

# 実数シフトのレゾルベント少数から構成されたフィルタによる実対称定値一般固有値問題の下端側固有値を持つ固有対の解法について

村上 弘<sup>1,a)</sup>

**概要:** 実対称定値一般固有値問題の固有対で固有値が下端付近にあるものを近似して解くためのフィルタを、実数シフトのレゾルベントを少数用いて構成する方法について考究する。既に我々は、シフトを複素数の範囲で選ぶのであれば、少数のレゾルベントを用いて容易に優れた特性のフィルタを構成することができて、しかも求めたい固有対の固有値が入る区間は任意に設定できることを示した。しかしシフトが虚数の場合には、レゾルベントの作用を与える連立1次方程式は係数が複素対称行列であり、それを解くには複素数の記憶と演算が必要になる。応用上は、固有値が下端付近にある固有対だけを求めたいことがよくあり、そのような場合に、最小固有値の下側の実数をシフトとするレゾルベントを少数用いてフィルタを構成すると、各レゾルベントに対応する連立1次方程式の係数行列は実対称定値になる。我々は既に、そのようなレゾルベントを1つあるいは2つ用いたフィルタの構成法を示した。今回はそのようなフィルタを簡単な固有値問題の例題に適用して、下端付近の固有値を持つ固有対を近似して求めた実験の例を示す。

**キーワード:** フィルタ, 対角化法, 固有値問題, レゾルベント, 多項式, 実数シフト, 伝達関数

## 1. はじめに

行列  $A$  と  $B$  は実対称で、 $B$  は正定値である実対称定値一般固有値問題 (1) の固有対  $(\lambda, \mathbf{v})$  で指定された区間  $[a, b]$  に固有値  $\lambda$  があるものの近似を求めることにする。

$$A\mathbf{v} = \lambda B\mathbf{v} \quad (1)$$

そのためのフィルタとして、少数  $k$  個のレゾルベントの線形結合の作用の Chebyshev 多項式の作用 (2) を採用する。

$$\mathcal{F} = g_s T_n(\mathcal{Y}) \quad (2)$$

ここで作用素  $\mathcal{Y}$  は、シフトが  $\rho_i$  のレゾルベント  $\mathcal{R}(\rho_i)$ ,  $i = 1, 2, \dots, k$  と恒等作用素  $I$  の線形結合とする (式 (3))。

$$\mathcal{Y} = c_\infty I + \sum_{i=1}^k c_i \mathcal{R}(\rho_i). \quad (3)$$

既に我々は文献 [29] において、複素数をレゾルベントのシフトとして用いる場合に、レゾルベントの数を増すことでフィルタの伝達関数  $f(\lambda)$  の形状を系統的に良くできる

<sup>1</sup> 東京都立大学・数理学専攻  
Department of Mathematical Sciences, Tokyo Metropolitan University

a) mrkhrsh@tmu.ac.jp

方法を示した。その方法ではレゾルベントの数  $k$  がいくつでも、各レゾルベントのシフトとその線形結合の係数は数式に数値を入れて計算することで求められることを示した。特に、すべてのシフトが虚数になる場合は、区間  $[a, b]$  の位置は自由に設定できて、中間固有値の固有対も特に困難なく求められる。

しかし、必要な固有対は固有値が固有値分布の端付近のものであることが応用上はよくある。そのような場合には、レゾルベントのシフト  $\rho_i$  を実数に制限することにより、計算に必要な記憶量と演算量を減らせる可能性がある。

なぜならば、計算を複素数で行う場合は実数で行う場合に比べると必要な記憶量は倍になる。また演算については、複素数の加算は実数の加算2つで構成され、通常の方法による複素数の乗算は実数の乗算4つと実数の加算2つで構成される。行列分解や前進後退代入などの計算を行う場合に、乗加算の形で加算と乗算がほぼ1対1の割合で含まれているとし、また計算機の演算装置も実数の加算器と乗算器が1対1の割合で備わっていて実数の加算と乗算の演算の手間は同じであると仮定すると、計算を複素数で行う場合は実数で行う場合に比べて演算の手間は4倍になる。

さらに今回は一般固有値問題を実対称定値の場合に限っているが、複素エルミート定値の一般固有値問題も固有値は実数に限られて、フィルタ構成の議論は実対称定値の場合とほぼ同様にできるが、レゾルベントの作用を実現する連立1次方程式の係数行列は、シフトが実数である場合は複素エルミートであり、行列分解で対称性を利用できるが、シフトが虚数である場合は特別な対称性のない複素行列になり、行列分解の計算量や分解の記憶量などがシフトが実数である場合と比べて増える。

レゾルベントのシフトに実数を用いる場合は、もしもシフトと一致もしくは近接する固有値が存在すると、求めたい固有対全体に対するフィルタによる伝達率の最大最小比が極端に大きくなり、精度の限られた数値と演算を用いる通常の計算では、得られる近似固有対の精度のばらつきが大きくなるリスクがある [10]。そのようなリスクは、必要とする固有対の固有値が固有値分布の下端付近である場合には、シフトを最小固有値よりも下側にとれば回避できる。しかしシフトを実数に制限することで選択範囲を狭めているので、シフトを複素数から選べる場合に比べて達成可能なフィルタの特性は必ず劣ったものになる。

我々はこれまで、レゾルベント1つだけから構成された極めて簡易なフィルタを扱ってきた (文献 [11], [12], [13], [14], [23])。たとえばレゾルベントの作用を実現する連立1次方程式を行列分解を利用して解くことを前提にすると、使用するレゾルベントの数が  $k$  ならば行列分解を  $k$  通り行なう必要があるため、行列分解に掛かる演算量はレゾルベントが複数の場合に比べて1つの場合が最も少なくなる。さらに行列分解の結果を保持することにより、右辺だけが異なる連立1次方程式の組を解く処理を多項式の次数に等しい回数だけ繰り返すのに掛かる計算量も少なくでき、それに用いる行列の分解結果を保持するための記憶量も単一のレゾルベントを用いる場合が複数用いる場合に比べて少なくなる。そのため、単一のレゾルベントを用いるフィルタは、行列分解を格納するための記憶量が計算実行上の制約となるような大規模な問題では利点がある。しかし単一のレゾルベントで構成されたフィルタは、複数で構成されたものに比べて、伝達関数の特性があまり良くないという難点がある。たとえば通過域における伝達関数の値の最大最小比を抑えながら遷移域の幅を狭くすることは難しい。しかし我々は、特性があまり良くないフィルタであっても、再直交化とフィルタを組み合わせた処理を数回反復することで不変部分空間の基底の近似を改良して、必要な固有対の近似精度を一斉に向上できることを示した (文献 [22], [24], [25], [27])。ここまでを読むと、単一のレゾルベントを用いてフィルタを構成しても問題がうまく解けるのであれば、あえて複数を用いることは必要無いと思われるであろう。しかしそれでも、小さい並行数の分散計算が

できるシステムであって、少数の係数行列の分解結果が主記憶全体に収まり、各行列の分解や分解した後の前進後退代入に必要な処理をそれぞれほぼ独立に並行して計算ができる場合であるなら、以下のような利点が生じる可能性がある (この事情は、フィルタとして少数のレゾルベントの線形結合の実部の Chebyshev 多項式を用いた場合 (文献 [26], [28], [29]) と同様である。ただしこれらの文献はシフトに実数ではなく複素数を用いた場合である)。

フィルタとして少数のレゾルベントの線形結合の作用の Chebyshev 多項式を採用することで、単一のレゾルベントを用いた場合に比べて：

- Chebyshev 多項式の次数  $n$  を減らすことができれば、それにより (少数  $k$  個のレゾルベントの処理がそれぞれ並行して処理できると仮定すれば) フィルタを1回適用する処理の中でレゾルベントを逐次的に  $n$  回反復して適用する部分 (行列分解の後に、前進後退代入を  $n$  回逐次的に繰り返して実現する) の経過時間が減る。
- 伝達関数の遷移域の幅を狭めることができれば、フィルタをベクトルの組に適用する際にベクトルの数を減らせる。
- 伝達特性の閾値を改善できれば ( $g_s$  をより微小にする、あるいは  $g_p$  をより大きく1に近づける)、それにより近似固有対の精度が向上する。たとえばフィルタの適用1回の段階で既に、用途が要求する精度の水準を近似固有対が満たしていれば、フィルタの反復を追加して近似を改良する手間が省ける。

そこで固有値が下端付近の固有対を近似して求める場合について、フィルタは実数シフトの少数のレゾルベントの線形結合の Chebyshev 多項式であるとして、なるべくその特性が良くなるように構成する方法を導くことにする。

我々は以前の文献 [15] において、フィルタを2つのレゾルベントの線形結合の (実部の) Chebyshev 多項式として構成するのに、シフト2つが共に実数である場合と虚数である場合 (ただし複数共役性を利用してシフトの虚部が正のもの2つ) のそれぞれについて考察を行った。文献 [15] で扱っている実数シフトのレゾルベントを2つ用いる場合のフィルタの「方式1」と「方式2」は、本報告の「方式I」、「方式II」と同じものである。

その文献 [15] のシフトが虚数のレゾルベントを2つ用いる場合の構成法は、その後の我々の一連の研究 (文献 [16], [17], [18], [19], [20], [21], [26], [28], [29]) において拡張され、アナログ電気回路におけるフィルタの設計手法を模倣して、最良近似理論に現れる有理関数を利用してフィルタの伝達関数を関数合成の手法で設計する方法を示した。その方法では、シフトは複素数であり、レゾルベントの数をいくつにしても数式に数値を入れて計算することでフィルタを具体的に決定できる。そうして優れた特性を持つフィルタが少数3-4個のレゾルベントを用いて構成

できることを示し、その確認のための実験も行なった [29]. しかしその方法はシフトを実数に限定することはできないので、新たな方法が必要になる. 本報告では (複数の始まりは 2 つであるから), 既に文献 [15][30] の中で扱った実数シフトのレゾルベント 2 つで構成される簡易型のフィルタを, 文献 [31] に続いて実際に用いて実験してみた. (実数シフトのレゾルベントを 3 つあるいは 4 つ用いる場合の構成は, まだ十分に検討できておらず, 今後の課題とする).

## 2. 実験について

### 2.1 例題に用いた一般固有値問題

例題とした実対称定値一般固有値問題 (1) は, 1 辺の長さ  $\pi$  の 3 次元立方体の内部を領域として, その表面において零-ディリクレ境界条件を課したときの (符号が逆の) 3 次元ラプラシアン  $-\Delta$  の固有値問題, それを有限要素法 (FEM) で離散化したものである.

FEM の要素分割は立方体の各辺方向をそれぞれ  $N_1 + 1$ ,  $N_2 + 1$ ,  $N_3 + 1$  の等間隔の小区間に分割したもので, 要素内での展開基底関数には各辺方向の 3 重線形関数を用いた. この FEM の離散化で得られる行列  $A$  と  $B$  の次数は  $N = N_1 N_2 N_3$  となり, ( $N_1 \leq N_2 \leq N_3$  であるとして) 基底関数にうまく番号を付けると, 各行列の (対角を含まない) 半帯幅 (下帯幅) は  $w_L = 1 + N_1 + N_1 N_2$  になる.

このようにして得られた実対称定値一般固有値問題 (1) に対して, フィルタ対角化法を適用して, 区間  $[a, b]$  に固有値  $\lambda$  が含まれる固有対を近似して求めた. このテスト例題の固有値は簡単な数式で表せるので, 式に値を入れて計算すれば厳密な固有値が簡単に求められる. またそのことを用いて, 固有値の厳密値を列挙して値の大小順に並べて数え上げることで, 固有値が区間  $[a, b]$  にある固有対の正しい数も求まる.

### 2.2 近似固有対の評価に用いた相対残差

計算で求めた近似固有対の品質の評価には相対残差を用いた. 近似固有対  $(\lambda, \mathbf{v})$  に対する相対残差  $\Theta$  を式 (4) で定義する. ベクトルのノルム  $\|\cdot\|$  には 2-ノルムを使用した. この  $\Theta$  の値はベクトル  $\mathbf{v}$  の規格化には依らず, また共通の非零の値で行列  $A$  と  $B$  をスケールしても不変である. 幾何学的には,  $N$  次元ユークリッド空間内で 2 つのベクトル  $A\mathbf{v}$  と  $\lambda B\mathbf{v}$  が挟む角の大きさを  $\phi$  とするとき, 不等式  $\sin \phi \leq \Theta$  が成り立つ.

$$\Theta \equiv \frac{\|A\mathbf{v} - \lambda B\mathbf{v}\|}{\|\lambda B\mathbf{v}\|}. \quad (4)$$

相対残差の計算では, まず複数の近似固有対のベクトルを列として並べた行列を  $V$  とし, それから行列  $A$  と  $B$  の対称性や帯性または疎性も利用して計算して行列積  $AV$  と  $BV$  を作る. このように行列積を計算する形にまとめることで,  $A$  と  $B$  への記憶参照は全体で 1 回ずつになり, 相対

残差を個別に求めるよりも計算の効率が良くなる.

### 2.3 直交化付きフィルタの反復による近似固有対の改良

計量  $B$  の正規直交化とフィルタの適用を組み合わせたものを少数 (IT\_MAX) 回反復することで, フィルタの伝達特性の形状の悪さを補って, 不変部分空間の基底の近似を改善できる [20], [27]. その処理の概要は以下ようになる.

- (1) 通過域が  $[a, b]$  であるフィルタ  $\mathcal{F}$  を用意
- (2) 乱数から作成した  $m$  個のベクトルの組を  $Y$  とする.
- (3) for IT=1, IT\_MAX do
- (4)  $X \leftarrow$  「 $Y$  の切断付き  $B$ -正規直交化」
- (5)  $Y \leftarrow \mathcal{F} X$
- (6) enddo
- (7)  $X$  と  $Y$  とフィルタの特性を考慮して,  $Y$  の線形結合で不変部分空間の基底の近似  $Z$  を作成 [3].
- (8) 式 (1) の一般固有値問題に対応する Rayleigh-Ritz 法を  $Z$  に適用して, 得られた Ritz 対を近似固有対とする.

上記処理中の (4) 「 $Y$  の切断付き  $B$ -正規直交化」の処理で  $Y$  の実効階数の低下を検出したら,  $X$  の持つベクトルの数を減らす (上記処理中の (5) により,  $Y$  の持つベクトルの数も  $X$  と同じになる).

### 2.4 実験に用いた計算機システム

実験の計算に用いたシステムは, 東京大学情報基盤センターの Oakbridge-CX の 1 ノード (CPU は Dual で Intel Xeon 8280 (2.7GHz, 28cores), 共有メモリは 192GiB, 倍精度演算でのピーク性能は 4.8TFLOPS) である.

プログラムは Fortran90 を用いて記述し, 並列化のための OpenMP の指示行を適宜追加した. 計算に用いた数値と演算は IEEE 754 の倍精度浮動小数点 (2 進, 64bit) である. コンパイラは intel fortran (version 19.0.5.281) で, コンパイラのオプションとして “-fast -qopenmp -xCORE-AVX512 -align array64byte” を指定した.

## 3. 実験 1

以下の各例では有限要素法 (FEM) で立方体の立方体への要素分割の方式を  $(N_1, N_2, N_3) = (40, 50, 60)$  とした. それにより得られる式 (1) の実対称定値一般固有値問題の行列  $A$  と  $B$  は帯行列で, 行列次数が  $N = 120,000$ , 下帯幅が  $w_L = 2,041$  となる.

この一般固有値問題の最小固有値は, 変分原理から要素分割には依らずに常に 3 よりも大きいので (いまの場合は有効数字 5 桁では 3.0010 である), 固有値が区間  $[a, b] = [3, 30]$  にある固有対を近似して求めることにした. その区間に固有値がある固有対の数は 54 である.

フィルタを適用する最初のベクトルの数  $m$  は、 $\mu = 2.0$  の場合には  $m = 200$ 、 $\mu = 1.5$  の場合には  $m = 125$ 、 $\mu = 1.25$  の場合には  $m = 100$  とした

そのようにした理由は、伝達関数  $g(t)$  の通過域  $t \in [0, 1]$  と遷移域  $t \in (1, \mu)$  を併せた区間  $t \in [0, \mu)$  に対応する固有値の区間  $\lambda \in [a, b)$  は、 $\mu = 2.0$  の場合には  $[3, 57)$ 、 $\mu = 1.5$  の場合には  $[3, 43.5)$ 、 $\mu = 1.25$  の場合には  $[3, 36.75)$  であり、それぞれの場合に、区間  $[a, b)$  に固有値がある固有対の数は 163 個、105 個、78 個であり、それよりも少し大きい値に  $m$  を設定した。

計算で得られた近似固有対で固有値が区間  $[a, b)$  にあるものすべてについて、式 (4) で定義した「相対残差」を計算した。

### 3.1 「方式 I」の伝達関数の例

伝達関数  $g(t)$  の形状パラメタ 3 つと次数の組  $(\mu, g_p, g_s, n)$  を指定して、それを満たす「方式 I」の伝達関数を決定した。対応するフィルタの作用は伝達関数から容易に導ける (付録 A.1.6)。

実験に用いた「方式 I」の 6 通りのフィルタ (I-1, I-2, I-3, I-4, I-5, I-6) について、指定した 4 つのパラメタ  $\mu, g_p, g_s, n$  の値をそれぞれ表 1 に示す (各例において、次数  $n$  以外の 3 つのパラメタ  $\mu, g_p, g_s$  を先に指定して、 $n$  の値は「方式 I」のフィルタが構成可能となる最小のものにしている)。そうして各例について、指定した 4 つのパラメタから決定された付録 A.1 の式 (A.2) の  $\sigma_k, \alpha_k, k = 1, 2$  の値をそれぞれ表 2 に示す。

「方式 I」の 6 通りのフィルタについて、正規化座標  $t$  を横軸にとり、伝達関数の大きさ  $|g(t)|$  の常用対数を縦軸にとってプロットしたグラフをそれぞれ図 A.1 から図 A.6 までに示す。

### 3.2 「方式 II」の伝達関数の例

伝達関数  $g(t)$  の形状パラメタ 3 つと次数の組  $(\mu, g_p, g_s, n)$  を指定して、それを満たす「方式 II」の伝達関数を決定した。対応するフィルタの作用は伝達関数から容易に導ける (付録 A.1.6)。

実験に用いた 6 通りのフィルタ (II-1, II-2, II-3, II-4, II-5, II-6) について、指定した 4 つのパラメタ  $\mu, g_p, g_s, n$  の値をそれぞれ表 5 に示す (各例において、次数  $n$  以外の 3 つのパラメタ  $\mu, g_p, g_s$  を先に指定して、 $n$  の値は「方式 II」のフィルタが構成可能となる最小のものにしている)。そうして各例について、指定した 4 つのパラメタから決定された付録 A.1 の式 (A.2) の  $\sigma_k, \alpha_k, k = 1, 2$  の値をそれぞれ表 6 に示す。

「方式 II」の 6 通りのフィルタについて、正規化座標  $t$  を横軸にとり、伝達関数の大きさ  $|g(t)|$  の常用対数を縦

軸にとってプロットしたグラフをそれぞれ図 A.13 から図 A.18 までに示す。

### 3.3 各フィルタによる計算結果

「方式 I」の 6 通りのフィルタのそれぞれについて、 $B$ -正規直交化と組み合わせた反復 3 回までについて、得られた近似固有対の固有値を横軸にとり、その相対残差の常用対数を縦軸にとってプロットしたグラフを図 A.7 から図 A.12 までに示す。

同様に、「方式 II」の 6 通りのフィルタのそれぞれについて、 $B$ -正規直交化と組み合わせた反復 3 回までについて、得られた近似固有対の固有値を横軸にとり、その相対残差の常用対数を縦軸にとってプロットしたグラフを図 A.19 から図 A.24 までに示す。

フィルタの反復回数に対する近似固有対の相対残差の最大値を「方式 I」のフィルタの各例について表 3 に、同様に「方式 II」のフィルタの各例について表 7 に示す。またフィルタの反復回数に対する対角化までの経過時間を「方式 I」のフィルタの各例について表 4 に、同様に「方式 II」のフィルタの各例について表 8 に示す。どの場合にもシフト行列 2 つの分解に要した経過時間は約 14 秒であった。反復あたりの経過時間は、フィルタの次数  $n$  が高いほど、ベクトルの数  $m$  が多いほど増える (ただしフィルタの適用によりベクトルの組の階数が低下すると  $B$ -正規直交化の際の閾値による切断を受けるので、一般にはベクトルの数  $m$  は反復ごとに減る可能性がある)。

今回扱った問題の規模は 1 ノードのシステムで計算できる最大のものではなく、示した経過時間についても並列化に対して最善の努力を傾けたものではない。またこれら各例の計算では、行列 2 つの分解や分解結果の 2 つを用いた前進後退代入の計算も、どちらも 2 つは独立なので並行に処理できるが、今回の実験では順次に行っている。

## 4. 実験 2: 実数シフトのレゾルベントを 1 つ用いるフィルタと 2 つ用いるフィルタの比較

まず比較の便利のために、シフトが実数のレゾルベントを 2 つ用いたフィルタの場合と同様に、シフトが実数のレゾルベントを 1 つ用いたフィルタについても 4 つのパラメタ  $\mu, g_p, g_s, n$  の値を指定してそれらを正しく満たすものとして構成できるようにするために、今回は実数シフトの単一のレゾルベントからなるフィルタの構成法として、従来のものを少し拡張したものを用いる (構成法は付録 A.2 に記述)。そうして、その単一のレゾルベントを使う場合である「単一」と、レゾルベントを 2 つ用いる場合の 2 通りの方式である「方式 I」と「方式 II」(構成法は付録 A.1 に記述) について若干の比較実験を行なった。

例題とする式 (1) の実対称定値一般固有値問題は、再び前節 3 と全く同じ有限要素法の問題から導かれたものであり、

表 1 実験 1: 「方式 I」の各例の伝達関数に指定した 4 つのパラメタ

フィルタ	$\mu$	$g_p$	$g_s$	$n$
I-1	2.0	1E-2	1E-09	25
I-2	2.0	1E-2	1E-10	35
I-3	2.0	1E-3	1E-12	25
I-4	2.0	1E-3	1E-13	32
I-5	2.0	1E-3	1E-14	40
I-6	1.5	1E-4	1E-11	30

表 5 実験 1: 「方式 II」の各例の伝達関数に指定した 4 つのパラメタ

フィルタ	$\mu$	$g_p$	$g_s$	$n$
II-1	2.0	1E-2	1E-13	30
II-2	2.0	1E-2	1E-14	35
II-3	2.0	1E-3	1E-13	21
II-4	1.5	1E-4	1E-12	24
II-5	1.5	1E-4	1E-13	28
II-6	1.25	1E-6	1E-13	29

表 2 実験 1: 「方式 I」の伝達関数の  $x(t)$  に用いた  $\sigma_k$  と  $\alpha_k$

フィルタ	$k$	$\sigma_k$	$\alpha_k$
I-1	1	4.09068 41137 85926 9	9.68147 36896 33707 0
	2	2.02528 07667 67491 7	2.37312 19592 31734 7
I-2	1	5.19655 07817 65392 2	15.25918 03018 57066
	2	3.21576 96254 85300 8	5.84346 85487 09282 1
I-3	1	2.22755 26153 98233 9	10.70208 67035 10560
	2	1.59850 75775 766164	5.51114 60688 83539 0
I-4	1	3.32580 23062 73146 3	8.98973 04258 55874 8
	2	1.79146 09244 00880 6	2.60836 57440 39891 1
I-5	1	3.99137 37417 64652 6	11.75250 98719 03449
	2	2.39289 28457 85695 5	4.22408 19519 01427 9
I-6	1	2.69117 50089 59303 0	8.93745 60356 09324 4
	2	1.71861 35211 28330 2	3.64490 72765 50080 1

表 6 実験 1: 「方式 II」の伝達関数の  $x(t)$  に用いた  $\sigma_k$  と  $\alpha_k$

フィルタ	$k$	$\sigma_k$	$\alpha_k$
II-1	1	1.67933 35315 46617 8	12.84712 18363 24346
	2	1.25898 93885 43740 0	8.12041 76097 42180 1
II-2	1	1.92356 13781 91710 9	14.18630 98321 53896
	2	1.45862 38171 49344 4	9.04662 44340 09778 8
II-3	1	1.22291 68196 12936 5	4.32668 10367 40262 2
	2	0.37200 77616 25172 68	0.81235 02518 27033 46
II-4	1	1.23356 16207 65095 2	3.93345 42009 89467 5
	2	0.41603 30166 83183 49	0.84103 96834 36731 41
II-5	1	0.96499 058641 91108 4	12.22386 05471 97841
	2	0.78605 22624 66379 16	9.05045 51521 88670 0
II-6	1	0.97498 17452 41140 78	4.55966 85101 81800 2
	2	0.51619 30340 47827 13	1.85327 70031 67030 3

表 3 実験 1: 「方式 I」のフィルタの反復回数と相対残差の最大値

フィルタ	反復 1 回	反復 2 回	反復 3 回
I-1	1.0E-04	5.5E-12	2.6E-12
I-2	1.2E-05	4.2E-12	4.3E-12
I-3	1.2E-06	3.8E-12	3.9E-12
I-4	9.5E-08	2.5E-12	2.6E-12
I-5	1.2E-08	3.2E-12	3.6E-12
I-6	8.7E-05	4.5E-12	3.8E-12

表 7 実験 1: 「方式 II」のフィルタの反復回数と相対残差の最大値

フィルタ	反復 1 回	反復 2 回	反復 3 回
II-1	4.8E-07	8.1E-12	9.6E-12
II-2	4.1E-08	9.8E-12	1.0E-11
II-3	1.3E-06	2.2E-12	2.7E-12
II-4	1.1E-04	3.6E-12	4.0E-12
II-5	1.7E-05	1.8E-11	1.9E-11
II-6	2.7E-03	1.8E-10	9.8E-12

表 4 実験 1: 「方式 I」: フィルタ反復回数と対角化に要した経過時間 (秒)

フィルタ	$n$	$m$	反復 1 回	反復 2 回	反復 3 回
I-1	25	200	100.2	180.5	256.7
I-2	35	200	128.9	237.2	343.8
I-3	25	200	99.8	180.6	257.3
I-4	32	200	121.1	222.1	318.9
I-5	40	200	143.7	248.8	346.5
I-6	30	125	91.5	164.3	236.5

表 8 実験 1: 「方式 II」: フィルタ反復回数と対角化に要した経過時間 (秒)

フィルタ	$n$	$m$	反復 1 回	反復 2 回	反復 3 回
II-1	30	200	116.0	210.2	301.1
II-2	35	200	129.5	220.8	307.2
II-3	21	200	88.9	159.5	224.4
II-4	24	125	80.3	137.9	195.2
II-5	28	125	88.3	157.0	224.0
II-6	29	100	75.6	131.5	186.8

FEM の要素分割の方式も同じ  $(N1, N2, N3) = (40, 50, 60)$  として、固有値が区間  $[a, b] = [3, 30]$  に含まれる全部で 54 個の固有対を求めた。

近似対を求めるのに掛かった経過時間を直交化付きフィルタの適用の回数 3 回までについて、 $\mu$  の値が 2.0, 1.5, 1.25 に対してそれぞれ表 15, 表 16, 表 17 に一応示したが、レゾルベントを 2 つ用いる場合には、2 つのレゾルベントの構成の準備やベクトル  $m$  個の組に 2 つのレゾルベントを作用させる処理は、それぞれ 2 つを並行に計算でき

るが、今回の実験ではそのようにはせずに、2 つを順次に処理している。また OpenMP によるスレッド並列化についても最善の努力によるものではない。

#### 4.1 実数シフトの単一のレゾルベントで構成されたフィルタの実験

これは、実数シフトの単一のレゾルベントの Chebyshev 多項式 (付録 A.2 の式 (A.67)) をフィルタに用いた場合の実験の例である。

この「単一」のフィルタで、 $\mu$  の値を 2.0, 1.5, 1.25 と指定した各場合の計算結果をそれぞれ表 9, 表 10, 表 11 の上段に示す。どの場合についてもパラメタ  $g_s$  の値は  $1E-13$  と指定している。表中では、次数  $n$  は 10 から 40 まで 5 刻みで変え、直交化付きフィルタの反復回数を 1 から 3 までとして、得られた近似固有対の相対残差の最大値を各場合について示している。このように 3 つのパラメタ  $\mu$ ,  $n$ ,  $g_s$  については値を直接指定し、残りのパラメタ  $g_p$  の値は  $0.5^j$  の形 ( $j$  は正の整数) に制限して、そのうちでフィルタが実現可能である場合の最大値に設定した。表中には実際に用いた  $g_p$  の値を有効数字 3 桁に丸めたものを載せている。

そうして、 $\mu$  の値を 2.0, 1.5, 1.25 とした各場合について、フィルタに最初に適用するベクトルの数  $m$  をそれぞれ 200, 125, 100 と指定した。

$\mu$  の値を 2.0, 1.5, 1.25 とした各場合について、今回の「単一」のフィルタで、近似固有対を求めるために掛かった経過時間をそれぞれ表 15, 表 16, 表 17 の上段に示す。

$\mu$  の値を 2.0, 1.5, 1.25 とした各場合について、「単一」のフィルタの適用回数が 1 から 3 まで (IT1, IT2, IT3) について横軸にフィルタの次数  $n$  をとり、縦軸に相対残差の最大値の常用対数の値をとってそれぞれ赤, 緑, 青の線を用いたグラフをプロットしたものが図 A-25, 図 A-28, 図 A-31 である。

$\mu$  の値を 2.0, 1.5, 1.25 とした各場合について、横軸にフィルタの次数  $n$  をとり、縦軸に「単一」のフィルタを 1 回適用して得られた近似固有対の相対残差の最大値の常用対数をとってプロットしたグラフを図 A-34 に示す。

#### 4.2 実数シフトのレゾルベント 2 つで構成される「方式 I」と「方式 II」のフィルタの実験

これは実数シフトのレゾルベント 2 つの線形結合の Chebyshev 多項式をフィルタとした場合の実験結果の例である (付録 A.1 に構成を記述)。

「方式 I」のフィルタで  $\mu$  の値を 2.0, 1.5, 1.25 とした各場合が、それぞれ表 9, 表 10, 表 11 の中段であり、同様に「方式 II」のフィルタではそれぞれの下段である。簡単のためパラメタ  $g_s$  の値はどの場合にも  $1E-13$  と指定した。各表中では、次数  $n$  を 10 から 40 まで 5 刻みで変えて、直交化付きフィルタの反復回数が 1 から 3 までについて、得られた近似固有対の相対残差の最大値を示している。このように 3 つのパラメタ  $\mu$ ,  $n$ ,  $g_s$  については値を直接指定し、残り 1 つのパラメタ  $g_p$  はその値を  $0.5^j$  ( $j$  は正の整数) の形に制限して、そのなかからフィルタが実現可能である最大のものを選択した。この探索でフィルタが実現不能な場合は、表中で横線を引いて示している (表中に実際に採択された  $g_p$  の値を有効数字 3 桁に丸めたものを載せている)。

$\mu$  の値を 2.0, 1.5, 1.25 とした各場合について、フィル

タに適用する最初のベクトルの数  $m$  はそれぞれ 200, 125, 100 とした。

「方式 I」で  $\mu$  の値を 2.0, 1.5, 1.25 とした各場合の経過時間をそれぞれ表 15, 表 16, 表 17 の中段に示し、同様に「方式 II」ではそれぞれ表 15, 表 16, 表 17 の下段に示す。フィルタに適用するベクトルの数  $m$  が同じである場合には、「方式 I」も「方式 II」も行列の分解や連立 1 次方程式の組を解く計算の作業内容はまったく同じであるので、経過時間もほぼ同じになっている。ただし、フィルタ特性が異なることから、フィルタ操作と組み合わせた  $B$ -正規直交化で検出される、ベクトルの組の実効階数の低下の状況に相違が生じると、反復の 2 回目以降からは作業対象となるベクトルの数が同じではなくなり、それにより経過時間の違いを生む可能性がある。

$\mu$  の値を 2.0, 1.5, 1.25 とした各場合について、レゾルベントを 2 つ用いる「方式 I」のフィルタの適用回数が 1 から 3 まで (IT1, IT2, IT3) について、横軸にフィルタの次数  $n$  をとり、縦軸に相対残差の最大値の常用対数の値をとって、それぞれ赤, 緑, 青の線を用いたグラフをプロットしたものが図 A-26, 図 A-29, 図 A-32 である。 $\mu = 2.0$  のときは次数  $n$  が 15 以上の場合には、近似固有対の精度の改良はフィルタの反復 2 回目で終了している。そうして  $\mu = 1.5$  のときは次数  $n$  が 25 以上の場合には反復 2 回目で終了している。さらに  $\mu = 1.25$  のときも反復 3 回目で改良はほぼ終了している。

同様に、 $\mu$  の値を 2.0, 1.5, 1.25 とした各場合について、レゾルベントを 2 つ用いる「方式 II」のフィルタの適用回数が 1 から 3 まで (IT1, IT2, IT3) について、横軸にフィルタの次数  $n$  をとり、縦軸に相対残差の最大値の常用対数の値をとって、それぞれ赤, 緑, 青の線を用いたグラフをプロットしたものが図 A-27, 図 A-30, 図 A-33 である。これも  $\mu = 2.0$  のときは次数  $n$  が 15 以上の場合には近似固有対の精度の改良はフィルタの反復 2 回目で終了している。そうして  $\mu = 1.5$  のときは次数  $n$  が 20 以上の場合には反復 2 回目で終了している。さらに  $\mu = 1.25$  のときも反復 3 回目で改良はほぼ終了している。

$\mu$  の値を 2.0, 1.5, 1.25 とした各場合について、横軸にフィルタの次数  $n$  をとり、縦軸に「方式 I」のフィルタの適用 1 回で得られた近似固有対の相対残差の最大値の常用対数をとってプロットしたグラフを図 A-35 に示す。同様に、縦軸に「方式 II」のフィルタの適用 1 回で得られた近似固有対の相対残差の最大値の常用対数をとってプロットしたグラフを図 A-36 に示す。

#### 4.3 実験 2 の結果

いま  $g_s$  の値は  $1E-13$  に固定して、 $\mu$  の値を 2.0, 1.5, 1.25 とした各場合の  $g_p$  の値をフィルタの次数  $n$  とフィルタの種類 (「単一」, 「方式 I」, 「方式 II」) について集めたものを、

表 9 実験 2 : 最大相対残差,  $\mu=2.0$  ( $m=200$ ),  $g_s=1E-13$

	$n$	$g_p$	反復 1 回	反復 2 回	反復 3 回
単一	10	1.19E-07	5.7E-04	8.7E-11	4.4E-13
	15	3.81E-06	2.0E-05	5.6E-13	5.6E-13
	20	1.53E-05	7.4E-06	7.1E-13	7.0E-13
	25	3.05E-05	5.1E-06	8.2E-13	8.3E-13
	30	3.05E-05	3.7E-06	7.0E-13	6.8E-13
	35	6.10E-05	2.9E-06	9.4E-13	9.4E-13
	40	6.10E-05	2.9E-06	7.6E-13	7.6E-13
方式 I	10	2.38E-07	3.3E-04	8.9E-11	4.6E-13
	15	1.53E-05	6.0E-06	1.2E-12	1.2E-12
	20	1.22E-04	8.6E-07	1.9E-12	2.0E-12
	25	2.44E-04	3.9E-07	1.1E-12	1.1E-12
	30	4.88E-04	2.3E-07	1.3E-12	1.3E-12
	35	9.77E-04	1.1E-07	1.5E-12	1.6E-12
	40	1.95E-03	7.3E-08	3.5E-12	3.6E-12
方式 II	10	4.77E-07	4.4E-03	2.9E-10	4.7E-12
	15	6.10E-05	3.4E-05	3.1E-12	4.6E-12
	20	4.88E-04	6.5E-06	1.9E-12	2.0E-12
	25	3.91E-03	7.1E-07	4.3E-12	5.0E-12
	30	7.81E-03	2.5E-07	2.6E-12	2.8E-12
	35	1.56E-02	1.8E-07	3.8E-12	3.7E-12
	40	1.56E-02	1.4E-07	2.3E-12	2.3E-12

表 10 実験 2 : 最大相対残差,  $\mu=1.5$  ( $m=125$ ),  $g_s=1E-13$

	$n$	$g_p$	反復 1 回	反復 2 回	反復 3 回
単一	10	3.73E-09	4.5E-02	2.3E-07	3.7E-12
	15	5.96E-08	1.0E-03	9.2E-10	5.4E-13
	20	2.38E-07	3.2E-04	4.9E-11	6.3E-13
	25	4.77E-07	2.0E-04	1.7E-11	7.3E-13
	30	4.77E-07	2.1E-04	1.2E-11	6.2E-13
	35	9.54E-07	8.7E-05	4.6E-12	8.8E-13
	40	9.54E-07	8.6E-05	6.2E-12	7.5E-13
方式 I	10	7.45E-09	1.6E-02	6.7E-08	7.1E-13
	15	1.19E-07	6.5E-04	3.1E-10	6.2E-13
	20	9.54E-07	9.2E-05	3.7E-12	8.9E-13
	25	3.81E-06	1.7E-05	1.3E-12	1.3E-12
	30	7.63E-06	8.0E-06	1.4E-12	1.4E-12
	35	1.53E-05	4.6E-06	1.8E-12	1.9E-12
	40	1.53E-05	4.7E-06	1.5E-12	1.6E-12
方式 II	10	7.45E-09	8.0E-02	2.1E-07	4.5E-12
	15	4.77E-07	9.2E-03	5.2E-10	4.2E-12
	20	7.63E-06	1.2E-04	4.4E-12	5.2E-12
	25	3.05E-05	4.4E-05	3.2E-12	3.3E-12
	30	1.22E-04	1.3E-05	4.5E-12	4.6E-12
	35	2.44E-04	1.4E-05	3.9E-12	4.0E-12
	40	4.88E-04	8.2E-06	5.8E-12	4.9E-12

それぞれ表 12, 表 13, 表 14 に掲げる. 通過域におけるフィルタの伝達率の最大最小比は  $1/g_p$  であるから,  $g_p$  の値が大きくて 1 に近いほど, 得られる近似固有対の精度の一樣性が良くなるのが期待できる. これらの 3 つの表では,  $\mu$  の値とフィルタの次数  $n$  が揃っているときには, 「単一」, 「方式 I」, 「方式 II」 の後のものほど  $g_p$  の値は大きく

表 11 実験 2 : 最大相対残差,  $\mu=1.25$  ( $m=100$ ),  $g_s=1E-13$

	$n$	$g_p$	反復 1 回	反復 2 回	反復 3 回
単一	10	2.33E-10	7.4E-02	5.2E-05	1.0E-08
	15	1.86E-09	6.2E-02	1.2E-06	4.2E-11
	20	3.73E-09	1.8E-02	2.5E-07	3.3E-12
	25	7.45E-09	8.9E-03	6.7E-08	6.3E-13
	30	1.49E-08	3.2E-03	1.3E-08	8.0E-13
	35	1.49E-08	3.7E-03	1.2E-08	6.7E-13
	40	1.49E-08	2.9E-03	8.8E-09	6.0E-13
方式 I	10	---	---	---	---
	15	3.73E-09	2.7E-02	2.1E-07	3.7E-12
	20	1.49E-08	3.0E-03	1.1E-08	8.4E-13
	25	5.96E-08	9.9E-04	8.0E-10	1.5E-12
	30	1.19E-07	4.3E-04	1.6E-10	1.8E-12
	35	1.19E-07	4.0E-04	1.5E-10	1.4E-12
	40	2.38E-07	2.3E-04	3.3E-11	1.8E-12
方式 II	10	---	---	---	---
	15	7.45E-09	7.1E-01	2.2E-06	3.8E-11
	20	5.96E-08	5.3E-02	2.8E-08	4.7E-12
	25	2.38E-07	8.6E-03	1.9E-09	4.1E-12
	30	9.54E-07	1.4E-03	1.3E-10	5.4E-12
	35	1.91E-06	2.4E-03	3.1E-11	4.9E-12
	40	3.81E-06	5.8E-04	9.0E-12	5.3E-12

なっている. つまりフィルタの伝達特性を  $\mu$ ,  $g_p$ ,  $g_s$  の 3 つで特徴付ける場合には, 「単一」, 「方式 I」, 「方式 II」 の後のものほど良い特性であることがわかる.

実験 2 で, フィルタを 1 回適用して得られた近似固有対の相対残差の最大値を比較したグラフ図 A.34, 図 A.35, 図 A.36 からは, 単一のレゾルベントで構成したフィルタは, 2 つのレゾルベントで構成した「方式 I」や「方式 II」のフィルタに比べて, 多項式の次数  $n$  を増したときの相対残差の最大値の減少傾向が緩やかである. このことはフィルタの次数  $n$  を増したときに, 通過域における伝達率の最小値  $g_p$  として選んだ値が増加している傾向と概ね一致している ( $\mu = 2.0$  の場合は表 12,  $\mu = 1.5$  の場合は表 13,  $\mu = 1.25$  の場合は表 14). どのフィルタの場合についても,  $\mu$  の値を小さくするほど (伝達関数は遷移域で急峻な変化を強いられてその結果として通過域での伝達関数の最大最小比が大きくなるので), 近似固有対の残差の最大値は増加している. 通過域に隣接する遷移域に含まれる固有値を持つ不要な固有対がたくさん存在する場合には, フィルタの種類を固定したときに, パラメタ  $\mu$  を小さくすれば遷移域が狭くなるので, フィルタに適用する初期ベクトルの数  $m$  を少なくできて, フィルタの中で  $n$  回繰り返してレゾルベントを適用する際のベクトルの数  $m$  に比例する記憶参照や演算量を減らせる利点があるが, フィルタの特性が悪くなる (今の場合は  $g_p$  が小さくなる) ので, 計算で得られる近似固有対の精度は低下する可能性がある. レゾルベント 2 つで構成されたフィルタの方が 1 つで構成されたものに比べて  $\mu$  の値をより小さく指定できる. そうして

表 12 実験 2 : 次数  $n$  の各種フィルタで選ばれた  $g_p$  の値 ( $\mu=2.0$ )

$n$	「単一」	「方式 I」	「方式 II」
10	1.19E-07	2.38E-07	4.77E-07
15	3.81E-06	1.53E-05	6.10E-05
20	1.53E-05	1.22E-04	4.88E-04
25	3.05E-05	2.44E-04	3.91E-03
30	3.05E-05	4.88E-04	7.81E-03
35	6.10E-05	9.77E-04	1.56E-02
40	6.10E-05	1.95E-03	1.56E-02

表 13 実験 2 : 次数  $n$  の各種フィルタで選ばれた  $g_p$  の値 ( $\mu=1.5$ )

$n$	「単一」	「方式 I」	「方式 II」
10	3.73E-09	7.45E-09	7.45E-09
15	5.96E-08	1.19E-07	4.77E-07
20	2.38E-07	9.54E-07	7.63E-06
25	4.77E-07	3.81E-06	3.05E-05
30	4.77E-07	7.63E-06	1.22E-04
35	9.54E-07	1.53E-05	2.44E-04
40	9.54E-07	1.53E-05	4.88E-04

表 14 実験 2 : 次数  $n$  の各種フィルタで選ばれた  $g_p$  の値 ( $\mu=1.25$ )

$n$	「単一」	「方式 I」	「方式 II」
10	2.33E-10	---	---
15	1.86E-09	3.73E-09	7.45E-09
20	3.73E-09	1.49E-08	5.96E-08
25	7.45E-09	5.96E-08	2.38E-07
30	1.49E-08	1.19E-07	9.54E-07
35	1.49E-08	1.19E-07	1.91E-06
40	1.49E-08	2.38E-07	3.81E-06

「方式 2」は「方式 1」に比べて他の 3 つのパラメタが同じ場合には  $\mu$  の値をより小さく指定できる。

実験 2 で、フィルタを 1 回だけ適用した場合の近似固有対の相対残差の最大値を比較したグラフ (図 A-37, 図 A-38, 図 A-39) からは、単一のレゾルベントで構成されたフィルタを用いた場合 (赤線) よりも、2 つのレゾルベントから構成された「方式 I」 (緑線) と「方式 II」 (青線) のフィルタを用いた場合の方が近似固有対の精度が少し高いことがわかる。そうして、「方式 I」と「方式 II」を比べると、「方式 I」の方が精度が高い。これは「方式 II」では区間 [3, 30] の下端に近い固有値 3.00102667 を持つ固有対は伝達率が小さいので精度が悪くなり、相対残差の最大値で比較をする際に不利になっている。

#### 4.3.1 フィルタ反復との併用の考察

$\mu = 2.0$  の場合 ( $m = 200$ )

- フィルタが「単一」の場合には (表 9 と表 15 のそれぞれ上段から) :  
たとえば次数  $n = 20$  のフィルタを 1 回適用した場合の相対残差の最大値  $7.4E-06$  よりも、次数  $n = 10$  のフィルタを 2 回適用した場合の相対残差の最大値

表 15 実験 2 : 経過時間 (秒),  $\mu=2.0$  ( $m=200$ )

$n$	反復 1 回	反復 2 回	反復 3 回
10	37.3	65.2	87.0
15	47.4	82.9	114.8
20	56.1	100.4	140.4
25	64.8	118.1	166.1
30	73.8	134.9	191.1
35	82.8	153.6	217.9
40	90.9	170.9	244.4
10	56.9	95.3	129.3
15	70.8	124.6	172.1
20	85.3	153.7	214.7
25	101.5	183.3	259.8
30	114.9	211.4	301.9
35	129.4	239.1	344.1
40	144.3	266.9	389.2
10	57.0	96.0	129.6
15	72.4	125.2	171.6
20	87.0	152.8	215.2
25	100.5	183.0	258.2
30	115.0	210.8	301.9
35	129.5	238.5	344.0
40	144.3	268.3	387.6

表 16 実験 2 : 経過時間 (秒),  $\mu=1.5$  ( $m=125$ )

$n$	反復 1 回	反復 2 回	反復 3 回
10	29.1	48.0	64.2
15	36.1	61.6	84.8
20	43.5	75.1	105.4
25	50.4	89.1	126.2
30	57.4	103.4	147.5
35	65.1	117.1	168.3
40	70.8	130.5	189.8
10	46.3	73.1	99.8
15	57.4	95.7	134.6
20	68.6	119.7	167.8
25	80.7	141.7	202.7
30	92.6	166.1	236.8
35	103.5	188.0	271.6
40	114.6	211.9	306.3
10	45.7	72.5	98.0
15	59.0	97.0	134.8
20	69.4	120.5	168.6
25	81.5	142.6	202.5
30	92.4	165.5	237.4
35	104.7	188.6	271.2
40	115.6	211.7	305.5

$8.7E-11$  の方が精度が 5 桁程度高いことがわかる。ただし経過時間は前者が 56.1 秒に対して後者は (途中で  $B$ -正規直交化が 1 回追加されるので) 65.2 秒で、少し長くなっている。

同様に、次数  $n = 30$  のフィルタを 1 回適用した場合の相対残差の最大値  $3.7E-06$  よりも、次数  $n = 15$

表 17 実験 2: 経過時間 (秒),  $\mu=1.25$  ( $m=100$ )

	$n$	反復 1 回	反復 2 回	反復 3 回
単一	10	24.4	38.9	51.6
	15	30.5	50.1	68.3
	20	36.5	60.9	85.3
	25	41.7	73.0	101.8
	30	47.3	84.2	118.4
	35	52.9	94.3	135.0
	40	58.2	106.0	151.7
	方式 I	10	---	---
15		49.0	79.4	108.6
20		59.1	98.1	137.1
25		66.6	116.8	163.8
30		76.7	133.7	192.1
35		85.3	153.2	220.1
40		95.1	171.3	245.8
方式 II		10	---	---
	15	49.4	79.8	110.3
	20	58.6	98.8	138.0
	25	68.1	116.6	165.6
	30	77.5	134.9	191.9
	35	86.4	153.8	222.0
	40	96.5	173.1	246.8

のフィルタを 2 回適用した場合の相対残差の最大値  $5.6E-13$  の方が精度が 7 桁程度高い。そうして経過時間は前者が 73.8 秒に対して後者は 82.9 秒である。

- フィルタが「方式 I」の場合には (表 9 と表 15 のそれぞれ中段から) :

たとえば  $n = 20$  のフィルタを 1 回適用した場合の相対残差の最大値  $8.6E-07$  よりも,  $n = 10$  のフィルタを 2 回適用した場合の相対残差の最大値  $8.9E-11$  の方が精度が 4 桁高い。そうして経過時間は前者が 85.3 秒に対して後者は 95.3 秒である。

同様に,  $n = 30$  のフィルタを 1 回適用した場合の相対残差の最大値  $2.3E-07$  よりも,  $n = 15$  のフィルタを 2 回適用した場合の相対残差の最大値  $1.2E-12$  の方が精度が 5 桁程度高い。そうして経過時間は前者が 114.9 秒に対して後者は 124.6 秒である。

- フィルタが「方式 II」の場合には (表 9 と表 15 のそれぞれ下段から) :

たとえば  $n = 20$  のフィルタを 1 回適用した場合の相対残差の最大値  $6.5E-06$  よりも,  $n = 10$  のフィルタを 2 回適用した場合の相対残差の最大値  $2.9E-10$  の方が精度が 4 桁程度高い。そうして経過時間は前者が 87.0 秒に対して後者は 96.0 秒である。

同様に  $n = 30$  のフィルタを 1 回適用した場合の相対残差の最大値  $2.5E-07$  よりも,  $n = 15$  のフィルタを 2 回適用した場合の相対残差の最大値  $3.1E-12$  の方が精度が 5 桁程度高い。そうして経過時間は前者が 115.0 秒に対して後者は 125.2 秒である。

$\mu = 1.5$  の場合 ( $m = 125$ )

- フィルタが「単一」のとき (表 10 と表 16 のそれぞれ上段から) :

たとえば  $n = 20$  のフィルタを 1 回適用した場合の相対残差の最大値  $3.2E-04$  よりも,  $n = 10$  のフィルタを 2 回適用した場合の相対残差の最大値  $2.3E-07$  の方が精度が 3 桁程度高い。そうして経過時間は前者が 43.5 秒に対して後者は 48.0 秒である。

同様に  $n = 30$  のフィルタを 1 回適用した場合の相対残差の最大値  $2.1E-04$  よりも,  $n = 15$  のフィルタを 2 回適用した場合の相対残差の最大値  $9.2E-10$  の方が精度が 5 桁程度高い。そうして経過時間は前者が 57.4 秒に対して後者は 61.6 秒である。

- フィルタが「方式 I」のとき (表 10 と表 16 のそれぞれ中段から) :

たとえば  $n = 20$  のフィルタを 1 回適用した場合の相対残差の最大値  $9.2E-05$  よりも,  $n = 10$  のフィルタを 2 回適用した場合の相対残差の最大値  $6.7E-08$  の方が精度が 3 桁程度高い。そうして経過時間は前者が 68.6 秒に対して後者は 73.1 秒である。

同様に  $n = 30$  のフィルタを 1 回適用した場合の相対残差の最大値  $8.0E-06$  よりも,  $n = 15$  のフィルタを 2 回適用した場合の相対残差の最大値  $3.1E-10$  の方が精度が 4 桁程度高い。そうして経過時間は前者が 92.6 秒に対して後者は 95.7 秒である。

- フィルタが「方式 II」のとき (表 10 と表 16 のそれぞれ下段から) :

たとえば  $n = 20$  のフィルタを 1 回適用した場合の相対残差の最大値  $1.2E-04$  よりも,  $n = 10$  のフィルタを 2 回反復した場合の相対残差の最大値  $2.1E-07$  の方が精度が 3 桁程度高い。そうして経過時間は前者が 69.4 秒に対して後者は 72.5 秒である。

同様に  $n = 30$  のフィルタを 1 回適用した場合の相対残差の最大値  $1.3E-05$  よりも,  $n = 15$  のフィルタを 2 回適用した場合の相対残差の最大値  $5.2E-10$  の方が精度が 4 桁程度高い。そうして経過時間は前者が 92.4 秒に対して後者は 97.0 秒である。

$\mu = 1.25$  の場合 ( $m = 100$ )

- フィルタが「単一」のとき (表 11 と表 17 の上段から) :

たとえば  $n = 20$  のフィルタを 1 回適用した場合の相対残差の最大値  $1.8E-02$  よりも,  $n = 10$  のフィルタを 2 回適用した場合の相対残差の最大値  $5.2E-05$  の方が精度が 2 桁程度高い。そうして経過時間は前者が 36.5 秒に対して後者は 38.9 秒である。

同様に  $n = 30$  のフィルタを 1 回適用した場合の相対残差の最大値  $3.2E-03$  よりも,  $n = 15$  のフィルタを 2 回適用した場合の相対残差の最大値  $1.2E-06$  の方が精度が 3 桁程度高い。そうして経過時間は前者が 47.3

秒に対して後者は 50.1 秒である。

- フィルタが「方式 I」のとき（表 11 と表 17 のそれぞれ中段から）：

たとえば  $n = 30$  のフィルタを 1 回適用した場合の相対残差の最大値  $4.3E-04$  よりも、 $n = 15$  のフィルタを 2 回適用した場合の相対残差の最大値  $2.1E-07$  の方が精度が 3 桁程度高い。そうして経過時間は前者が 76.7 秒に対して後者は 79.4 秒である。

- フィルタが「方式 II」のとき（表 11 と表 17 のそれぞれ下段から）：

たとえば  $n = 30$  のフィルタを 1 回適用した場合の相対残差の最大値  $1.4E-03$  よりも、 $n = 15$  のフィルタを 2 回適用した場合の相対残差の最大値  $2.2E-06$  の方が精度が 3 桁程度高い。そうして経過時間は前者が 77.5 秒に対して後者は 79.8 秒である。

これらの結果から、ランダムなベクトルの組から始める場合に、多項式の次数  $n$  の高いフィルタを 1 回適用するよりも、次数が  $n$  の半分のフィルタを途中で  $B$ -正規直交化をはさんで 2 回適用する方が、得られる近似固有対の精度がずっと良いので有利であることがわかる。ただし、 $B$ -正規直交化の計算時間（いまの例の場合にはあまり大した割合ではない）が追加になる。分散並列計算の場合には  $B$ -正規直交化を行うところでの同期待ちも必要になる。

## 5. おわりに

式 (1) の実対称定値一般固有値問題の固有対で固有値が固有値分布の下端付近のものを近似して求めるためのフィルタを、シフトが実数のレゾルベント 2 つの線形結合の Chebyshev 多項式として構成する方法を示し、それを簡単な例題に適用してある程度うまく働くことを確認した。

フィルタの伝達関数の形状の設計方式として本報告で提案している「方式 I」はいわゆる Butterworth 特性のものであり、固有値が小さい固有ベクトルほどフィルタを良く通過するので、固有値が下側のものほど得られる近似固有対の精度が良くなる傾向を持つ。この性質は固有値が固有値分布の下端付近にある固有対だけが必要で、しかも固有値が下側のものほど高い近似精度が必要とされる用途に対しては適しているといえる。今回の「方式 II」のフィルタは「方式 I」に比べて通過域全体での伝達率の一様性の向上（最大最小比  $1/g_p$  の低減）を狙ったものであり、フィルタを 4 つのパラメタ  $\mu$ ,  $g_p$ ,  $g_s$ ,  $n$  の値の組で指定する場合に、フィルタを実現可能なパラメタの選択範囲が「方式 I」に比べて広がる。たとえば 4 つのうちで  $\mu$  以外の 3 つのパラメタの値が共通である場合には、「方式 I」に比べて「方式 II」の方が  $\mu$  の値をより小さくできる。しかし「方式 II」のフィルタの伝達率は通過域がなるべく広くなるように、通過域における伝達率の最小値  $g_p$  を通過域の両端でとるように構成したものであるため、固有値が通過

域の両端に近い固有対であるほどそれだけ近似が悪くなる傾向を持つ。このような「方式 II」のフィルタの性質は、有限要素法などの変分的な離散化手法により導かれた固有値問題のように、固有値が最小固有値付近の固有対だけが重要であって、固有値が下側のものほど精度良く求める価値があるような場合には、あまり望ましいものではないかもしれない。

この研究は、大規模の実対称定値一般固有値問題を限られた計算資源で解くことを想定して、必要な記憶量と計算量をなるべく少なく抑えようとするものである。すべての固有対ではなくて、固有値が限られた範囲にある比較的少数の固有対だけを選んで解くが、固有値が下端付近の固有対だけが必要であることはしばしばあり、そのような場合には、比較的容易に固有値の存在範囲にはない実数をレゾルベントのシフトとして選ぶことができる。しかもシフトが実数であれば全体の計算は実数だけを用いて行える。そこでシフトとして実数だけを用いることにしてみる。

シフトを最小固有値未満の実数に選ぶと、レゾルベントに対応する連立 1 次方程式は係数行列が実対称正定値になり、その解法として直接法を用いることにすると、対称性を生かしてピボット選択なしで行列分解を行って数値的にも安定に解くことができ、さらに係数が帯行列の場合には行列分解は帯幅を拡げずに計算できる。

レゾルベントをたくさん（たとえば 10 個程度）使うことは諦めて、数個程度のごく少数に限れば、問題が大規模な場合に行列分解の計算を主記憶上で行い、得られた分解結果を主記憶に置いて前進後退代入を行うことが容易になる。その反面として、実数シフトのレゾルベントを少数（1 つあるいは 2 つ）用いて構成できるフィルタの特性はあまり良くなく、得られる近似固有対の精度も良くない。

用途によっては近似固有対の精度は数桁程度で構わない場合もありうる。しかし要求される精度がかなり高い場合（たとえば相対残差が  $1E-10$  以下）であると、「単一」、「方式 I」、「方式 II」はフィルタ 1 回だけの適用では得られる近似固有対の精度は十分でない。そのような場合には、直交化付きフィルタを（たとえば全部で 2~3 回）反復することで近似固有対の精度を改良できる。これは直交化とフィルタの処理を組み合わせることで、通過域におけるフィルタの伝達率の一様性があまり良くなくても、あるいは阻止域において伝達関数の大きさがそれほど微小でなくても、近似固有対の精度を反復により改良できるので、とても強力な方法である。

また今後は、用いる実数シフトのレゾルベントの数をもう少し増やして 3 つあるいは 4 つにした場合についてもフィルタを構成して検討してみたい。

## 付 録

### A.1 実数シフトのレゾルベント 2 つによる簡易型フィルタの設計法

フィルタをシフトが相異なる実数であるレゾルベント 2 つから作られる多項式とすると、その伝達関数はその 2 つの実数だけを (多重の) 極とする有理関数である。以下では、その関数形を強く制限した簡易型の設計法だけについて考察する。

これまで Chebyshev 多項式を用いた簡易型の設計法では、実数の極を 1 つだけを持つ伝達関数  $g(t)$  を制限して式 (A.1) の形にしてきた (極は  $n$  位で実数  $-\sigma$  だけである)。

$$g(t) \equiv g_s T_n(y(t)), y(t) \equiv 2x(t) - 1, x(t) \equiv \frac{\mu + \sigma}{t + \sigma}. \quad (\text{A.1})$$

この形の伝達関数は、( $n$  を増やせば) 阻止域における伝達率を容易に微小にできるが、通過域における伝達率の最大最小比を小さくすることは (遷移域の幅  $\mu - 1$  を大きくしなければ) できない。そこでこの簡易設計の手法を極が実数 1 つだけの場合から相異なる実数 2 つの場合へと拡張し、この拡張によって増えた自由度を用いて通過域における伝達率の最大最小比を小さくすることを試みる。

そのためには、式 (A.2) で表される実有理関数  $x(t)$  を新たに採用する。これは無限遠での値は零であり、極として相異なる負の実数 2 つだけを持つ。以下では、相異なる 2 つの極の符号を逆にした正の実数を  $\sigma_1$  と  $\sigma_2$  として、それらは常に  $\sigma_1 > \sigma_2 > 0$  を満たすと決める。

$$x(t) \equiv \frac{\alpha_1}{t + \sigma_1} - \frac{\alpha_2}{t + \sigma_2}. \quad (\text{A.2})$$

いま  $\mu > 1$  であるとする。式 (A.2) の関数  $x(t)$  は  $t \geq 0$  では連続である。そうしていま  $x_H$  と  $x_L$  は  $x_H > x_L > 1$  を満たす実数の未知数として、阻止域  $t \in [\mu, \infty)$  では  $x(t)$  は値が 1 以下で正、遷移域  $t \in (1, \mu)$  では  $x(t)$  は値が  $x_L$  未満で 1 より大きいとし、通過域  $t \in [0, 1]$  では  $x(t)$  は最大値が  $x_H$  で最小値は  $x_L$  であると仮定する。すると関数  $x(t)$  の連続性から  $x(\mu) = 1$  と  $x(1) = x_L$  が成り立つことが必要である。そうして伝達関数  $g(t)$  は Chebyshev 多項式を用いた簡易型であるとして、 $x(t)$  の多項式として従来と同様に式 (A.3) を採用する。

$$g(t) \equiv g_s T_n(y(t)), y(t) \equiv 2x(t) - 1. \quad (\text{A.3})$$

すると阻止域において  $x(t)$  の値は 1 以下で正と仮定したから、阻止域において式 (A.3) の  $g(t)$  の大きさ  $|g(t)|$  は  $g_s$  以下になる。さらに通過域  $t \in [0, 1]$  における  $g(t)$  の最大値と最小値がそれぞれ 1 と  $g_p$  であると仮定すると、 $x_H > x_L > 1$  より  $2x_H - 1 > 2x_L - 1 > 1$  であること、

Chebyshev 多項式は引数が 1 以上で単調増加性を持つこと、通過域において  $x(t)$  の最大値は  $x_H$  で最小値は  $x_L$  と仮定したこと、を合わせると式 (A.4) の 2 つの関係が得られる。

$$\begin{cases} 1 &= g_s T_n(2x_H - 1), \\ g_p &= g_s T_n(2x_L - 1). \end{cases} \quad (\text{A.4})$$

さらに  $z$  が実数の場合の恒等式 (A.5) を用いると、式 (A.4) を双曲線関数を用いて表した式 (A.6) が得られる。

$$\cosh^{-1}(2z^2 - 1) = 2 \cosh^{-1} |z| \quad (\text{A.5})$$

$$\begin{cases} x_H = \cosh^2 \left( \frac{1}{2n} \cosh^{-1} \frac{1}{g_s} \right), \\ x_L = \cosh^2 \left( \frac{1}{2n} \cosh^{-1} \frac{g_p}{g_s} \right). \end{cases} \quad (\text{A.6})$$

すると 3 つのパラメタ  $g_p, g_s, n$  の組が指定された場合に、式 (A.6) の中の各式の右辺をそれぞれ計算すると、 $x(t)$  の通過域  $t \in [0, 1]$  における最大値  $x_H$  と最小値  $x_L$  が決まる。

4 つのパラメタ  $\mu, g_p, g_s, n$  が指定されたときに、式 (A.2) の  $x(t)$  に含まれている 4 つの実数の値である  $\sigma_1, \sigma_2, \alpha_1, \alpha_2$  の組 (ただし  $\sigma_1 > \sigma_2 > 0$ ) がうまくとれることが伝達関数  $g(t)$  およびフィルタの構成には必要であり、そうでなければ構成は不可能である。

#### A.1.1 「方式 I」: 通過域の左端で最大かつ停留になる伝達関数

いま通過域  $[0, 1]$  における伝達率の最大最小比を小さくすることを狙って、通過域の左端である原点  $t = 0$  において伝達関数が最大値 1 をとり、しかもそこで値が停留する (微分の値が零になる) という条件を課して構成してみる。(伝達関数が原点  $t = 0$  において最大でかつ平坦な特性を持てば、固有値が指定区間の下側にある固有ベクトルであるほど良く通過させるフィルタになるので、得られる固有対も固有値が下側のものであるほど精度が良いものになると期待できる。そのような性質は、必要な固有対の固有値は下端側から順に少数だけの場合には、通常最も望ましいものである。) )

いま 4 つのパラメタの値の組  $(\mu, g_p, g_s, n)$  を指定するとき、その組を持つ伝達関数が実現可能な場合には、以下に示す手順に従って式 (A.2) の  $x(t)$  に含まれる 4 つの値  $\sigma_1, \sigma_2, \alpha_1, \alpha_2$  を決定できる。

まず式 (A.6) により、 $x_L$  と  $x_H$  の値を求めておく。そうして  $t = \mu$  における値の条件  $x(\mu) = 1$ 、 $t = 1$  における値の条件  $x(1) = x_L$ 、さらに通過域の左端  $t = 0$  における値の条件  $x(0) = x_H$  とそこで停留になるという条件  $\frac{d}{dt} x(t) \Big|_{t=0} = 0$  の全部で 4 つを順番に並べて数式で表すと、式 (A.7) になる。

$$\begin{cases} 1 &= \frac{\alpha_1}{\mu + \sigma_1} - \frac{\alpha_2}{\mu + \sigma_2}, \\ x_L &= \frac{\alpha_1}{1 + \sigma_1} - \frac{\alpha_2}{1 + \sigma_2}, \\ x_H &= \frac{\alpha_1}{\sigma_1} - \frac{\alpha_2}{\sigma_2}, \\ 0 &= -\frac{\alpha_1}{\sigma_1^2} + \frac{\alpha_2}{\sigma_2^2}. \end{cases} \quad (\text{A.7})$$

この連立方程式 (A.7) を既知である 3 つの値  $x_L$ ,  $x_H$ ,  $\mu$  を用いて解いて 4 つの未知数  $\sigma_1$ ,  $\alpha_1$ ,  $\sigma_2$ ,  $\alpha_2$  を実数の範囲で求める. そのような解が存在すれば, 指定されたパラメタの 4 つ組  $(\mu, g_p, g_s, n)$  に対応する「方式 I」の伝達関数  $g(t)$  は実現が可能であり, そうでなければ実現は不可能である.

まず連立式 (A.7) の第 4 番目の式から式 (A.8) の最初の等式が得られるが, その値は  $t$  によらない定数なのでそれをここでは  $C$  とおく.

$$\frac{\alpha_1}{\sigma_1^2} = \frac{\alpha_2}{\sigma_2^2} = C. \quad (\text{A.8})$$

すると, 2 つの極の係数はそれぞれ式 (A.9) で表せるので, このことを用いて  $\alpha_1$  と  $\alpha_2$  を消去できる.

$$\begin{cases} \alpha_1 = C\sigma_1^2, \\ \alpha_2 = C\sigma_2^2. \end{cases} \quad (\text{A.9})$$

連立式 (A.7) の第 1 番目の式と式 (A.9) から, 定数  $C$  の値の逆数は式 (A.10) で表せる.

$$\frac{1}{C} = \frac{\sigma_1^2}{\mu + \sigma_1} - \frac{\sigma_2^2}{\mu + \sigma_2} = (\sigma_1 - \sigma_2) \times \frac{\mu(\sigma_1 + \sigma_2) + \sigma_1\sigma_2}{(\mu + \sigma_1)(\mu + \sigma_2)}. \quad (\text{A.10})$$

連立式 (A.7) の第 3 番目の式と式 (A.9) と式 (A.10) をあわせると, 式 (A.11) が得られる.

$$x_H = C \times (\sigma_1 - \sigma_2) = \frac{(\mu + \sigma_1)(\mu + \sigma_2)}{\mu(\sigma_1 + \sigma_2) + \sigma_1\sigma_2}. \quad (\text{A.11})$$

この式 (A.11) から (構成が可能な場合には)  $\sigma_1 > \sigma_2$  であり,  $x_H$  が正であることから  $C > 0$  であること, そのことと式 (A.9) により  $\alpha_1 > \alpha_2 > 0$  であることもわかる.

連立式 (A.7) の第 2 番目の式と式 (A.9) と式 (A.10) と (A.11) をあわせると, 式 (A.12) が得られる.

$$\begin{aligned} x_L &= C \times \left( \frac{\sigma_1^2}{1 + \sigma_1} - \frac{\sigma_2^2}{1 + \sigma_2} \right) \\ &= \frac{(\mu + \sigma_1)(\mu + \sigma_2)}{\mu(\sigma_1 + \sigma_2) + \sigma_1\sigma_2} \times \frac{(\sigma_1 + \sigma_2) + \sigma_1\sigma_2}{(1 + \sigma_1)(1 + \sigma_2)} \\ &= x_H \times \frac{(\sigma_1 + \sigma_2) + \sigma_1\sigma_2}{(1 + \sigma_1)(1 + \sigma_2)}. \end{aligned} \quad (\text{A.12})$$

式 (A.2) に式 (A.9) と式 (A.11) をあわせると, 式 (A.13) が得られる.

$$\begin{aligned} x(t) &= \frac{(\mu + \sigma_1)(\mu + \sigma_2)}{\mu(\sigma_1 + \sigma_2) + \sigma_1\sigma_2} \times \frac{t(\sigma_1 + \sigma_2) + \sigma_1\sigma_2}{(t + \sigma_1)(t + \sigma_2)} \\ &= x_H \times \frac{t(\sigma_1 + \sigma_2) + \sigma_1\sigma_2}{(t + \sigma_1)(t + \sigma_2)}. \end{aligned} \quad (\text{A.13})$$

すると与えられた 3 つの値  $\mu$ ,  $x_H$ ,  $x_L$  の組から以下の手順で  $\sigma_1$  と  $\sigma_2$  の値が求められる. まず (A.11) と (A.12) から以下の関係 (A.14) が得られる.

$$\begin{cases} \frac{1}{x_H} &= \frac{(\mu + \sigma_1)(\mu + \sigma_2) - \mu^2}{(\mu + \sigma_1)(\mu + \sigma_2)} = 1 - \frac{\mu^2}{(\mu + \sigma_1)(\mu + \sigma_2)}, \\ \frac{x_L}{x_H} &= \frac{(1 + \sigma_1)(1 + \sigma_2) - 1}{(1 + \sigma_1)(1 + \sigma_2)} = 1 - \frac{1}{(1 + \sigma_1)(1 + \sigma_2)}. \end{cases} \quad (\text{A.14})$$

この式 (A.14) を式 (A.15) の置き換えを用いて書き直すと式 (A.16) が得られる.

$$\begin{cases} p \equiv \frac{\mu^2 x_H}{x_H - 1}, \\ q \equiv \frac{x_H}{x_H - x_L}. \end{cases} \quad (\text{A.15})$$

$$\begin{cases} (\mu + \sigma_1)(\mu + \sigma_2) = p, \\ (1 + \sigma_1)(1 + \sigma_2) = q. \end{cases} \quad (\text{A.16})$$

さらに式 (A.16) を少し変形して, 式 (A.17) を得る.

$$\begin{cases} \sigma_1\sigma_2 + \mu(\sigma_1 + \sigma_2) = p - \mu^2, \\ \sigma_1\sigma_2 + (\sigma_1 + \sigma_2) = q - 1. \end{cases} \quad (\text{A.17})$$

式 (A.17) は  $\sigma_1$  と  $\sigma_2$  の基本対称式である  $S_1 \equiv \sigma_1 + \sigma_2$  と  $S_2 \equiv \sigma_1\sigma_2$  についての連立 1 次方程式であり, それを解いて式 (A.18) が得られる.

$$\begin{cases} S_1 = \frac{p - q}{\mu - 1} - (\mu + 1), \\ S_2 = \mu + \frac{\mu q - p}{\mu - 1}. \end{cases} \quad (\text{A.18})$$

すると 2 次方程式 (A.19) の 2 根が相異なる正の実数である場合に限り, それらは  $\sigma_1$  と  $\sigma_2$  ( $\sigma_1 > \sigma_2 > 0$ ) である.

$$w^2 - S_1 w + S_2 = 0. \quad (\text{A.19})$$

この 2 次方程式 (A.19) の 2 根が相異なる正の実数であるための必要十分条件は,  $S_1 > 0$  かつ  $S_2 > 0$  かつ  $D \equiv S_1^2 - 4S_2 > 0$  である.

こうして相異なる正の実数  $\sigma_1$  と  $\sigma_2$  が求まれば, 式 (A.11) と式 (A.9) から得られる式 (A.20) を順に計算することで, 2 つの極の係数  $\alpha_1$  と  $\alpha_2$  がそれぞれ求まる.

$$\begin{cases} C \leftarrow \frac{x_H}{\sigma_1 - \sigma_2}, \\ \alpha_1 \leftarrow C\sigma_1^2, \\ \alpha_2 \leftarrow C\sigma_2^2. \end{cases} \quad (\text{A.20})$$

以上が式 (A.2) の  $x(t)$  を決定する手順である.

#### A.1.1.1 「方式 I」の設計法のまとめ

4 つのパラメタ  $(\mu, g_p, g_s, n)$  が与えられたとき, 式 (A.21) を順に計算する (ここで  $x'_H = x_H - 1$  の意味である).

$$\left\{ \begin{array}{l} x_H \leftarrow \cosh^2 \left( \frac{1}{2n} \cosh^{-1} \frac{1}{g_s} \right), \\ x'_H \leftarrow \sinh^2 \left( \frac{1}{2n} \cosh^{-1} \frac{1}{g_s} \right), \\ x_L \leftarrow \cosh^2 \left( \frac{1}{2n} \cosh^{-1} \frac{g_p}{g_s} \right), \\ p \leftarrow \frac{x_H}{x'_H} \times \mu^2, \\ q \leftarrow \frac{x_H}{x_H - x_L}, \\ S_1 \leftarrow \frac{p - q}{\mu - 1} - (\mu + 1), \\ S_2 \leftarrow \mu + \frac{\mu q - p}{\mu - 1}, \\ D \leftarrow S_1^2 - 4S_2. \end{array} \right. \quad (\text{A.21})$$

そうして、3つの条件  $S_1 > 0$ ,  $S_2 > 0$ ,  $D > 0$  をすべて満たしていれば構成は可能であり、そうでなければ構成は不可能である。そうして、構成が可能な場合には、式 (A.22) を順に計算する (なお  $\sigma_1 - \sigma_2 = \sqrt{D}$  であることを用いている)。

$$\left\{ \begin{array}{l} \sigma_1 \leftarrow \frac{1}{2} (S_1 + \sqrt{D}), \\ \sigma_2 \leftarrow \frac{S_2}{\sigma_1}, \\ C \leftarrow \frac{x_H}{\sqrt{D}}, \\ \alpha_1 \leftarrow C\sigma_1^2, \\ \alpha_2 \leftarrow C\sigma_2^2. \end{array} \right. \quad (\text{A.22})$$

こうして構成が可能である場合には  $(\sigma_1, \alpha_1, \sigma_2, \alpha_2)$  が求まり、 $\alpha_1 > \alpha_2 > 0$  かつ  $\sigma_1 > \sigma_2 > 0$  となる

### A.1.2 「方式 II」: 通過域の両端で値の等しい伝達関数

「方式 II」では通過域に於ける伝達関数の最大最小比を小さくすることを容易にする手段として、通過域の両端で伝達関数の値が等しいという条件  $g(0) = g(1) = g_p$  を課すことにする。そうして伝達関数は通過域での最大値 1 を内部のある 1 点  $t = t_p$  ( $0 < t_p < 1$ ) においてとるとする。

簡易設計による伝達関数の関数形は「方式 I」の場合と同じく (A.2) と (A.3) で与えられるとする ( $\sigma_1 > \sigma_2 > 0$  も仮定する)。そうしてパラメタの値の 4 つ組  $(\mu, g_p, g_s, n)$  を指定して、その組を持つ伝達関数の実現が可能であれば、式 (A.2) の  $x(t)$  が含む 4 つの実数値  $\sigma_1, \sigma_2, \alpha_1, \alpha_2$  を以下の手順で具体的に決定できる。

まず前節 A.1.1 と同様に、通過域に於ける  $x(t)$  の最小値  $x_L$  と最大値  $x_H$  の値はそれぞれ、3 つのパラメタ  $g_p, g_s, n$  の値から式 (A.6) を計算して求める。

すると  $x(t)$  の満たすべき条件は  $x(0) = x_L$  と  $x(1) = x_L$  と  $x(\mu) = 1$ , それと最大点における条件  $x(t_p) = x_H$  と  $\frac{d}{dt}x(t)|_{t=t_p} = 0$  の全部で 5 つであり、それらの条件を表す式を順に並べると式 (A.23) が得られる。

$$\left\{ \begin{array}{l} x_L = \frac{\alpha_1}{\sigma_1} - \frac{\alpha_2}{\sigma_2}, \\ x_L = \frac{\alpha_1}{1 + \sigma_1} - \frac{\alpha_2}{1 + \sigma_2}, \\ 1 = \frac{\alpha_1}{\mu + \sigma_1} - \frac{\alpha_2}{\mu + \sigma_2}, \\ x_H = \frac{\alpha_1}{t_p + \sigma_1} - \frac{\alpha_2}{t_p + \sigma_2}, \\ 0 = -\frac{\alpha_1}{(t_p + \sigma_1)^2} + \frac{\alpha_2}{(t_p + \sigma_2)^2}. \end{array} \right. \quad (\text{A.23})$$

この連立方程式 (A.23) を 3 つの値  $x_L, x_H, \mu$  を与えて解いて、5 つの未知数  $\sigma_1, \alpha_1, \sigma_2, \alpha_2, t_p$  を実数の範囲で求める。それが可能である場合には、指定されたパラメタの 4 つ組  $(\mu, g_p, g_s, n)$  に対応する「方式 II」の伝達関数  $g(t)$  は実現可能であり、そうでなければ不可能である。

まず連立方程式 (A.23) の第 1 番目と第 2 番目の式をあわせると式 (A.24) の最初の等式が得られる。その等式の値は  $t$  にはよらない定数なので、それを  $C$  とおく。

$$\frac{\alpha_1}{\sigma_1(1 + \sigma_1)} = \frac{\alpha_2}{\sigma_2(1 + \sigma_2)} = C. \quad (\text{A.24})$$

すると極の係数はそれぞれ式 (A.25) で表される。

$$\left\{ \begin{array}{l} \alpha_1 = C\sigma_1(1 + \sigma_1), \\ \alpha_2 = C\sigma_2(1 + \sigma_2). \end{array} \right. \quad (\text{A.25})$$

連立式 (A.23) の第 3 番目の式から式 (A.25) を用いて  $\alpha_1$  と  $\alpha_2$  を消去すれば、式 (A.26) が得られる。

$$1 = \frac{\alpha_1}{\mu + \sigma_1} - \frac{\alpha_2}{\mu + \sigma_2} = C \left\{ \frac{\sigma_1(1 + \sigma_1)}{\mu + \sigma_1} - \frac{\sigma_2(1 + \sigma_2)}{\mu + \sigma_2} \right\}. \quad (\text{A.26})$$

すると  $C$  の逆数の値は式 (A.27) で表される。

$$\begin{aligned} \frac{1}{C} &= \frac{\sigma_1(1 + \sigma_1)}{\mu + \sigma_1} - \frac{\sigma_2(1 + \sigma_2)}{\mu + \sigma_2} \\ &= (\sigma_1 - \sigma_2) \times \frac{\mu(1 + \sigma_1 + \sigma_2) + \sigma_1\sigma_2}{(\mu + \sigma_1)(\mu + \sigma_2)}. \end{aligned} \quad (\text{A.27})$$

そうして  $x_L$  の値は式 (A.28) で表される。

$$\begin{aligned} x_L &= \frac{\alpha_1}{\sigma_1} - \frac{\alpha_2}{\sigma_2} \\ &= C \times (\sigma_1 - \sigma_2) \\ &= \frac{(\mu + \sigma_1)(\mu + \sigma_2)}{\mu(1 + \sigma_1 + \sigma_2) + \sigma_1\sigma_2}. \end{aligned} \quad (\text{A.28})$$

なお、この式 (A.28) から  $C > 0$  であることがわかる。そのことと式 (A.25) をあわせると、 $\alpha_1 > \alpha_2 > 0$  であることもわかる。

つぎに極大点の位置  $t_p$  についての条件は連立式 (A.23) の第 5 番目の式から導かれる式 (A.29) である。

$$0 = \frac{\alpha_1}{(t_p + \sigma_1)^2} - \frac{\alpha_2}{(t_p + \sigma_2)^2} = C \left\{ \frac{\sigma_1(1 + \sigma_1)}{(t_p + \sigma_1)^2} - \frac{\sigma_2(1 + \sigma_2)}{(t_p + \sigma_2)^2} \right\}. \quad (\text{A.29})$$

そうして  $t_p > 0$ ,  $\sigma_1 > \sigma_2 > 0$  であることを用いて式 (A.29) の平方根を開くと、式 (A.30) の最初の等式が得られる。そ

の等式の値は  $t$  によらない定数であるのでそれを  $\Gamma$  とおいた。

$$\frac{\sqrt{\sigma_1(1+\sigma_1)}}{t_p + \sigma_1} = \frac{\sqrt{\sigma_2(1+\sigma_2)}}{t_p + \sigma_2} = \Gamma. \quad (\text{A.30})$$

すると  $\sigma_1 \neq \sigma_2$  であるから、式 (A.31) が得られる。

$$\begin{aligned} \Gamma &= \frac{\sqrt{\sigma_1(1+\sigma_1)} - \sqrt{\sigma_2(1+\sigma_2)}}{\sigma_1 - \sigma_2} \\ &= \frac{1 + \sigma_1 + \sigma_2}{\sqrt{\sigma_1(1+\sigma_1)} + \sqrt{\sigma_2(1+\sigma_2)}}. \end{aligned} \quad (\text{A.31})$$

そうして、式 (A.30) から式 (A.32) が得られる。

$$\begin{cases} t_p + \sigma_1 = \frac{1}{\Gamma} \times \sqrt{\sigma_1(1+\sigma_1)}, \\ t_p + \sigma_2 = \frac{1}{\Gamma} \times \sqrt{\sigma_2(1+\sigma_2)}. \end{cases} \quad (\text{A.32})$$

式 (A.32) の中の 2 つの式を連立させると、 $t_p$  を表す式 (A.33) が得られる。

$$\begin{aligned} t_p &= \frac{1}{\Gamma} \times \frac{\sigma_1 \sqrt{\sigma_2(1+\sigma_2)} - \sigma_2 \sqrt{\sigma_1(1+\sigma_1)}}{\sigma_1 - \sigma_2} \\ &= \frac{1}{\Gamma} \times \frac{\sigma_1 \sigma_2}{\sigma_1 \sqrt{\sigma_2(1+\sigma_2)} + \sigma_2 \sqrt{\sigma_1(1+\sigma_1)}}. \end{aligned} \quad (\text{A.33})$$

よって、式 (A.32)、式 (A.24)、式 (A.31)、式 (A.28) を用いて、式 (A.23) の第 4 番目である  $x_H$  を表す式を書き換えると、式 (A.34) が得られる。

$$\begin{aligned} x_H &= \frac{\alpha_1}{t_p + \sigma_1} - \frac{\alpha_2}{t_p + \sigma_2} \\ &= \Gamma \left\{ \frac{\alpha_1}{\sqrt{\sigma_1(1+\sigma_1)}} - \frac{\alpha_2}{\sqrt{\sigma_2(1+\sigma_2)}} \right\} \\ &= C \Gamma \left\{ \sqrt{\sigma_1(1+\sigma_1)} - \sqrt{\sigma_2(1+\sigma_2)} \right\} \\ &= C \Gamma^2 \times (\sigma_1 - \sigma_2) \\ &= \Gamma^2 x_L. \end{aligned} \quad (\text{A.34})$$

すると式 (A.34)、式 (A.31)、および (A.28) を用いて、式 (A.35) が導かれる。

$$\begin{cases} \frac{1}{\Gamma} = \sqrt{\frac{x_L}{x_H}} = \frac{\sqrt{\sigma_1(1+\sigma_1)} + \sqrt{\sigma_2(1+\sigma_2)}}{1 + \sigma_1 + \sigma_2}, \\ \frac{1}{x_L} = \frac{\mu(1 + \sigma_1 + \sigma_2) + \sigma_1 \sigma_2}{(\mu + \sigma_1)(\mu + \sigma_2)} = 1 - \frac{\mu(\mu - 1)}{(\mu + \sigma_1)(\mu + \sigma_2)}. \end{cases} \quad (\text{A.35})$$

この式 (A.35) から  $\sigma_1$  と  $\sigma_2$  についての連立方程式 (A.36) が得られるので、それを解いて  $\sigma_1$  と  $\sigma_2$  を求めればよい (ただし  $\sigma_1 > \sigma_2 > 0$  である)。

$$\begin{cases} \frac{\sqrt{\sigma_1(1+\sigma_1)} + \sqrt{\sigma_2(1+\sigma_2)}}{1 + \sigma_1 + \sigma_2} = \sqrt{\frac{x_L}{x_H}}, \\ (\mu + \sigma_1)(\mu + \sigma_2) = \mu(\mu - 1) \times \frac{x_L}{x_L - 1}. \end{cases} \quad (\text{A.36})$$

そこでいま (A.36) の上側の式に含まれる平方根をはずす

ために、式 (A.37) で表される変数の置換を行う。

$$\begin{cases} \sigma_1 \equiv \frac{z_1^2}{1 - z_1^2}, \\ \sigma_2 \equiv \frac{z_2^2}{1 - z_2^2}. \end{cases} \quad (\text{A.37})$$

ただし  $0 < z_1 < 1$ ,  $0 < z_2 < 1$  で、さらに  $\sigma_1 > \sigma_2$  であるから、 $1 > z_1 > z_2 > 0$  である。

するとこの置換により (A.36) の上側の式の左辺は、式 (A.38) の右辺に書き換えられる。

$$\frac{\sqrt{\sigma_1(1+\sigma_1)} + \sqrt{\sigma_2(1+\sigma_2)}}{1 + \sigma_1 + \sigma_2} = \frac{z_1 + z_2}{1 + z_1 z_2}. \quad (\text{A.38})$$

よって (A.36) の上側の式は、式 (A.39) に書き換えられる。

$$z_1 + z_2 = (1 + z_1 z_2) \sqrt{\frac{x_L}{x_H}}. \quad (\text{A.39})$$

(さらに  $t_p = \frac{z_1 z_2}{1 + z_1 z_2}$  であることもわかる)。

いま  $S_1$  と  $S_2$  をそれぞれ式 (A.40) で表される  $z_1$  と  $z_2$  の基本対称式とする。

$$\begin{cases} S_1 \equiv z_1 + z_2, \\ S_2 \equiv z_1 z_2. \end{cases} \quad (\text{A.40})$$

そうすると関係式 (A.39) は式 (A.41) に書き換えられる。

$$S_1 = (1 + S_2) \sqrt{\frac{x_L}{x_H}}. \quad (\text{A.41})$$

さらに (A.36) の下側の式について、 $\sigma_1$  と  $\sigma_2$  を  $z_1$  と  $z_2$  を用いて書き換えて、式 (A.42) が得られる。

$$(z_1^2 - \kappa)(z_2^2 - \kappa) = \nu \kappa (z_1^2 - 1)(z_2^2 - 1). \quad (\text{A.42})$$

ただしここで導入した記号  $\kappa$  と  $\nu$  は、式 (A.43) で表されるものである。

$$\begin{cases} \kappa \equiv \frac{\mu}{\mu - 1}, \\ \nu \equiv \frac{x_L}{x_L - 1}. \end{cases} \quad (\text{A.43})$$

式 (A.42) を整理することで式 (A.44) が得られる。

$$\eta_0 (z_1 z_2)^2 + \eta_1 (z_1^2 + z_2^2) + \eta_2 = 0. \quad (\text{A.44})$$

ただしここで導入した 3 つの係数  $\eta_0$ ,  $\eta_1$ ,  $\eta_2$  はそれぞれ、式 (A.45) により与えられる。

$$\begin{cases} \eta_0 \equiv 1 - \nu \kappa, \\ \eta_1 \equiv (\nu - 1) \kappa, \\ \eta_2 \equiv (\kappa - \nu) \kappa. \end{cases} \quad (\text{A.45})$$

式 (A.44) を  $z_1$  と  $z_2$  の基本対称式 (A.40) を用いて書き換えて、式 (A.46) が得られる。

$$\eta_0 S_2^2 + \eta_1 (S_1^2 - 2S_2) + \eta_2 = 0. \quad (\text{A.46})$$

式 (A.41) の関係を用いて、式 (A.46) から  $S_1$  を消去する

と、 $S_2$  についての2次方程式 (A.47) が得られる。

$$\zeta_0 S_2^2 + \zeta_1 S_2 + \zeta_2 = 0. \quad (\text{A.47})$$

ただしこの方程式の各係数は式 (A.48) により与えられる。

$$\begin{cases} \zeta_0 \equiv \eta_0 + \frac{x_L}{x_H} \times \eta_1, \\ \zeta_1 \equiv 2 \left( \frac{x_L}{x_H} - 1 \right) \eta_1, \\ \zeta_2 \equiv \eta_2 + \frac{x_L}{x_H} \times \eta_1. \end{cases} \quad (\text{A.48})$$

式 (A.47) の係数  $\zeta_0$  と  $\zeta_1$  は共に負である。

実際  $\mu > 1$ ,  $x_L > 1$  であることから,  $\kappa \equiv \frac{\mu}{\mu-1} > 1$ ,  $\nu \equiv \frac{x_L}{x_L-1} > 1$  であるから,  $\nu\kappa > 1$  であるので,  $\zeta_0 \equiv 1 - \nu\kappa < 0$  である。

また  $\nu-1 > 0$  であることから  $\eta_1 \equiv (\nu-1)\kappa > 0$  であり, さらに  $x_H > x_L > 1$  より  $\frac{x_L}{x_H} < 1$  であるから,  $\frac{x_L}{x_H} - 1 < 0$  なので  $\zeta_1 \equiv 2 \left( \frac{x_L}{x_H} - 1 \right) \eta_1 < 0$  である。

すると  $S_2$  についての2次方程式 (A.47) が正根を持つためには  $\zeta_2$  が正であることが必要である。そうして  $\zeta_2$  が正であるときには, 2次方程式 (A.47) の判別式  $D_1 \equiv \zeta_1^2 - 4\zeta_0\zeta_2$  は正で実根は2つあるが, 根と係数の関係から正根は単一である。

いま  $S_2$  についての2次方程式 (A.47) が区間 (0,1) に入る実根を持つとする (そうでなければ  $\sigma_1$  と  $\sigma_2$  には適切な解は無い)。そのような  $S_2$  が存在するとき, 式 (A.41) を用いて  $S_2$  から  $S_1$  を作る。

そうして, 2次方程式  $w^2 - S_1w + S_2 = 0$  の相異なる2つの実根がどちらも区間 (0,1) にあるとき, それらを  $z_1$  と  $z_2$  ( $1 > z_1 > z_2 > 0$ ) とする (そのような2つの根  $z_1$  と  $z_2$  が無ければ, 適切な  $\sigma_1$  と  $\sigma_2$  も存在しない)。そうして式 (A.37) の関係を用いて,  $z_1$  と  $z_2$  の値から  $\sigma_1$  と  $\sigma_2$  の値を求める。式 (A.49) を順に計算することで2つの極の係数  $\alpha_1$  と  $\alpha_2$  が求まる。

$$\begin{cases} C \leftarrow \frac{x_L}{\sigma_1 - \sigma_2}, \\ \alpha_1 \leftarrow C\sigma_1(1 + \sigma_1), \\ \alpha_2 \leftarrow C\sigma_2(1 + \sigma_2). \end{cases} \quad (\text{A.49})$$

以上の手順により, 式 (A.2) の  $x(t)$  が決定される。

#### A.1.2.1 「方式 II」の設計法のまとめ

4つのパラメタ ( $\mu, g_p, g_s, n$ ) が与えられたとき, 以下の式 (A.50) を順に計算する。(ここで  $x'_H = x_H - 1$ ,  $x'_L = x_L - 1$  の意味である。)

$$\begin{cases} x_H \leftarrow \cosh^2 \left( \frac{1}{2n} \cosh^{-1} \frac{1}{g_s} \right), \\ x'_H \leftarrow \sinh^2 \left( \frac{1}{2n} \cosh^{-1} \frac{1}{g_s} \right), \\ x_L \leftarrow \cosh^2 \left( \frac{1}{2n} \cosh^{-1} \frac{g_p}{g_s} \right), \\ x'_L \leftarrow \sinh^2 \left( \frac{1}{2n} \cosh^{-1} \frac{g_p}{g_s} \right), \\ \kappa \leftarrow \frac{\mu}{\mu-1}, \\ \zeta_0 \leftarrow 1 - \frac{x_L}{x'_L} \times \frac{x'_H}{x_H} \times \kappa, \\ \zeta_1 \leftarrow -2 \times \frac{\kappa}{x'_L} \times \frac{x_H - x_L}{x_H}, \\ \zeta_2 \leftarrow \left( \kappa - \frac{x_L}{x'_L} \times \frac{x'_H}{x_H} \right) \kappa, \\ D_1 \leftarrow \zeta_1^2 - 4\zeta_0\zeta_2. \end{cases} \quad (\text{A.50})$$

そうして,  $\zeta_2 \leq 0$  であれば構成不能として終了し,  $\zeta_2 > 0$  の場合にはさらに次の式 (A.51) を順に計算する。

$$\begin{cases} S_2 \leftarrow \frac{2\zeta_2}{-\zeta_1 + \sqrt{D_1}}, \\ S_1 \leftarrow (1 + S_2) \sqrt{\frac{x_L}{x_H}}, \\ D_2 \leftarrow S_1^2 - 4S_2. \end{cases} \quad (\text{A.51})$$

次に, 2次方程式  $w^2 - S_1w + S_2 = 0$  が相異なる正の実数 ( $z_1 > z_2 > 0$ ) を持つための必要十分条件「 $S_1 > 0$  かつ  $S_2 > 0$  かつ  $D_2 \equiv S_1^2 - 4S_2 > 0$ 」を調べて, 満たしていなければ構成は不可能なので終了する。満たしている場合は, 式 (A.52) を用いて2次方程式の相異なる2つの正の実数解  $z_1$  と  $z_2$  の値 ( $z_1 > z_2 > 0$ ) を計算する。

$$\begin{cases} z_1 \leftarrow \frac{1}{2} (S_1 + \sqrt{D_2}), \\ z_2 \leftarrow \frac{S_2}{z_1}. \end{cases} \quad (\text{A.52})$$

このとき  $z_1 \geq 1$  であれば (条件  $1 > z_1 > z_2 > 0$  を満たせない), 構成は不可能であり終了する。それ以外の場合には構成は可能であり, 式 (A.53) を順に計算する。(ここでは,  $\sigma_1 - \sigma_2 = (1 + \sigma_1)(1 + \sigma_2)S_1\sqrt{D_2}$  となることを用いている。)

$$\begin{cases} \sigma_1 \leftarrow \frac{z_1^2}{(1 - z_1)(1 + z_1)}, \\ \sigma_2 \leftarrow \frac{z_2^2}{(1 - z_2)(1 + z_2)}, \\ C \leftarrow \frac{x_L}{(1 + \sigma_1)(1 + \sigma_2)S_1\sqrt{D_2}}, \\ \alpha_1 \leftarrow C\sigma_1(1 + \sigma_1), \\ \alpha_2 \leftarrow C\sigma_2(1 + \sigma_2). \end{cases} \quad (\text{A.53})$$

以上の手順により, 構成が可能である場合には, 式 (A.2) の  $x(t)$  が含む値の組 ( $\sigma_1, \alpha_1, \sigma_2, \alpha_2$ ) を決定できる。そうして  $\sigma_1 > \sigma_2 > 0$  かつ  $\alpha_1 > \alpha_2 > 0$  となる。

### A.1.3 「方式 I」と「方式 II」で構成された関数 $g(t)$ の振る舞いの確認

まず式 (A.2) の  $x(t)$  は「方式 I」と「方式 II」のどちらの場合にも阻止域  $t \in [\mu, \infty)$  において  $1 \geq x(t) > 0$  を満たすことを確認する (この条件を満たせば、式 (A.3) の伝達関数の大きさ  $|g(t)|$  は阻止域において常に  $g_s$  以下となる). そのためには具体的に構成された  $x(t)$  は、「方式 I」と「方式 II」どちらの場合も  $\alpha_1 > \alpha_2 > 0$  を満たすことを利用する.

いま式 (A.2) を変形すれば式 (A.54) が得られる.

$$x(t) = \frac{(\alpha_1 - \alpha_2)t + (\alpha_1\sigma_2 - \alpha_2\sigma_1)}{(t + \sigma_1)(t + \sigma_2)}. \quad (\text{A.54})$$

いま  $x(0) \geq x_L > 0$  であることを式 (A.54) にあてはめてみると、 $\sigma_1$  と  $\sigma_2$  は正なので、不等式 (A.55) がなりたつことがわかる.

$$\alpha_1\sigma_2 - \alpha_2\sigma_1 > 0. \quad (\text{A.55})$$

式 (A.54) とこの不等式 (A.55), および  $\sigma_1$  と  $\sigma_2$  が共に正であり、また  $\alpha_1 > \alpha_2$  であることを使うと、 $t \geq 0$  であるときには  $x(t) > 0$  であることがわかる.

次に  $x(t)$  は「方式 I」と「方式 II」どちらの場合にも、阻止域  $[\mu, \infty)$  で単調減少であることを示すために、 $x(t)$  の導関数の式 (A.56) について調べてみる.

$$\begin{aligned} x'(t) &= -\frac{\alpha_1}{(t + \sigma_1)^2} + \frac{\alpha_2}{(t + \sigma_2)^2} \\ &= -\frac{Q(t)}{(t + \sigma_1)^2(t + \sigma_2)^2}. \end{aligned} \quad (\text{A.56})$$

ここで  $Q(t)$  は式 (A.57) で与えられる  $t$  の 2 次式多項式である.

$$Q(t) = (\alpha_1 - \alpha_2)t^2 + 2(\alpha_1\sigma_2 - \alpha_2\sigma_1)t + \alpha_1\sigma_2^2 - \alpha_2\sigma_1^2. \quad (\text{A.57})$$

この 2 次方程式  $Q(t) = 0$  の判別式を求めてみると  $4\alpha_1\alpha_2(\sigma_1 - \sigma_2)^2$  となり、その値は「方式 I」と「方式 II」どちらの場合にも必ず正になるので、2 次多項式  $Q(t)$  は必ず相異なる 2 つの実数  $t_1$  と  $t_2$  を零点として持つことがわかる (その大小関係を  $t_1 < t_2$  と決める). いま 2 次方程式の解と係数の関係から式 (A.58) が導かれる.

$$t_1 + t_2 = -2 \times \frac{\alpha_1\sigma_2 - \alpha_2\sigma_1}{\alpha_1 - \alpha_2}. \quad (\text{A.58})$$

この式 (A.58) の値は、「方式 I」と「方式 II」どちらの場合にも負であることが、式 (A.55) と  $\alpha_1 > \alpha_2$  であることからわかる. よって少なくとも  $t_1$  と  $t_2$  のうちのどちらかは値が負である. すると  $t_1 < t_2$  であると決めておいたので、 $t_1$  の値は必ず負である.

- 「方式 I」ではそれを構成したときの条件から  $x'(0) = 0$  である. つまり  $t = 0$  は  $Q(t)$  の零点である. すると ( $t_1$  の値は負であることから)  $t_2 = 0$  である. 主係数が正である 2 次式  $Q(t)$  の値は  $t > t_2 = 0$  のときには

正であるから、式 (A.56) により  $t > 0$  のときには必ず  $x'(t) < 0$  である. よって  $x(t)$  は  $t > 0$  で単調減少である.

- 「方式 II」ではその構成から  $x'(t_p) = 0$  ( $t_p \in (0, 1)$ ) である. つまり  $t = t_p$  が  $Q(t)$  の零点である. すると ( $t_1$  の値は負であるから)  $t_2 = t_p$  であることがわかる. 主係数が正である 2 次式  $Q(t)$  の値は  $t > t_2 = t_p$  のときには正であるから、式 (A.56) により  $t > t_p$  のときには必ず  $x'(t) < 0$  となる. よって  $x(t)$  は  $t_p < t$  のときには単調減少である (同様に  $x(t)$  は  $0 \leq t < t_p$  のときには単調増加であることが示せる).

さらに、 $x(t)$  を構成したときの条件から  $x(\mu) = 1$  であり、また式 (A.2) から  $\lim_{t \rightarrow \infty} x(t) = 0$  である. よって上記のこととあわせると、「方式 I」と「方式 II」どちらの場合にも、 $x(t)$  は阻止域  $t \in [\mu, \infty)$  において単調減少で  $(0, 1]$  への全射である

以上の結果をまとめると、「方式 I」と「方式 II」どちらの場合にも阻止域  $t \in [\mu, \infty)$  では  $|g(t)| \leq g_s$  であり、さらに遷移域  $t \in (1, \mu)$  では  $g(t)$  は単調減少で  $[g_s, g_p]$  への全射である. そうして、通過域  $t \in [0, 1]$  では「方式 I」の  $g(t)$  は単調減少であるが「方式 II」の  $g(t)$  は  $0 \leq t < t_p$  では単調増加で、 $t_p < t \leq 1$  では単調減少であり、 $g(t)$  はどちらの方式でも通過域では  $[g_p, 1]$  への全射である.

### A.1.4 「方式 I」と「方式 II」の設計の関係性

いま通過域における「方式 II」の伝達関数  $g(t)$  の最大点を  $t = t_p$  とするとき、「方式 II」の本来の通過域  $[0, 1]$  から  $[0, t_p)$  の部分を取り除いた区間である  $t \in [t_p, 1]$  の全体を  $\tilde{t} \in [0, 1]$  の全体に写す線形の座標変換  $t = \mathcal{L}(\tilde{t})$  により  $\tilde{g}(\tilde{t}) = g(\mathcal{L}(\tilde{t}))$  としたものは「方式 I」の伝達関数とみなせるので、「方式 I」と「方式 II」の 3 つのパラメータ  $g_p$  と  $g_s$  と次数  $n$  を共通に設定した場合に、「方式 II」において伝達関数の遷移域  $(1, \mu)$  の幅と通過域  $[0, 1]$  の幅の比が  $\mu - 1$  であるならば、「方式 I」の場合に対応する比は  $\tilde{\mu} - 1 = \frac{\mu - 1}{1 - t_p}$  である. 「方式 II」の場合の比の値である  $\mu - 1$  と「方式 I」の場合の比の値である  $\tilde{\mu} - 1$ , その前者に対する後者の比の値を計算すると  $\frac{\mu - 1}{\tilde{\mu} - 1} = 1 - t_p$  であり、1 よりも小さい.

以上のことから、3 つのパラメータ  $g_p$  と  $g_s$  と次数  $n$  を共通にとるとき、「方式 I」よりも「方式 II」の方がフィルタを実現できる  $\mu$  の値を小さくできることがわかる.

### A.1.5 パラメータ 4 つのうち 3 つだけを指定するやり方

フィルタを実現可能なパラメータの 4 つ  $\mu$ ,  $g_p$ ,  $g_s$ ,  $n$  をすべて直接指定する以外の方法としては、たとえば以下に述べるように、3 つの値だけを直接指定して、残りの 1 つについては探索によりフィルタが実現可能な範囲でなるべ

く都合の良い値となるように決めることができる。

### 次数 $n$ を最小にする場合

実用性の観点から探索する次数の上限  $n_{\max}$  をあらかじめ設定して (たとえば 50 とする) 次数  $n$  を 1 から始めて  $n_{\max}$  まで 1 ずつ増してフィルタを実現可能にする最初のもの ( $n$  が最小のもの) を探す。探しても無ければ指定した 3 つのパラメータを持つフィルタは ( $n$  が上限  $n_{\max}$  以下の範囲では) 実現不可能である。

### $g_p$ を最大にする場合

まず定義から  $1 > g_p > g_s$  である。  $g_p$  の最大値はたとえば 2 分法で精密に決めることもできるが、簡易には本当の最大値ではなくてきりの良い値、たとえば単調に減少する  $0.5^j$  の形に制限してその指数  $j$  を 1 ずつ増やして ( $g_p > g_s$  の範囲で) フィルタを実現可能にする最初のもの ( $g_p$  が大きいもの) を探す。探しても無ければ指定した 3 つのパラメータを持つフィルタの構成を諦める。

### $g_s$ を最小にする場合

まず定義から  $g_p > g_s > 0$  である。  $g_s$  の最小値はたとえば 2 分法で精密に決めることもできるが、簡易にはきりの良い値、たとえば単調に減少する  $0.5^j$  の形に制限して、  $g_p > g_s$  を満たしてさらに現実的な考慮からたとえば丸め誤差の単位を  $\epsilon_{\text{MAC}}$  とするとき  $g_s > \epsilon_{\text{MAC}}$  の範囲で  $j$  の値を 1 ずつ増やして、フィルタを実現可能にする最後のもの ( $g_s$  が小さいもの) を探す。探しても無ければ指定した 3 つのパラメータを持つフィルタの構成を諦める。

### $\mu$ を最小にする場合

まず定義から  $\mu > 1$  である。  $\mu$  の最小値はたとえば 2 分法で精密に決めることもできるが、簡易にはきりの良い値に制限して、たとえば単調に増加する  $1 + 0.05j$  の形とし、整数  $j$  を 1 から始めて 1 ずつ増して探索することができる。ただし実用性を考えて  $\mu$  の値にはある上限を設けて (たとえば  $\mu \leq 2$ ) それ以下の範囲でだけ探索をするなどとする。そうしてフィルタを実現可能にする最初のもの ( $\mu$  が小さいもの) を探す。探しても無ければ指定した 3 つのパラメータを持つフィルタの構成を諦める。

## A.1.6 伝達関数からのフィルタの構成

極として実数 2 つだけを持つ簡易型の伝達関数に対応するフィルタは、シフトが実数であるレゾルベントを 2 つ用いて構成できる (実対称定値一般固有値問題 (1) に対応するシフトが  $\rho$  のレゾルベント  $\mathcal{R}(\rho)$  は  $\mathcal{R}(\rho) \equiv (A - \rho B)^{-1} B$  である)。

いま採用する簡易型の設計では、伝達関数  $g(t)$  は式 (A.2) と (A.3) により与えられる。下端付近の固有値を扱うので、  $\lambda \in [a, b]$  と  $t \in [0, 1]$  の間の線形対応関係は式 (A.59) で与えられる。

$$t = \frac{\lambda - a}{b - a}. \quad (\text{A.59})$$

すると  $y$  を  $\lambda$  の関数として表すと、以下の式 (A.60) になる。

$$y = \frac{2\ell_1}{\lambda - \rho_1} - \frac{2\ell_2}{\lambda - \rho_2} - 1. \quad (\text{A.60})$$

ここでシフト  $\rho_k$  と係数  $\ell_k$  ( $k = 1, 2$ ) は式 (A.61) により与えられる。

$$\rho_k \equiv a - (b - a)\sigma_k, \quad \ell_k \equiv (b - a)\alpha_k \quad (\text{A.61})$$

いま  $\sigma_1 > \sigma_2 > 0$  であることから  $\rho_1 < \rho_2 < a$  である。さらに「方式 I」と「方式 II」の場合には  $\alpha_1 > \alpha_2 > 0$  になることから、  $\ell_1 > \ell_2 > 0$  となることもわかる。

式 (A.60) の  $y$  に対応する線形作用素  $\mathcal{Y}$  は、  $\frac{1}{\lambda - \rho_k}$  にはレゾルベント  $\mathcal{R}(\rho_k)$  を、定数 1 には恒等作用素  $I$  を、それぞれ対応させた式 (A.62) になる。

$$\mathcal{Y} \equiv 2\ell_1 \mathcal{R}(\rho_1) - 2\ell_2 \mathcal{R}(\rho_2) - I \quad (\text{A.62})$$

そうして伝達関数  $f(\lambda)$  に対応する作用素であるフィルタ  $\mathcal{F}$  は、作用素  $\mathcal{Y}$  の多項式として式 (2) で表せる。

ベクトルの組  $V$  に式 (2) の形のフィルタ  $\mathcal{F}$  を作用させる計算には、Chebyshev 多項式のもつ 3 項漸化式  $T_0(z) = I$ ,  $T_1(z) = z$ ,  $T_j(z) = 2zT_{j-1}(z) - T_{j-2}(z)$  ( $j \geq 2$ ) を利用する。具体的には、  $\mathcal{Y}$  の  $j$  次 Chebyshev 多項式  $T_j(\mathcal{Y})$  を  $V$  に作用させて得られるベクトルの組  $V^{(j)} \equiv T_j(\mathcal{Y})V$  を以下の漸化式 (A.63) を用いて計算する。

$$\begin{cases} V^{(0)} = V, \\ V^{(1)} = \mathcal{Y}V, \\ V^{(j)} = 2\mathcal{Y}V^{(j-1)} - V^{(j-2)} \quad (j \geq 2). \end{cases} \quad (\text{A.63})$$

すると  $V$  から始めて漸化式 (A.63) により  $V^{(n)}$  を求めれば、ベクトルの組  $V$  に式 (2) のフィルタ  $\mathcal{F}$  を作用させた結果であるベクトルの組は式 (A.64) の右辺で与えられる。

$$\mathcal{F}V = g_s V^{(n)} \quad (\text{A.64})$$

## A.2 今回の比較実験に用いた、実数シフトの単一のレゾルベントによるフィルタ (拡張版) の構成

シフトが実数であるレゾルベント 1 つから Chebyshev 多項式を用いた簡易構成で作られるフィルタの伝達関数  $g(t)$  としてこれまででは式 (A.1) で表されるものを用いてきた。

今回はそれを拡張して、新たに式 (A.65) で表される伝達関数を導入する ( $\beta$  の値を  $-1$  に制限した場合には従来のものに帰着する)。

$$g(t) = g_s T_n(y(t)), \quad y(t) = \frac{\alpha}{t + \sigma} + \beta. \quad (\text{A.65})$$

ただし  $\sigma > 0$ ,  $\alpha > 0$  であり、有理関数  $y(t)$  は非負領域  $t \in [0, \infty)$  において連続かつ単調減少になる。そうしてこの伝達関数の形状のパラメータ ( $\mu, g_p, g_s$ ) (ただし  $\mu > 1$ ,

$1 > g_p > g_s > 0$  である) についての条件として, 形状パラメタ  $\mu, g_p, g_s$  のそれぞれの意味に基づいてこれまでと同じように  $g(0) = 1, g(1) = g_p, g(\mu) = g_s$  であること, および  $t \in [\mu, \infty)$  において  $g_s \geq |g(t)|$  であることを要請する.

すると 4 つのパラメタの組  $(\mu, g_p, g_s, n)$  を指定したときに, 式 (A.65) の形の伝達関数  $g(t)$  の実現の可能性の判定と, 実現が可能なときに  $x(t)$  を決める 3 つの実数  $\sigma, \alpha, \beta$  の値を与える手順は以下のようにまとめられる.

まず (A.66) を順に計算する.

$$\begin{cases} y_H \leftarrow \cosh\left(\frac{1}{n} \cosh^{-1} \frac{1}{g_s}\right), \\ y_L \leftarrow \cosh\left(\frac{1}{n} \cosh^{-1} \frac{g_p}{g_s}\right), \\ \sigma \leftarrow \frac{(y_L - 1)\mu}{(y_H - y_L)\mu - (y_H - 1)}, \\ \alpha \leftarrow (y_H - y_L)\sigma \times (\sigma + 1), \\ \beta \leftarrow y_L - (y_H - y_L)\sigma. \end{cases} \quad (\text{A.66})$$

これにより得られた計算結果が  $\sigma > 0$  かつ  $\beta \geq -1$  であるときに限って要請を満たすフィルタが実現可能である (なぜならば  $y(t)$  が非負領域全体で連続関数になるためには  $\sigma > 0$  であることが必要十分で, そうして  $1 > g_p$  より  $y_H > y_L$  であり,  $y(0) = y_H, y(1) = y_L$  であるから  $y(t)$  は非負領域において単調減少関数となり, さらに  $y(\mu) = 1$  と決めるので, 阻止域で  $g_s \geq |g(t)|$  となるための必要十分条件は  $\beta \geq -1$  である).

以上の手順により, 4 つのパラメタの組  $(\mu, g_p, g_s, n)$  を指定したときに, シフトが実数である単一のレゾルベントを用いた式 (A.65) で表される伝達関数  $g(t)$  を持つフィルタの実現の可能性が判定できる. そうして実現が可能である場合には与えられた任意のベクトルに対するフィルタの作用を実現する具体的な手続きを構成できる. また, (レゾルベントを 2 つ使う場合と同様に) 4 つのパラメタのうちの 3 つの値だけを直接指定して, 残りの 1 つの値はフィルタが実現可能な範囲でもってフィルタの性質が最も都合の良いものとなるものを探索して決めることもできる.

式 (A.65) の有理関数  $y(t)$  を決定する 3 つの実数である  $\sigma, \alpha, \beta$  の値が決まれば, 求めようとする固有対の固有値が固有値分布の下端の区間  $[a, b]$  に含まれているとすると, その区間を通過域とするフィルタ  $\mathcal{F}$  は式 (A.67) で与えられる.

$$\begin{cases} \mathcal{F} \equiv g_s T_n(\mathcal{Y}), \quad \mathcal{Y} \equiv \gamma \mathcal{R}(\rho) + \beta I, \\ \rho = a - (b - a)\sigma, \quad \gamma = (b - a)\alpha. \end{cases} \quad (\text{A.67})$$

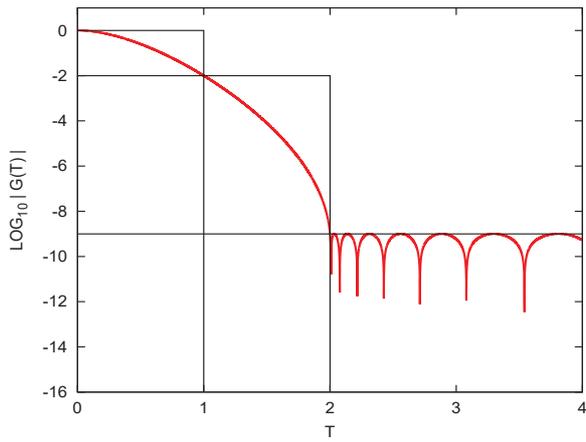


図 A・1 実験 1 : フィルタ I-1 : 伝達関数の大きさ  $|g(t)|$  ( $\mu=2.0$ ,  $g_p=1E-2$ ,  $g_s=1E-09$ ,  $n=25$ )

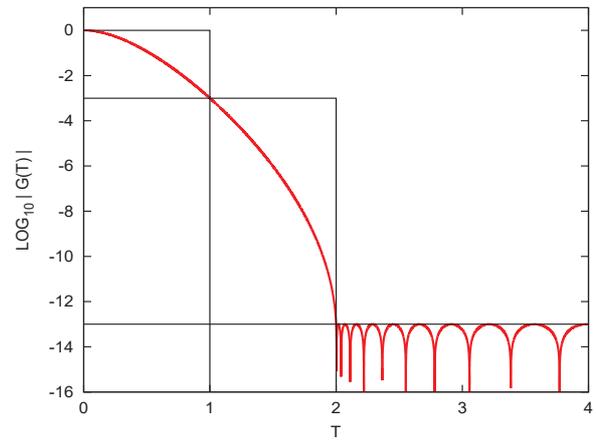


図 A・4 実験 1 : フィルタ I-4 : 伝達関数の大きさ  $|g(t)|$  ( $\mu=2.0$ ,  $g_p=1E-3$ ,  $g_s=1E-13$ ,  $n=32$ )

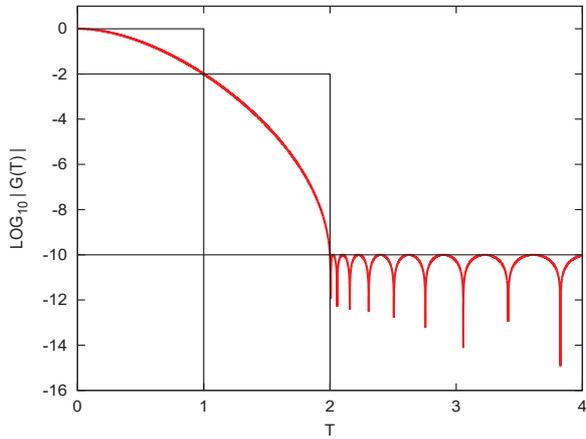


図 A・2 実験 1 : フィルタ I-2 : 伝達関数の大きさ  $|g(t)|$  ( $\mu=2.0$ ,  $g_p=1E-2$ ,  $g_s=1E-10$ ,  $n=35$ )

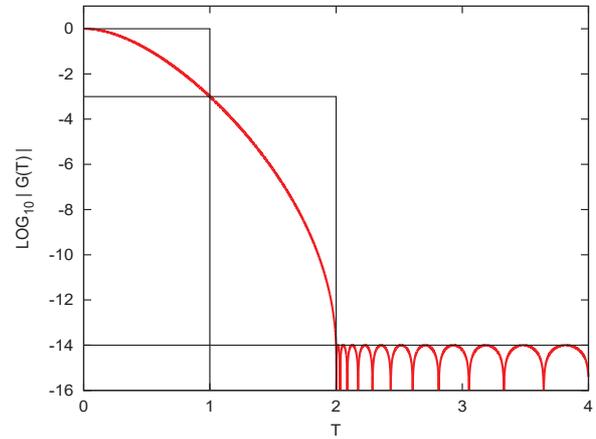


図 A・5 実験 1 : フィルタ I-5 : 伝達関数の大きさ  $|g(t)|$  ( $\mu=2.0$ ,  $g_p=1E-3$ ,  $g_s=1E-14$ ,  $n=40$ )

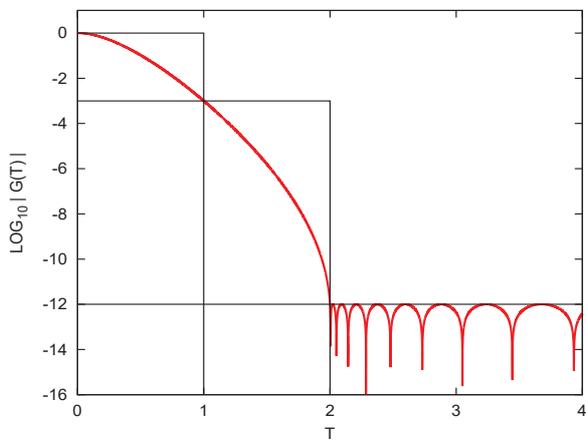


図 A・3 実験 1 : フィルタ I-3 : 伝達関数の大きさ  $|g(t)|$  ( $\mu=2.0$ ,  $g_p=1E-3$ ,  $g_s=1E-12$ ,  $n=25$ )

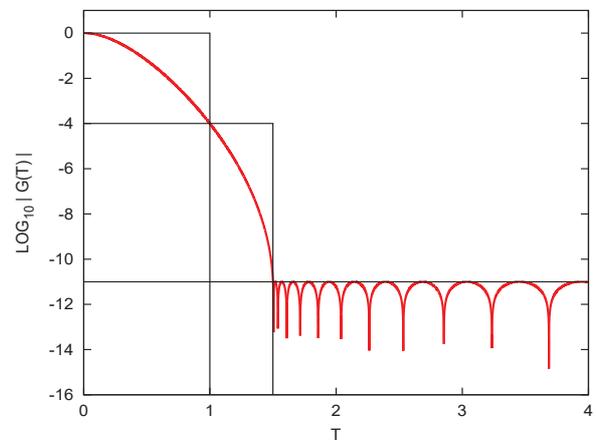


図 A・6 実験 1 : フィルタ I-6 : 伝達関数の大きさ  $|g(t)|$  ( $\mu=1.5$ ,  $g_p=1E-4$ ,  $g_s=1E-11$ ,  $n=30$ )

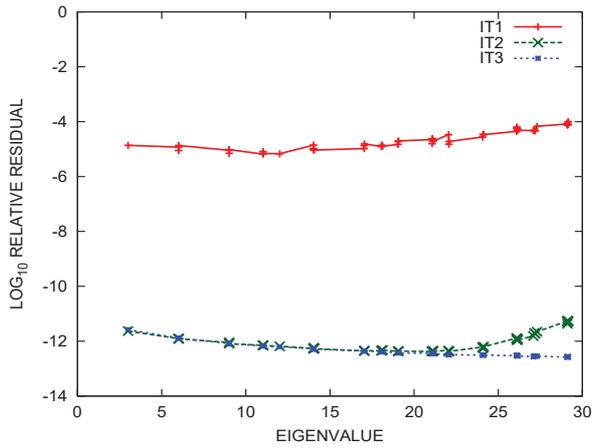


図 A-7 実験 1: フィルタ I-1: 反復回数ごとの各近似固有対の相対残差 ( $\mu=2.0$ ,  $g_p=1E-2$ ,  $g_s=1E-09$ ,  $n=25$ )

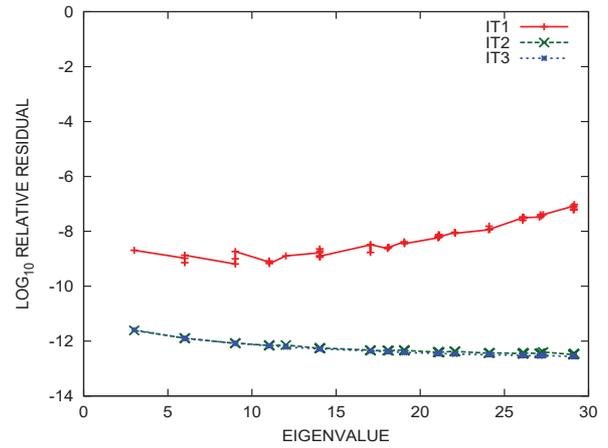


図 A-10 実験 1: フィルタ I-4: 反復回数ごとの各近似固有対の相対残差 ( $\mu=2.0$ ,  $g_p=1E-3$ ,  $g_s=1E-13$ ,  $n=32$ )

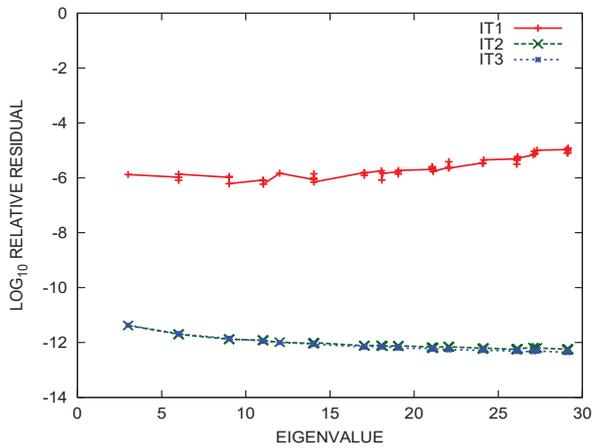


図 A-8 実験 1: フィルタ I-2: 反復回数ごとの各近似固有対の相対残差 ( $\mu=2.0$ ,  $g_p=1E-2$ ,  $g_s=1E-10$ ,  $n=35$ )

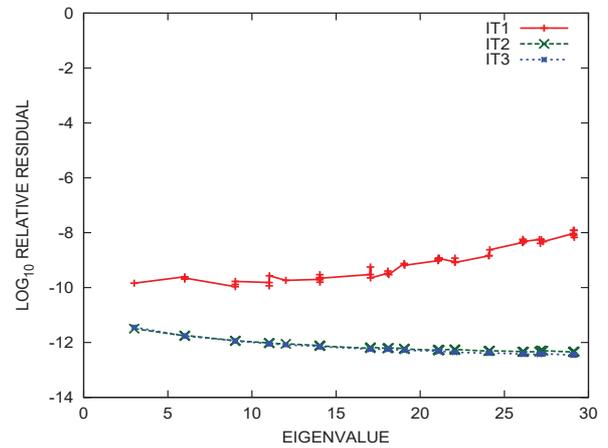


図 A-11 実験 1: フィルタ I-5: 反復回数ごとの各近似固有対の相対残差 ( $\mu=2.0$ ,  $g_p=1E-3$ ,  $g_s=1E-14$ ,  $n=40$ )

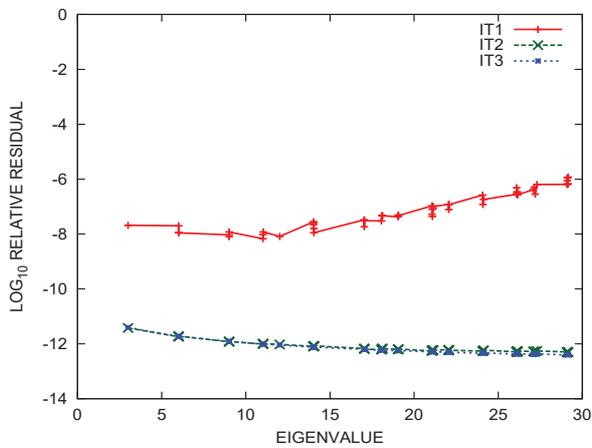


図 A-9 実験 1: フィルタ I-3: 反復回数ごとの各近似固有対の相対残差 ( $\mu=2.0$ ,  $g_p=1E-3$ ,  $g_s=1E-12$ ,  $n=25$ )

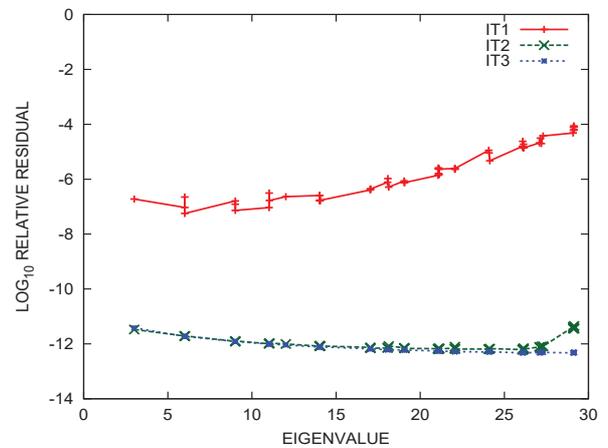


図 A-12 実験 1: フィルタ I-6: 反復回数ごとの各近似固有対の相対残差 ( $\mu=1.5$ ,  $g_p=1E-4$ ,  $g_s=1E-11$ ,  $n=30$ )

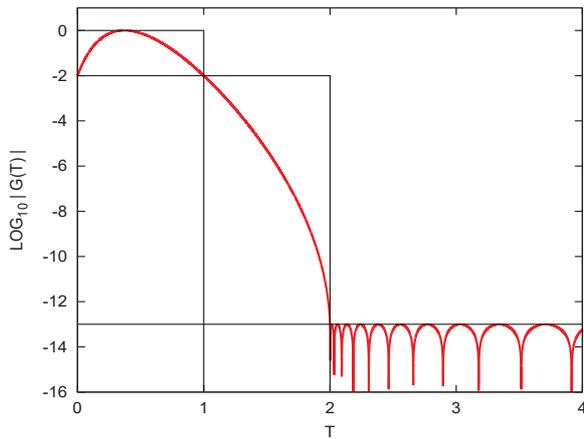


図 A-13 実験 1: フィルタ II-1: 伝達関数の大きさ  $|g(t)|$  ( $\mu=2.0$ ,  $g_p=1E-2$ ,  $g_s=1E-13$ ,  $n=30$ )

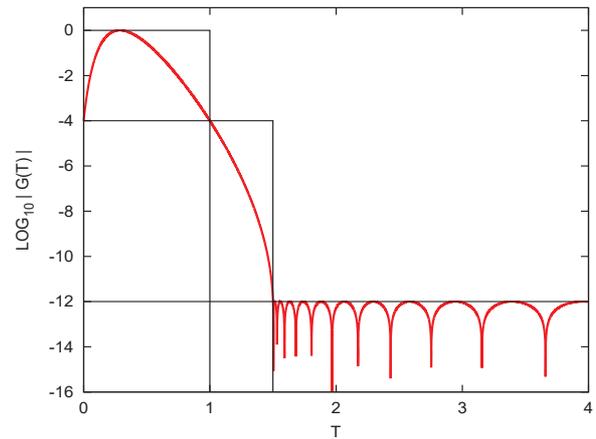


図 A-16 実験 1: フィルタ II-4: 伝達関数の大きさ  $|g(t)|$  ( $\mu=1.5$ ,  $g_p=1E-4$ ,  $g_s=1E-12$ ,  $n=24$ )

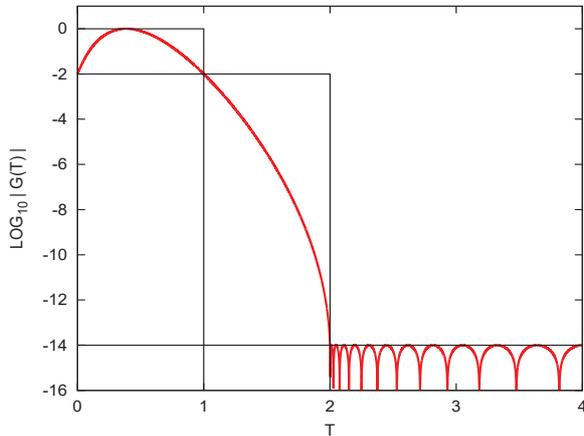


図 A-14 実験 1: フィルタ II-2: 伝達関数の大きさ  $|g(t)|$  ( $\mu=2.0$ ,  $g_p=1E-2$ ,  $g_s=1E-14$ ,  $n=35$ )

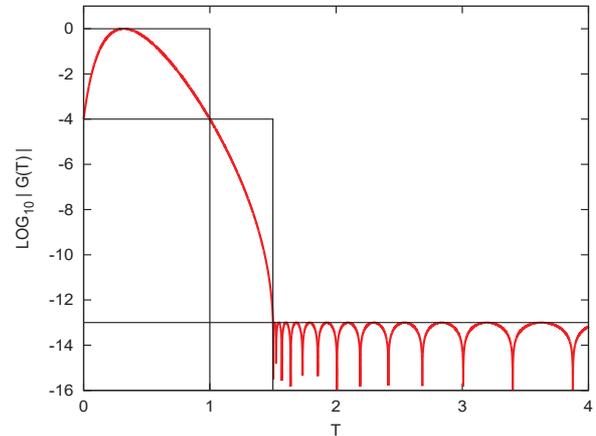


図 A-17 実験 1: フィルタ II-5: 伝達関数の大きさ  $|g(t)|$  ( $\mu=1.5$ ,  $g_p=1E-4$ ,  $g_s=1E-13$ ,  $n=28$ )

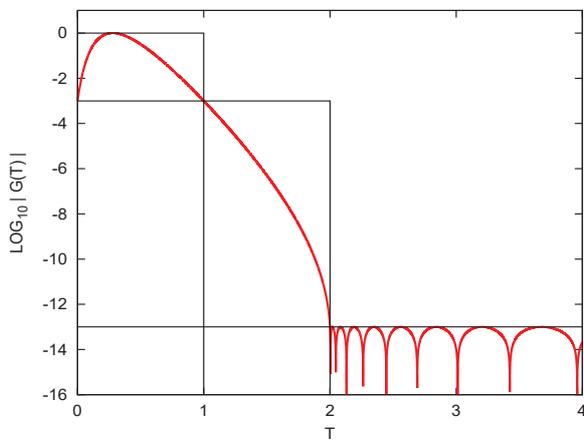


図 A-15 実験 1: フィルタ II-3: 伝達関数の大きさ  $|g(t)|$  ( $\mu=2.0$ ,  $g_p=1E-3$ ,  $g_s=1E-13$ ,  $n=21$ )

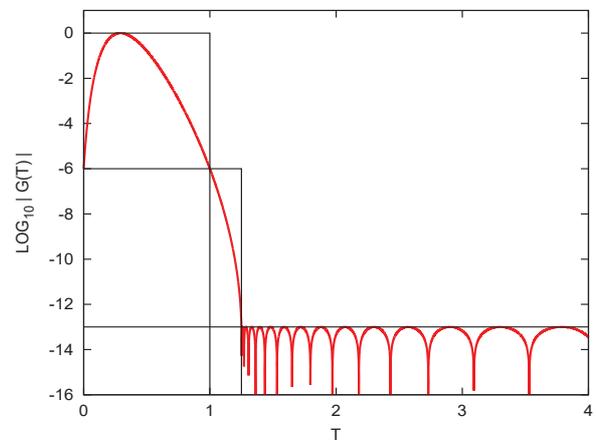


図 A-18 実験 1: フィルタ II-6: 伝達関数の大きさ  $|g(t)|$  ( $\mu=1.25$ ,  $g_p=1E-6$ ,  $g_s=1E-13$ ,  $n=29$ )

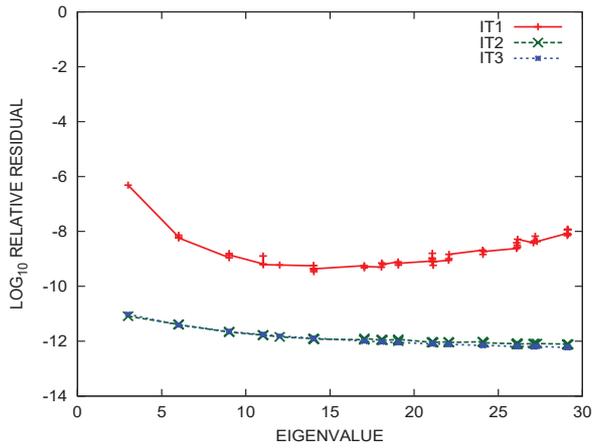


図 A-19 実験 1 : フィルタ II-1 : 反復回数ごとの各近似固有対の相対残差 ( $\mu=2.0$ ,  $g_p=1E-2$ ,  $g_s=1E-13$ ,  $n=30$ )

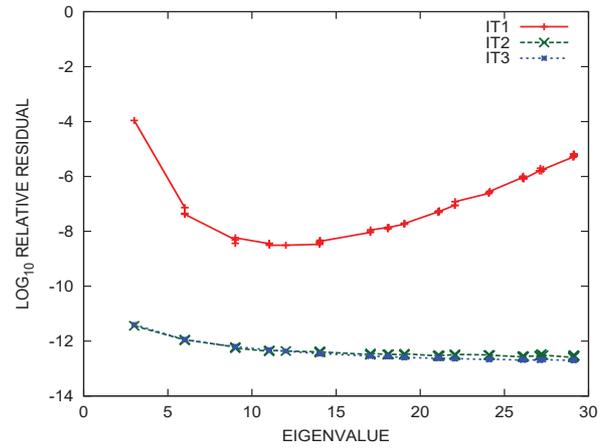


図 A-22 実験 1 : フィルタ II-4 : 反復回数ごとの各近似固有対の相対残差 ( $\mu=1.5$ ,  $g_p=1E-4$ ,  $g_s=1E-12$ ,  $n=24$ )

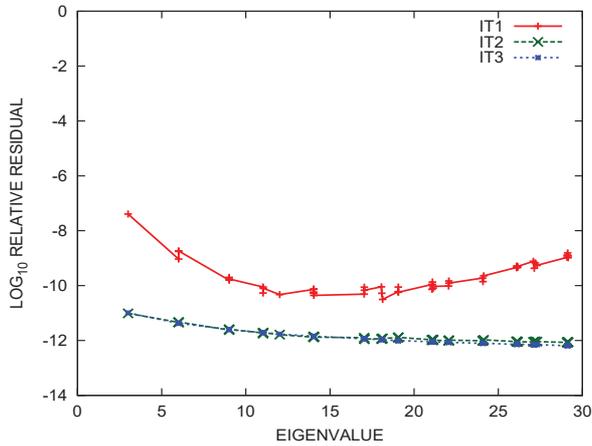


図 A-20 実験 1 : フィルタ II-2 : 反復回数ごとの各近似固有対の相対残差 ( $\mu=2.0$ ,  $g_p=1E-2$ ,  $g_s=1E-14$ ,  $n=35$ )

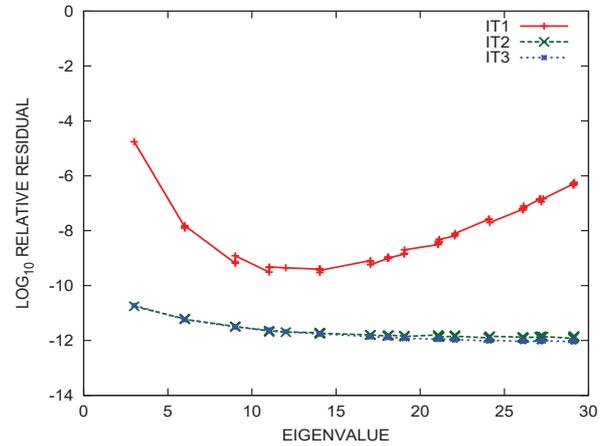


図 A-23 実験 1 : フィルタ II-5 : 反復回数ごとの各近似固有対の相対残差 ( $\mu=1.5$ ,  $g_p=1E-4$ ,  $g_s=1E-13$ ,  $n=28$ )

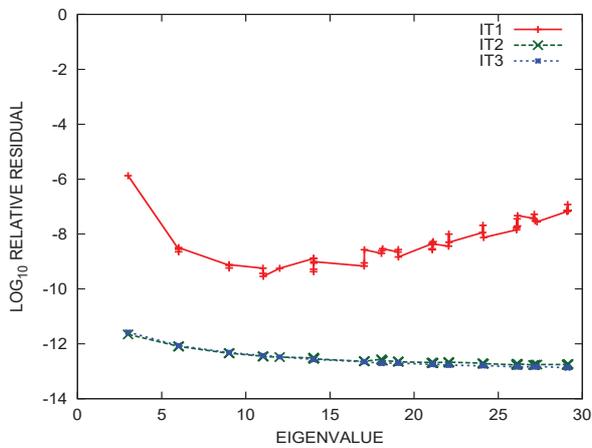


図 A-21 実験 1 : フィルタ II-3 : 反復回数ごとの各近似固有対の相対残差 ( $\mu=2.0$ ,  $g_p=1E-3$ ,  $g_s=1E-13$ ,  $n=21$ )

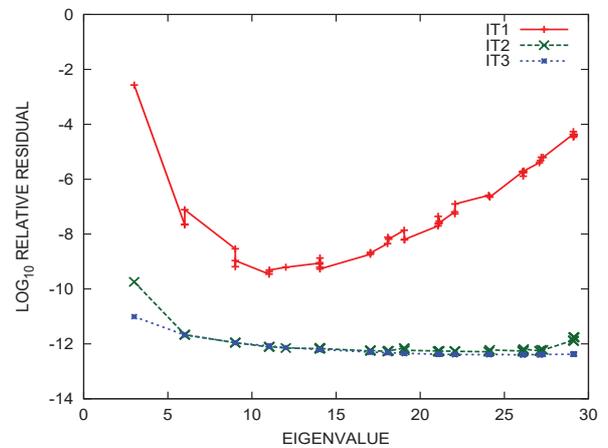


図 A-24 実験 1 : フィルタ II-6 : 反復回数ごとの各近似固有対の相対残差 ( $\mu=1.25$ ,  $g_p=1E-6$ ,  $g_s=1E-13$ ,  $n=29$ )

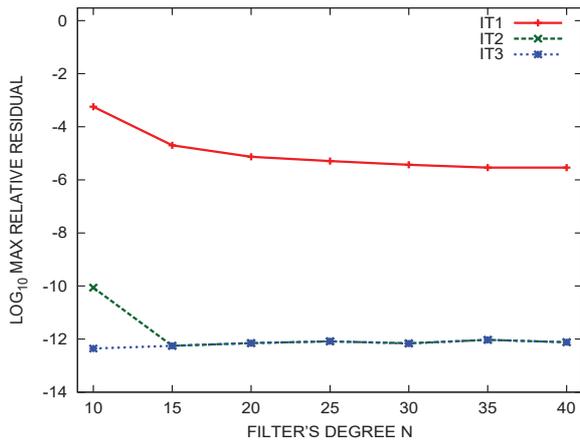


図 A-25 実験 2 : 「単一」フィルタの次数  $n$  と最大相対残差  
 $(\mu = 2.0 (m = 200), g_s = 1E-13)$

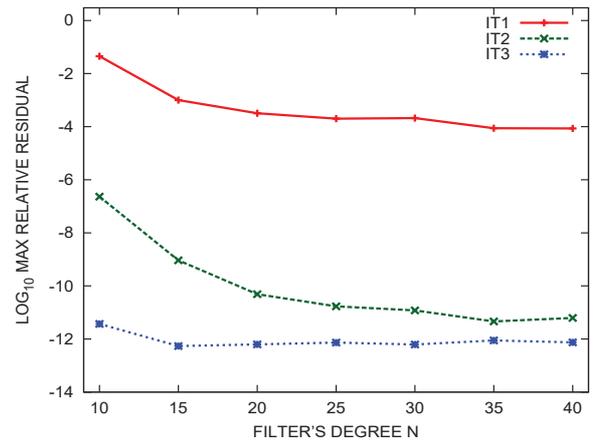


図 A-28 実験 2 : 「単一」フィルタの次数  $n$  と最大相対残差  
 $(\mu = 1.5 (m = 125), g_s = 1E-13)$

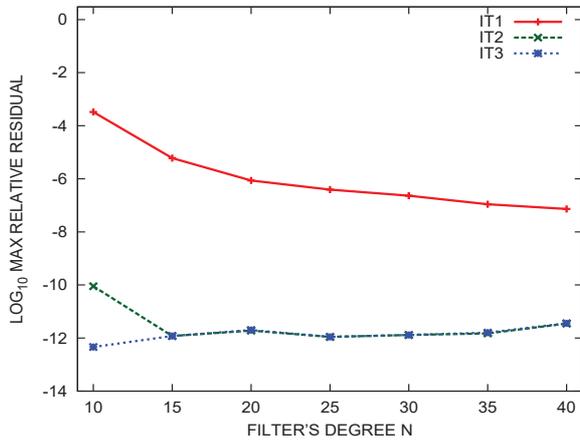


図 A-26 実験 2 : 「方式 I」フィルタの次数  $n$  と最大相対残差  
 $(\mu = 2.0 (m = 200), g_s = 1E-13)$

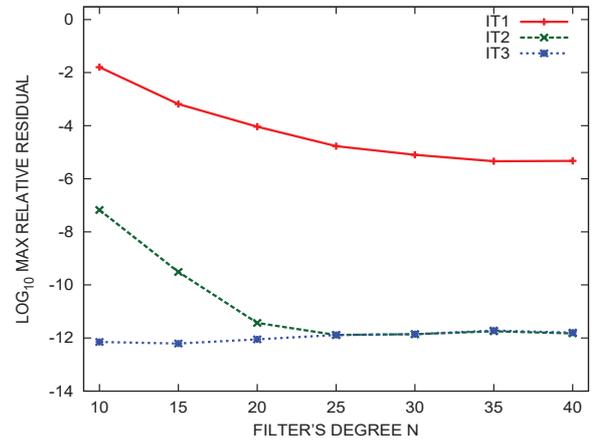


図 A-29 実験 2 : 「方式 I」フィルタの次数  $n$  と最大相対残差  
 $(\mu = 1.5 (m = 125), g_s = 1E-13)$

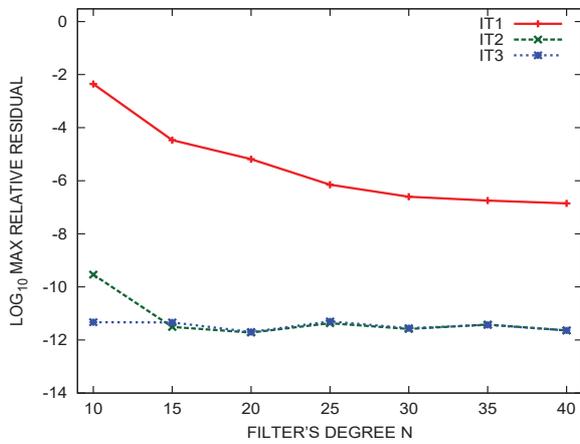


図 A-27 実験 2 : 「方式 II」フィルタの次数  $n$  と最大相対残差  
 $(\mu = 2.0 (m = 200), g_s = 1E-13)$

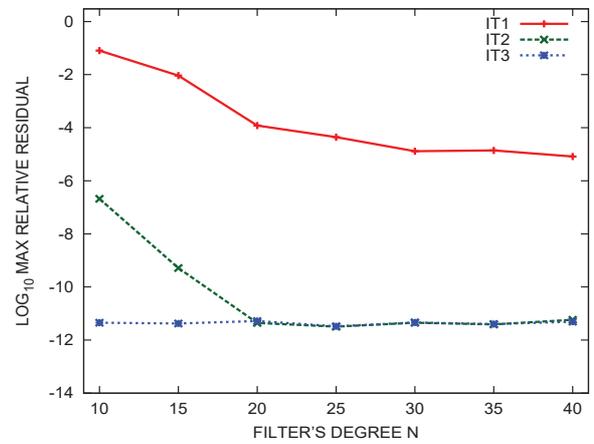


図 A-30 実験 2 : 「方式 II」フィルタの次数  $n$  と最大相対残差  
 $(\mu = 1.5 (m = 125), g_s = 1E-13)$

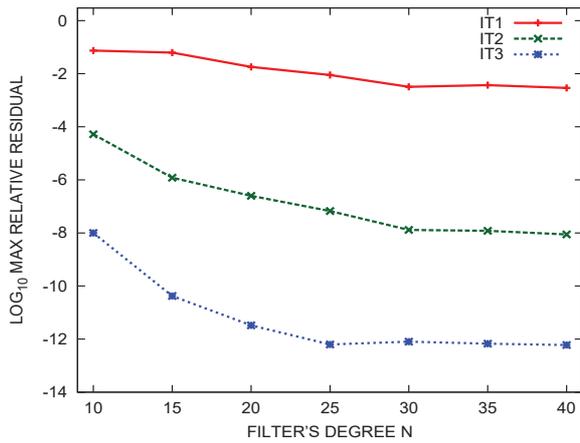


図 A-31 実験 2 : 「単一」フィルタの次数  $n$  と最大相対残差 ( $\mu = 1.25$  ( $m = 100$ ),  $g_s = 1E-13$ )

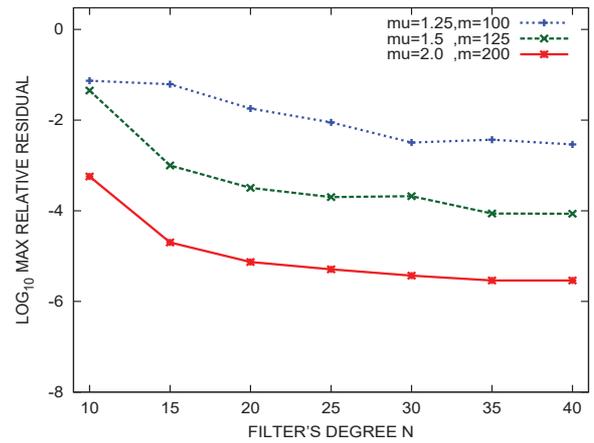


図 A-34 実験 2 : 「単一」フィルタの次数  $n$  と適用 1 回目の最大相対残差 ( $g_s = 1E-13$ )

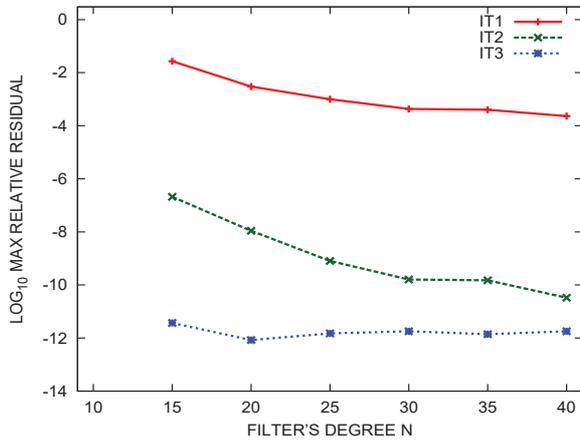


図 A-32 実験 2 : 「方式 I」フィルタの次数  $n$  と最大相対残差 ( $\mu = 1.25$  ( $m = 100$ ),  $g_s = 1E-13$ )

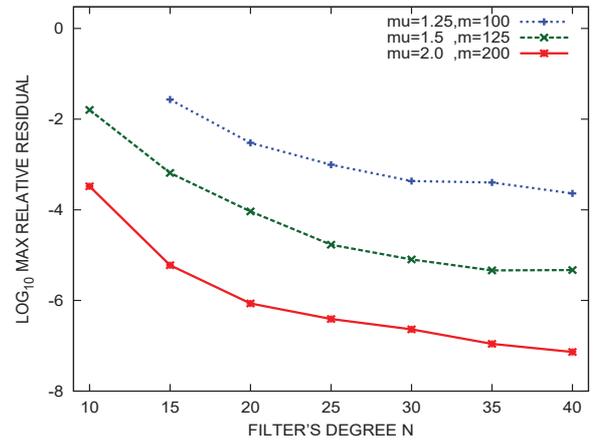


図 A-35 実験 2 : 「方式 I」フィルタの次数  $n$  と適用 1 回目の最大相対残差 ( $g_s = 1E-13$ )

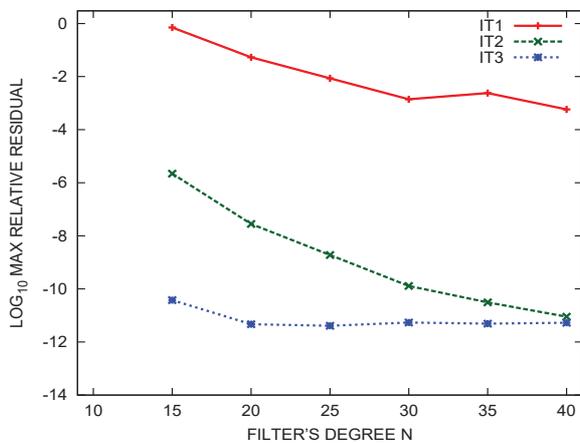


図 A-33 実験 2 : 「方式 II」フィルタの次数  $n$  と最大相対残差 ( $\mu = 1.25$  ( $m = 100$ ),  $g_s = 1E-13$ )

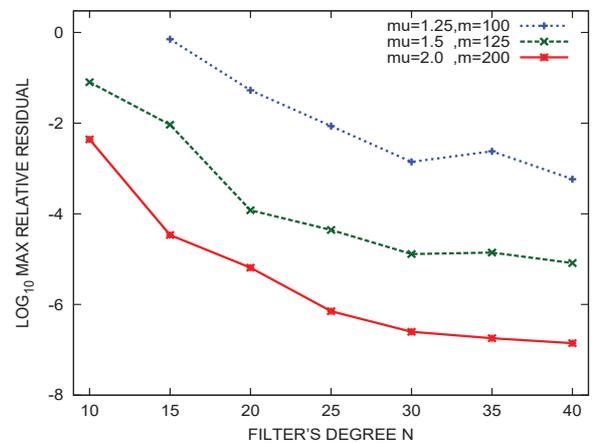


図 A-36 実験 2 : 「方式 II」フィルタの次数  $n$  と適用 1 回目の最大相対残差 ( $g_s = 1E-13$ )

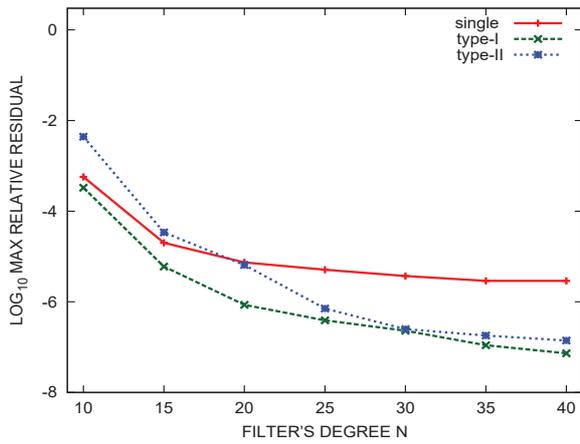


図 A-37 実験 2 : フィルタの次数  $n$  と適用 1 回目の最大相対残差  
 $(\mu = 2.0 (m = 200), g_s = 1E-13)$

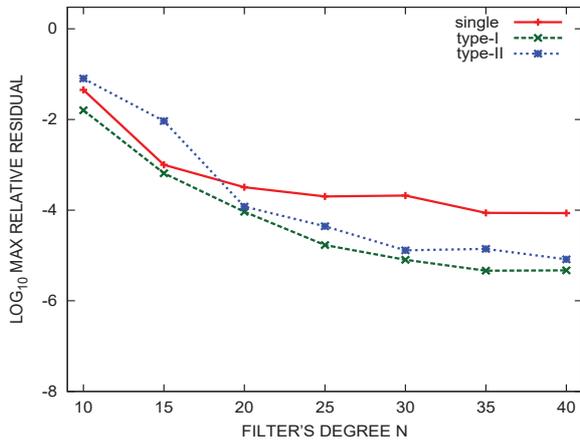


図 A-38 実験 2 : フィルタの次数  $n$  と適用 1 回目の最大相対残差  
 $(\mu = 1.5 (m = 125), g_s = 1E-13)$

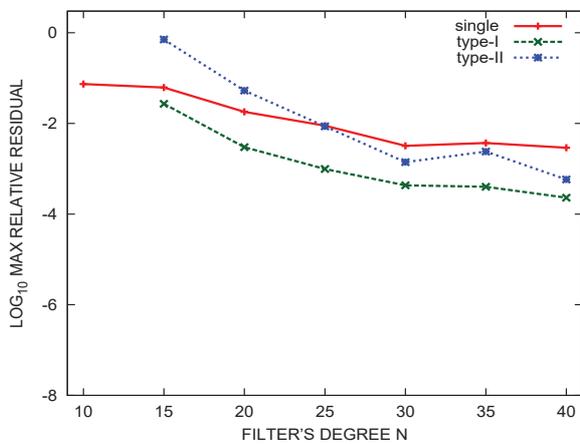


図 A-39 実験 2 : フィルタの次数  $n$  と適用 1 回目の最大相対残差  
 $(\mu = 1.25 (m = 100), g_s = 1E-13)$

参考文献

- [1] 村上弘：レゾルベントの線形結合によるフィルタ対角化法, *情報処理学会論文誌：コンピューティングシステム (ACS)*, Vol.49, No.SIG 2(ACS21), pp.66–87 (2008).
- [2] 村上弘：固有値が指定された区間内にある固有対を解くための対称固有値問題用のフィルタの設計, *情報処理学会論文誌：コンピューティングシステム (ACS)*, Vol.3, No.3(ACS31), pp.1–21 (2010).
- [3] 村上弘：対称一般固有値問題のフィルタ作用素を用いた不変部分空間の近似構成, *情報処理学会論文誌：コンピューティングシステム (ACS)*, Vol.4, No.4 (ACS35), pp.1–14 (2011).
- [4] 村上弘：レゾルベントを用いたフィルタによる固有値問題の解法について, *情報処理学会研究報告*, Vol.2012-HPC-133, No.22, pp.1–8 (2012).
- [5] 村上弘：実対称定値一般固有値問題の最小側固有値を持つ固有対に対する実数シフトのレゾルベントを組み合わせたフィルタによる解法, *先進的計算基盤システムシンポジウム論文集 2012*, pp.81–82 (2012).
- [6] 村上弘：レゾルベントの線形結合をフィルタに用いたエルミート定値一般固有値問題のフィルタ対角化法, *情報処理学会論文誌：コンピューティングシステム (ACS)*, Vol.7, No.1 (ACS45), pp.57–72 (2014).
- [7] 村上弘：レゾルベントの多項式をフィルタとして用いる対角化法について, *情報処理学会研究報告*, Vol.2014-HPC-146, No.13, pp.1–4 (2014).
- [8] 村上弘：実対称定値一般固有値問題に対するレゾルベントの多項式によるフィルタの構成法の検討, *情報処理学会研究報告*, Vol.2014-HPC-147, No.2, pp.1–10 (2014).
- [9] 村上弘：実数シフトのレゾルベントを組み合わせたフィルタによる実対称定値一般固有値問題の下端付近の固有値を持つ固有対の解法, *HPCS2015 シンポジウム論文集*, Vol.2015, pp.38–51 (2015).
- [10] Anthony P. Austin and Lloyd N. Trefethen: "Computing Eigenvalues of Real Symmetric Matrices with Rational Filters in Real Arithmetic", *SIAM J. Sci. Comput.*, Vol.37, No.3, pp.A1365–1387 (2015).
- [11] 村上弘：一つのレゾルベントから構成されたフィルタを用いた実対称定値一般固有値問題に対するフィルタ対角化法の実験, *情報処理学会研究報告*, Vol.2015-HPC-149, No.7, pp.1–16 (2015).
- [12] 村上弘：実対称定値一般固有値問題の最小側固有対を解くための実数シフトのレゾルベントの多項式によるフィルタの簡易な設計法, *情報処理学会研究報告集*, Vol.2016-HPC-155, No.44, pp.1–27 (2016).
- [13] 村上弘：レゾルベントの多項式によるフィルタを用いた実対称定値一般固有値問題の解法, *情報処理学会研究報告集*, Vol.2016-HPC-157, No.4, pp.1–15 (2016).
- [14] 村上弘：チェビシェフ展開形で表わされたレゾルベントの多項式によるフィルタの伝達特性の調整, *数理解析研究所講究録*, No.2019, pp.96–112 (2017).
- [15] 村上弘：実対称定値一般固有値問題を解くためのレゾルベントの多項式型フィルタの設計について, *情報処理学会研究報告*, Vol.2017-HPC-158, No.7, pp.1–10 (2017).
- [16] 村上弘：実対称定値一般固有値問題を解くための少数のレゾルベントの多項式を用いたフィルタの設計法, *情報処理学会研究報告*, Vol.2017-HPC-159, No.4, pp.1–13 (2017).
- [17] 村上弘：少数のレゾルベントから構成されたフィルタを用いた実対称定値一般固有値問題の解法, *情報処理学会研究報告*, Vol.2017-HPC-160, No.32, pp.1–32 (2017).
- [18] 村上弘：少数のレゾルベントで構成された多項式型フィルタによる対称定値一般固有値問題の解法, *情報処理学会研究報告*, Vol.2017-HPC-161, No.7, pp.1–13 (2017).
- [19] 村上弘：少数のレゾルベントから構成されたフィルタを用いた対称定値一般固有値問題の解法, *情報処理学会研究報告*, Vol.2017-HPC-162, No.21, pp.1–34 (2017).
- [20] 村上弘：少数のレゾルベントの多項式型フィルタを用いた一般固有値問題の解法, *情報処理学会研究報告*, Vol.2018-HPC-165, No.15, pp.1–21 (2018).
- [21] 村上弘：フィルタにレゾルベントの線形結合の多項式を用いた複素エルミート定値一般固有値問題の解法, *情報処理学会研究報告*, Vol.2018-HPC-166, No.10, pp.1–17 (2018).
- [22] 村上弘：フィルタ対角化法による近似固有対の精度の改良について, *情報処理学会研究報告*, Vol.2018-HPC-167, No.29, pp.1–31 (2018).
- [23] 村上弘：単一のレゾルベントのチェビシェフ多項式による実対称定値一般固有値問題の解法用の簡易型フィルタ, *情報処理学会論文誌：コンピューティングシステム (ACS)*, Vol.12, No.2 (ACS64), pp.1–26 (2019).
- [24] 村上弘：フィルタ対角化法による固有値問題の近似対の改良, *情報処理学会研究報告*, Vol.2019-HPC-168, No.18, pp.1–36 (2019).
- [25] 村上弘：直交化付きフィルタ適用による固有値問題の近似対の反復改良について, *情報処理学会研究報告*, Vol.2019-HPC-169, No.1, pp.1–31 (2019).
- [26] Hiroshi Murakami: Filters consist of a few resolvents to solve real symmetric-definite generalized eigenproblems, *JJIAM*, Vol.36, No.2, pp.579–618 (2019).
- [27] 村上弘：フィルタの反復適用による実対称定値一般固有値問題の近似対の改良, *情報処理学会論文誌：コンピューティングシステム (ACS)*, Vol.12, No.3(ACS65), pp.14–33 (2019).
- [28] 村上弘：少数のレゾルベントの線形結合の多項式をフィルタとして用いた実対称定値一般固有値問題の解法, *情報処理学会研究報告*, Vol.2019-HPC-171, No.7, pp.1–45 (2019).
- [29] 村上弘：少数のレゾルベントで構成されたフィルタを用いた実対称定値一般固有値問題の解法, *情報処理学会論文誌：コンピューティングシステム (ACS67)*, Vol.13, No.1, pp.1–27 (2020).
- [30] 村上弘：実数シフトのレゾルベントを少数用いて構成された実対称定値一般固有値問題の下端固有対を解くためのフィルタ, *情報処理学会研究報告*, Vol.2020-HPC-173, No.1, pp.1–10 (2020).
- [31] 村上弘：実数シフトのレゾルベント2つで構成されたフィルタによる実対称定値一般固有値問題の下端固有値を持つ固有対の解法, *情報処理学会研究報告*, Vol.2020-HPC-174, No.3, pp.1–17 (2020).