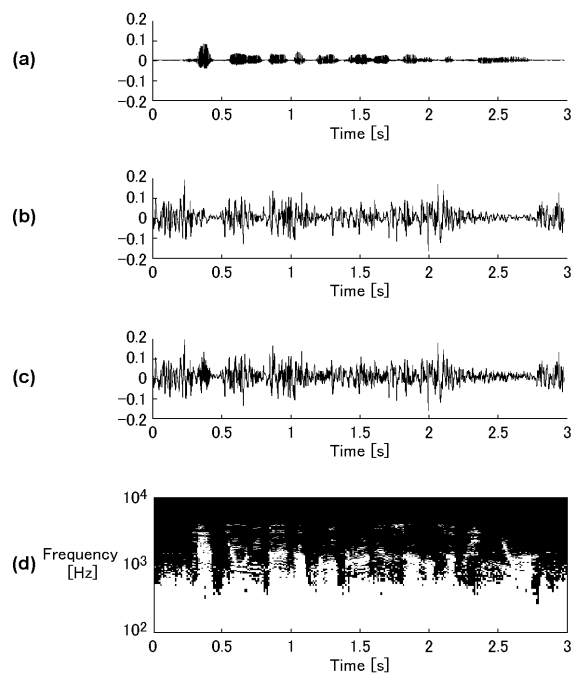
50

体に格納されたプログラムコードを読み出し実行する。この場合、記録媒体から読み出されたプログラムコード自体が前述した実施形態の機能を実現することになり、そのプログラムコードを記録した記録媒体は本発明を構成することになる。

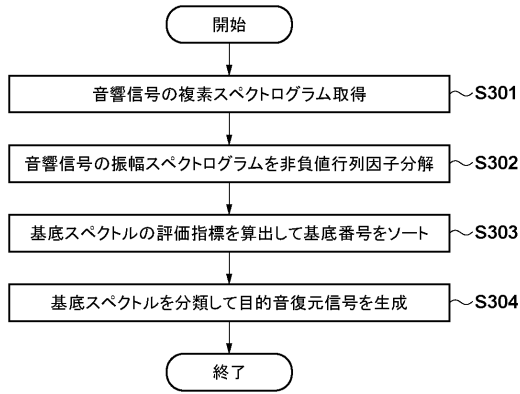
【図 1】



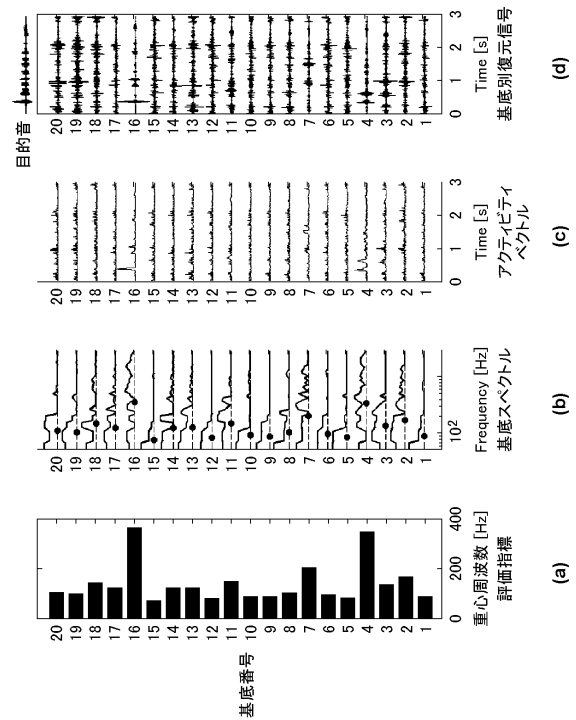
【図 2】



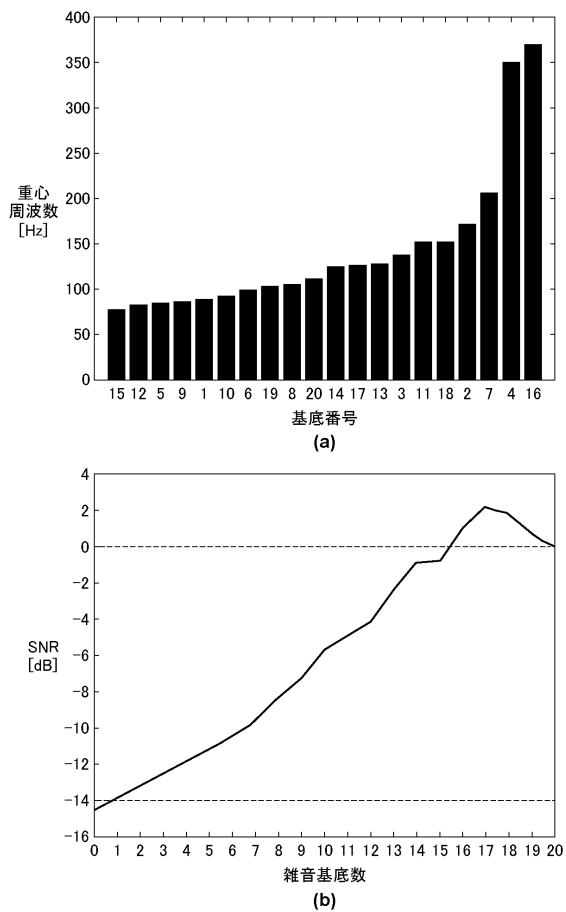
【図 3】



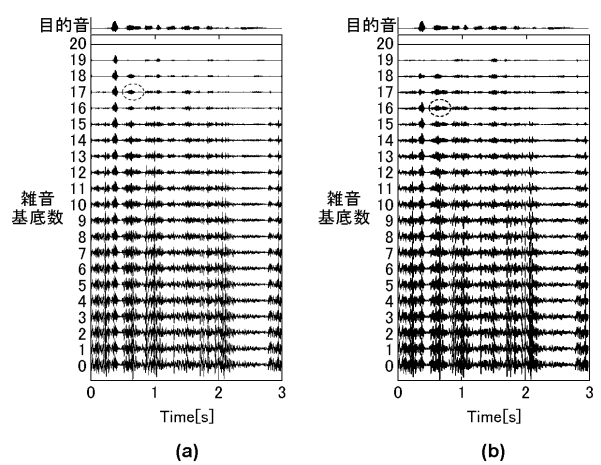
【図 4】



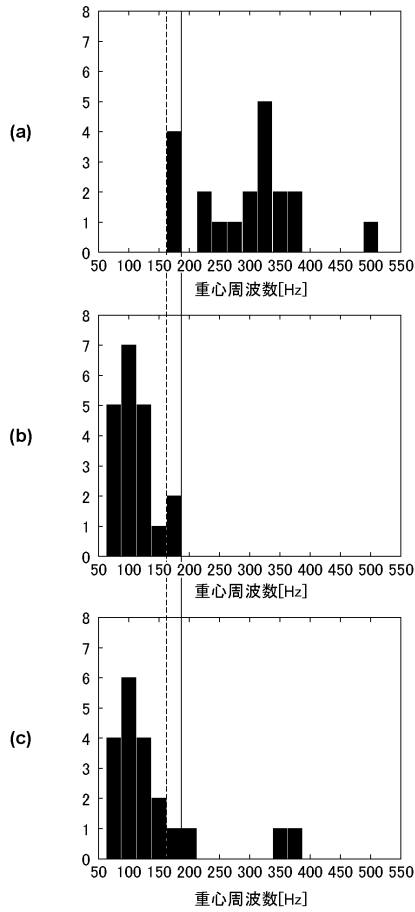
【図 5】



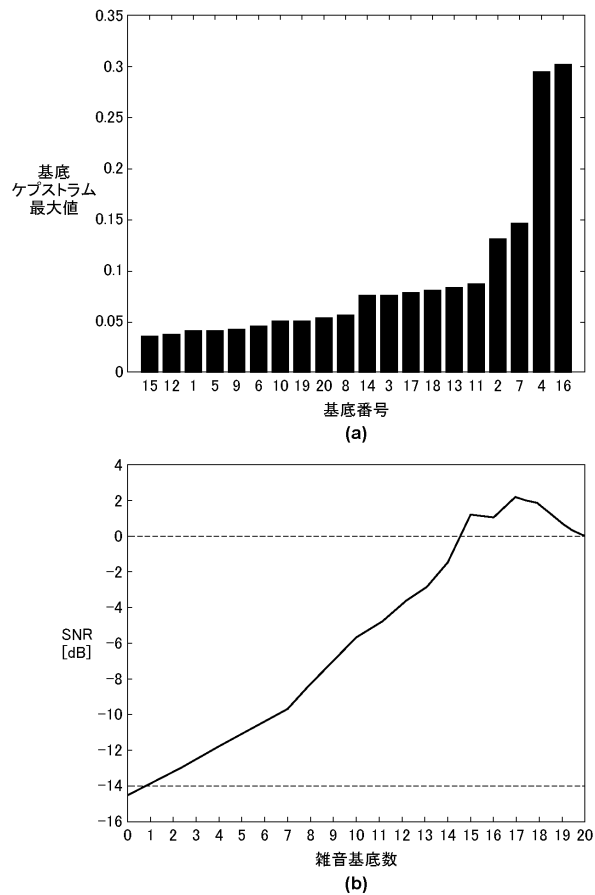
【図 6】



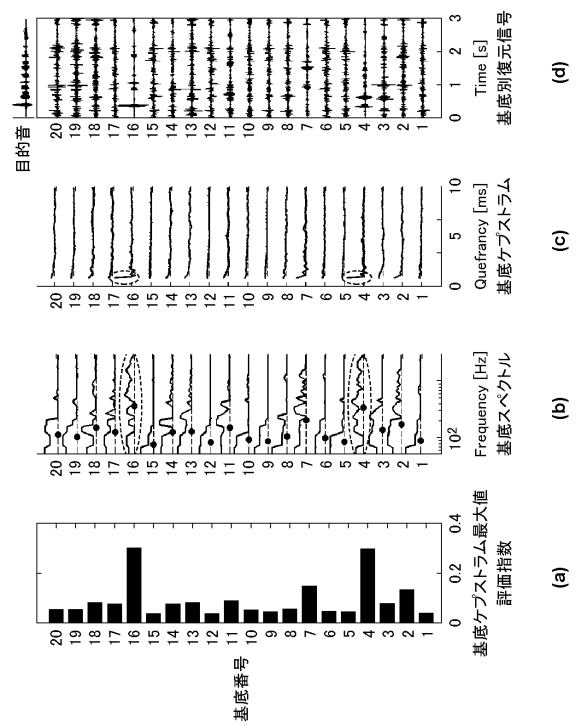
【図 7】



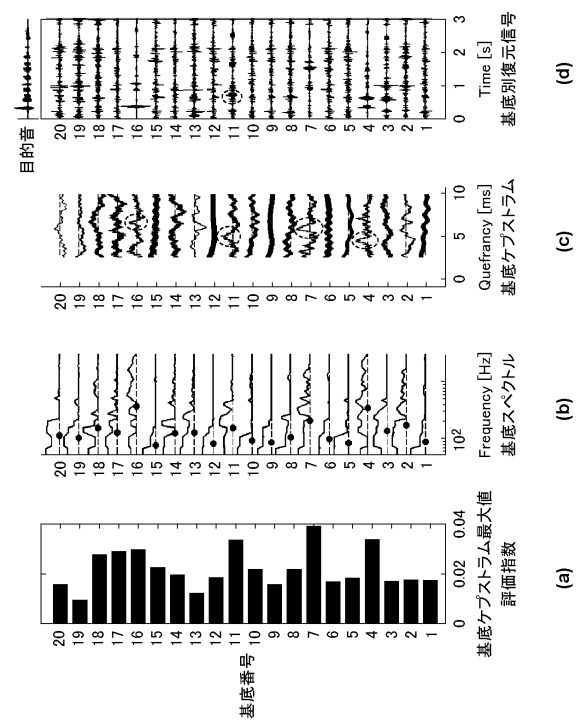
【図 9】



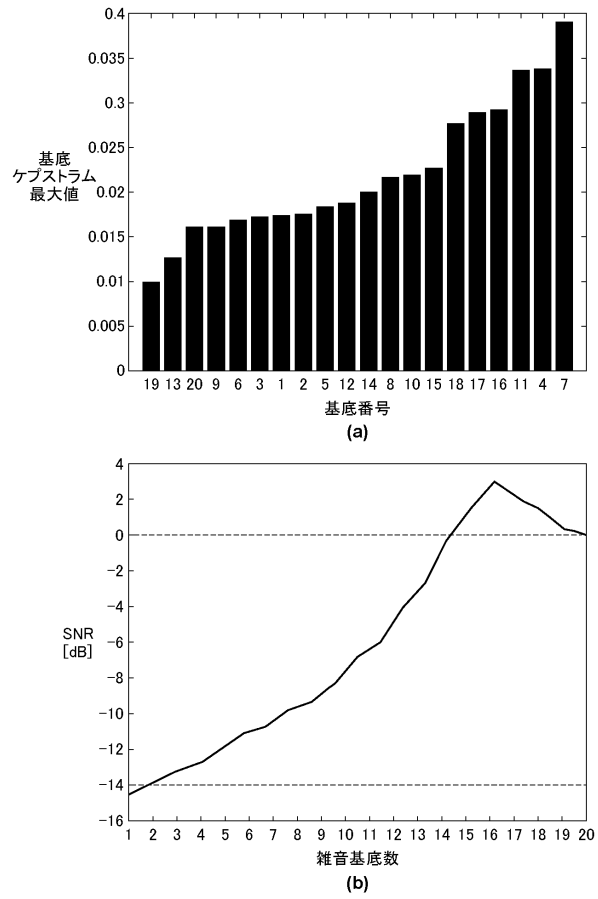
【図 8】



【図 10】



【図 11】



フロントページの続き

(72)発明者 多和田 典朗
東京都大田区下丸子3丁目30番2号 キヤノン株式会社内

審査官 上田 雄

(56)参考文献 特開2013-037152(JP,A)
特開2012-163918(JP,A)
特開2013-033196(JP,A)
特開2012-133346(JP,A)
特開2006-243178(JP,A)
Marko Helen and Tuomas Virtanen, SEPARATION OF DRUMS FROM POLYPHONIC MUSIC USING NON-NEGATIVE MATRIX FACTORIZATION AND SUPPORT VECTOR MACHINE, Proc. 13th European Signal Processing Conference, IEEE, 2005年 9月

(58)調査した分野(Int.Cl., DB名)
G10L 21/00-21/18
G10L 25/00-25/93