

ランキングデータにおけるクラスタリング手法の提案

金森 弘晃* 松田 眞一†

E-Mail: matsu@nanzan-u.ac.jp

本論文は、ランキングデータのような複数の時系列データに対して新たなクラスタリング手法を提案するものである。sliding window technique を複数時系列用に拡張した上でクラスタリングに必要な距離を考案した。結果として、視覚的に似ていると感じるデータがクラスタを構成し、ランキングデータの細かな性質を見極めることができた。

1 はじめに

時系列データの中からなんらかの特徴を導き出す研究はパターン抽出法として重要なものである。部分時系列クラスタリング (井手 [3], Das et al.[1]) はそのような方法の一つであるが、この方法は1本の長い時系列データを多数の部分時系列に分け、それらをそれぞれ独立なベクトルとしてクラスタリングしようというものである。このクラスタリングによってできたクラスタの平均を代表パターンとして取り出すことが目的である。

しかし、この方法はあくまでも1本の時系列上で連続して似たような反応が起こる場合しか扱っておらず、2本以上の時系列があった場合や連続して反応が起きない場合のクラスタリングまでは扱っていない。

本論文ではこれまで扱われなかった単反応複数時系列のクラスタリング手法を提案する。特にポイント変化のあるランキングデータにおける変動パターンをつかむ手法を考える。したがって、以下では時系列データの単位を便宜上ポイントと呼ぶことにする。

2 アプローチ

基本方針として2つの時系列データ間に適当な距離を定めた上で一般のクラスタリング手法を用いることとする。その距離の算出方法であるが、時系列の長さにはばらつきがある場合も考えたい。例えば、あるデータは18回しかランクインしなかったため18回分のデータしかないが、あるデータは26回分のデータがあるといった場合、時系列の長さには差があるものを比較する新たな方法が必要となるからである。

3 ポイント補正方法

2つの時系列データ間距離を計算する際に、前処理としてそれぞれのデータについてポイントの補正が必要となる場合がある。もしそれぞれの時系列データのポイントに全体的な大小の差があれば、そのままクラスタリングを行うと「高ポイントのクラスタ」「中ポイントのクラスタ」「低ポイントのクラスタ」の3つのクラスタに分かれてしまうのは明らかで、変動パターンをつかむという本論文の目的とは一致しない。補正には平均ポイントや

*南山大学数理情報学部数理科学科

†南山大学数理情報学部情報システム数理学科

最高ポイントなどそれぞれの時系列データに関する何らかの特性値を計算し、その特性値で各回のポイントを割ることによって全ての時系列データのポイントをそろえる方法がある。ポイント補正方法として以下の3種類を考える。

1. 最大値で割る
2. 平均値で割る
3. 上位値の平均で割る

3.1 最大値で割った場合

最大値で割った場合のクラスタ平均ポイントの例を図1・図2に示す。(以下どの場合も前者が元データでプロットしたもので後者が実際に割った後の値でプロットしたものである。また、図中の黒線は全体平均を表し、それ以外はクラスタ平均を表す。)

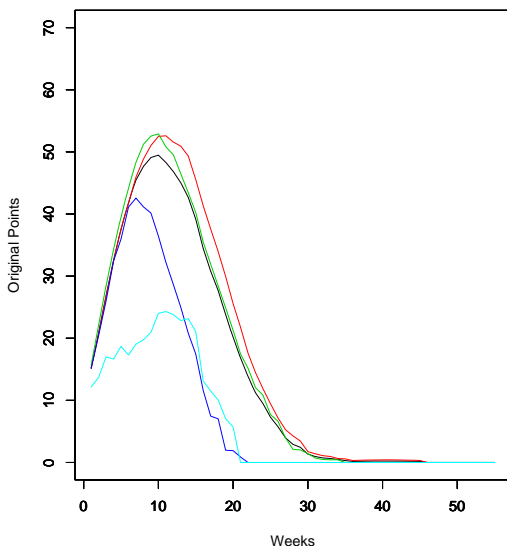


図 1: 最大値で割った場合の平均ポイント
(元データでプロット)

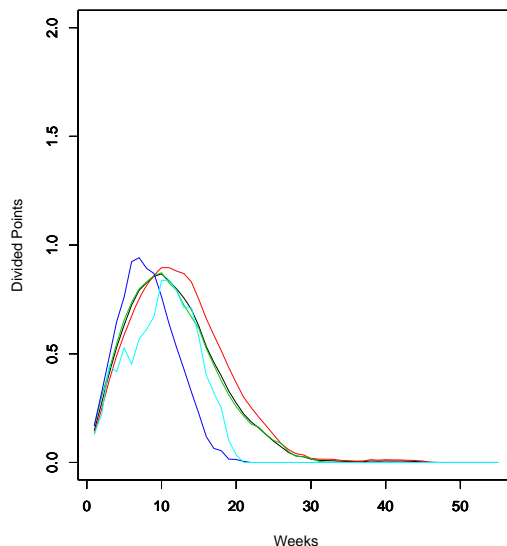


図 2: 最大値で割った場合の平均ポイント
(実際に割った後の値でプロット)

利点は最大値で割ることによってポイントの大小をある程度無視したクラスタリングができることである。図1を見ると、他の方法に比べどのクラスタも最高ポイントがまとまっていることが分かる。一番高いものでも50ポイントを少し超えたぐらいで、明確な「高ポイントのクラスタ」が存在していないことが分かる。

一方欠点は、ピークタイミングによる群分けはされにくいことである。最大値で割った後の値の図2ではピークタイミングに差があり、早いものは7週前後、遅いものでは12週前後にピークを迎えているが、これはピークタイミングによって群分けされたというよりは時系列の長さによって群分けされていると見る方が自然である。なぜならデータ数の少な

いくつかのクラスタを除きどれも逆U型・逆V型といった「ピーク維持期間」や、ランクイン後すぐに最大値を記録するのか、ランクイン後末期に最大値を記録するのかという傾向があまり表れていない。また、図1において上昇スピードがどのクラスタもほぼ同じであり、違いが分かりづらい。他には、最大値という1つの値を使い計算するのでこの値が少し違うだけで全体が大きく左右されやすいという面を持っている。

3.2 平均値で割った場合

平均値で割った場合のクラスタ平均ポイントの例を図3・図4に示す。

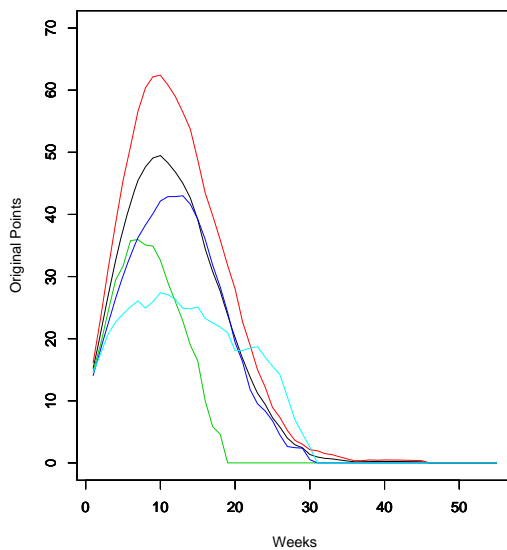


図3: 平均値で割った場合の平均ポイント
(元データでプロット)

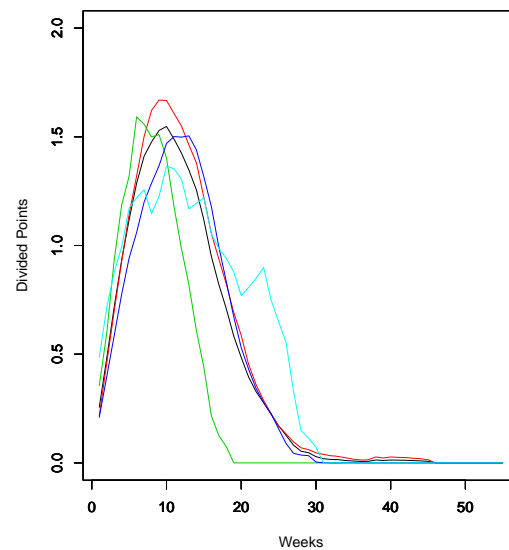


図4: 平均値で割った場合の平均ポイント
(実際に割った後のプロット)

利点はピークタイミングを考慮したクラスタリングができることである。図4を見ると時系列の長さがほぼ同じであるが、ピークタイミングや上昇スピードが異なるクラスタが存在する。このようにそれぞれのデータの平均ポイントが異なるため、割った際の最高ポイントもランク変動パターンによって違いがあり、全てを加味したクラスタリングができる。

一方、欠点はポイントの全体的な大小によって群分けされやすいことである。データの補正後最大値(補正前最大値/平均値)の最高は今回解析したデータでは2.24であるが、補正後最大値の最低は1.13と2倍近い差がある。さらに、補正前の各データの最大値と補正後の各データの最大値の間には0.72という正の相関があることが分かった。これでは2つの時系列データ間のポイント差を補正する方法としては不十分であり、結果的にポイントの大小で分かれてしまうことになる。また、低ポイントが長く続くとその結果に左右されてしまう。ピークポイントが普通であっても低いポイントのままランクインし続ける場合、平均値がとても小さくなり、平均値で割った補正後の最高ポイントが他のデータよりはるかに高くなってしまう場合も出てくる。

3.3 上位値の平均で割った場合

上位値の平均で割った場合のクラスタ平均ポイントの例(ここでは上位5値を用いた)を図5・図6に示す。

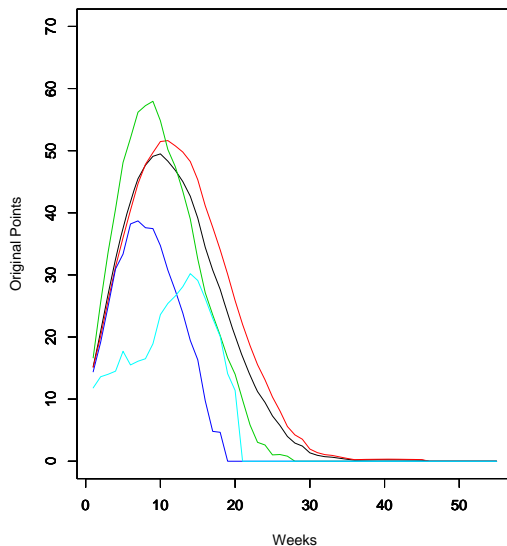


図 5: 上位値の平均で割った場合の平均ポイント(元データでプロット)

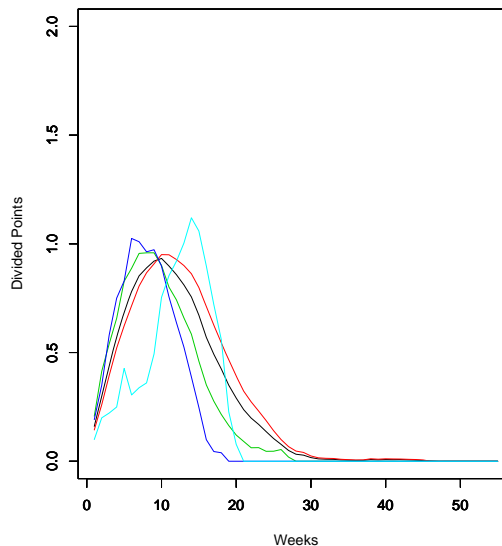


図 6: 上位値の平均で割った場合の平均ポイント(実際に割った後のプロット)

この上位値の平均で割る方法に明確な利点はない。しかし図5を見ると図1ほどではないがある程度ポイントの大小を無視したクラスタリングができていくことがわかる。実際にこの方法による補正前の各データの最大値と補正後の各データの最大値の相関係数は0.14であり、ほとんど影響はないことが分かる。また、図6においてはピークタイミングが同じであるがランクイン期間の異なるクラスタが存在するなど最低限の群分けができていく。1つの最大値やポイントの全体的な大小によって補正後の値が大きく変わってしまうようなこともなく、安定したクラスタリング結果が得られる。

3.4 まとめ

最大値や平均値で割る方法はどちらも利点はあるが欠点が非常に大きい極端な方法であった。その上、補正後の値が一つの最大値によって左右されたり、低ポイントが続くことによって正確な補正ができない例もあった。一方、上位値の平均で割る方法は両者の間をとったようなものではあるが、この方法が最も欠点の少ないものである。今回は上位5値の平均を使ったが、上位のいくつ(上位 x 値)まで使うかによって結果も変わってくる。つまり何を求めたいか、どんなクラスタリングをしたいかによってこの値を変えればこの方法の実用性はさらに高まる。 x を1に近づければ最大値で割る方法に近くなり、逆に x の値を

増やせば平均で割った場合に近くなる。なお、今回のように単反応で変動が激しくないものであれば上位5値ぐらいが妥当であろう。

4 2つの時系列データ間の距離行列

2つの時系列データ間の距離の算出にはまず、井出 [3] の sliding window technique を使っていくつかのパートに分ける必要がある。井出 [3] では対象となる時系列が1つだったため一方のみをスライドさせる方法であったが、ここでは2つの時系列の比較であるため両方のデータをスライドさせる点で違いがある。

次に2つの時系列データ a , b のそれぞれのパート間の距離を求める。sliding window technique の固定窓幅 ω を15とすると a のパートの1つ目 $A_{1r}(r = 1, 2, \dots, \omega)$ は a_1 から a_{15} まで、2つ目 $A_{2r}(r = 1, 2, \dots, \omega)$ は a_2 から a_{16} まで、3つ目 $A_{3r}(r = 1, 2, \dots, \omega)$ は a_3 から a_{17} までという具合に部分時系列を作っていく。最終的には長さ n の1つの時系列データから $n - \omega + 1$ 個の部分時系列が生成されることになる。 a の部分時系列 A と b の部分時系列 B を

$$\begin{cases} A_i = (A_{i1}, \dots, A_{i\omega})' & (i = 1, 2, \dots, n - \omega + 1) \\ B_j = (B_{j1}, \dots, B_{j\omega})' & (j = 1, 2, \dots, m - \omega + 1) \end{cases} \quad (1)$$

とおくことにする。ただし、 n , m はそれぞれ a , b の時系列の長さである。

最後にこの2つの時系列データ間の距離行列を計算するのであるが、それは部分時系列同士のユークリッド距離によって求めることとする。すなわち、 $d(x, y)$ を同じ長さの2つのベクトル x , y 間のユークリッド距離として、2つの時系列データ間の距離行列 $X = \{X_{ij}\}$ を以下のように定義する。

$$X_{ij} = d(A_i, B_j) \quad (2)$$

5 距離行列の意味

距離行列 X は時系列データ間の全ての一部分どうしの距離をあらわしたものである。このような操作が必要である理由はデータによってピークのタイミングが違うからである。例として図7を挙げる。縦軸はポイント、横軸は週単位の時間である。この2つの時系列データは最初の上昇スピードは違うがその後の推移はかなり似ている。これはただ単に急上昇した時期が少しずれていただけと考えるのが自然であり、データ間の距離は近くなるのが理想と考えられる。そこでこの距離行列は全ての一部分どうしの距離を計算しているため、2つの時系列データ間で最も距離の近くなった週どうしを選んで総合的な距離として適用することができる。

6 2つの時系列データ間距離

2つの時系列データ間距離 D は以下の式で定義する。

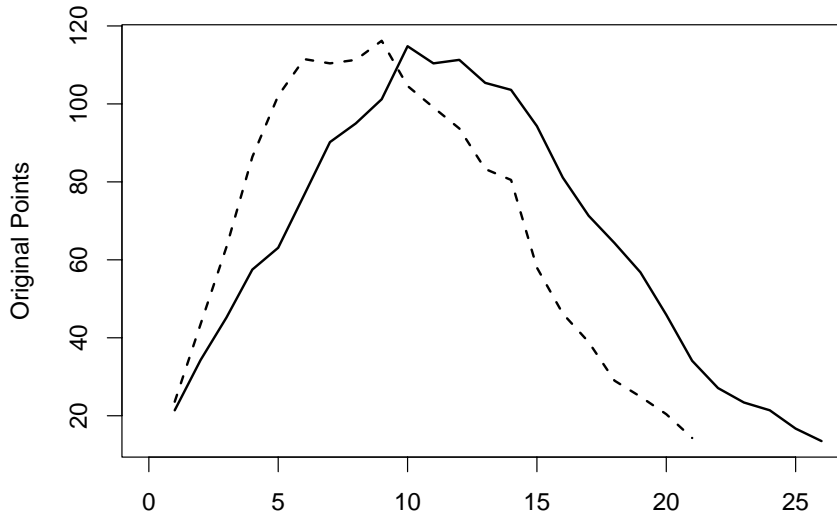


図 7: タイミングの差

$$D = \frac{1}{2} \left(\frac{\sum_{k=1}^{n-\omega+1} \min X_{kj}}{n-\omega+1} + \frac{\sum_{k=1}^{m-\omega+1} \min X_{ik}}{m-\omega+1} \right) \quad (3)$$

すなわち、まずそれぞれの行方向について最小値をとり、その最小値の平均をとる。同様にそれぞれの列方向についても最小値をとり、その平均をとる。最後に行方向の平均値と列方向の平均値の平均を2つの時系列データ間距離とする。

前節でずれを考慮して一番近くなる部分時系列を選ぶと述べたが、単純に距離行列の全体の最小値ではうまくいかなかった。全体の最小値は細かな凹凸の偶然の一致を表すことが多いのが原因でないかと思われる。そのため上記のようにずれを考慮して最小値を探すものの全体としての曲線の当てはまりも考慮してその平均をとるという方法を考案した。(距離行列の全平均では全然当てはまっていない部分時系列どうしも考慮されることになり、うまくはいかない。) 行方向の平均と列方向の平均でさらに平均を取る理由は対称性を保つためである。

7 クラスタリング

前節の距離を元にクラスタリングを行う。その際最長距離法を用いてクラスター分析をすることにする。(クラスター生成方法は様々あるが、実例で問題が生じなかったので本論文ではその点は追求しなかった。可能性としてはウォード法の方が向いている場合もあると思われる。)

8 固定窓幅 ω

固定窓幅 ω は一般的にウィンドウサイズとも呼ばれているもので、長い時系列データをどのくらいの幅で分けていくのかを決めるものである。井手ら [2] では固定窓幅 ω の設定について「変化点の発生間隔が ω より小さいと探知は難しいことに注意する必要がある。」と述べている。本論文では単反応を考えているので、変化点の発生間隔とは時系列の長さそのものであり、固定窓幅 ω より短い時系列データを解析に用いないため今回の手法では問題がない。結論として固定窓幅 ω が長すぎると解析可能なデータ数が減ってしまうことから、できる限り短い幅が望ましい。しかし、固定窓幅が短ければ良いというわけでもない。表 1 はある時系列データ a (時系列の長さ 40) と、時系列データ b (時系列の長さ 32) の距離行列の一部を固定窓幅 ω 別に並べてみたものである。

$\omega = 10$	$\omega = 11$	$\omega = 12$	$\omega = 13$	$\omega = 14$	$\omega = 15$
1.28	1.46	1.64	1.84	2.08	2.36
1.41	1.57	1.75	1.97	2.24	2.42
1.53	1.69	1.90	2.14	2.32	2.47
1.66	1.86	2.09	2.25	2.38	2.51
1.85	2.07	2.21	2.33	2.45	2.52
2.07	2.20	2.31	2.42	2.48	2.55
2.20	2.31	2.41	2.46	2.53	2.54
2.31	2.41	2.46	2.51	2.51	2.51
⋮	⋮	⋮	⋮	⋮	⋮
⋮	⋮	⋮	⋮	⋮	⋮
1.78	1.80	1.82	1.83	1.84	1.84
1.70	1.72	1.73	1.74	1.74	1.74
1.60	1.62	1.63	1.63	1.63	1.63
1.49	1.50	1.51	1.52	1.52	1.52
1.34	1.38	1.39	1.40	1.40	1.40
1.20	1.24	1.26	1.26	1.26	1.27
1.06	1.09	1.11	1.12	1.12	1.13
0.88	0.92	0.94	0.95	0.95	0.96

表 1: 固定窓幅 ω の違いによる距離行列の変化

表の上部は a の序盤と b の終盤との距離で、表の下部は a の終盤と b の終盤との距離である。本来ならば時系列データ同士の序盤は序盤と、終盤は終盤との距離が近くなければならないのであるが、 $\omega = 10$ の場合序盤と終盤の距離が 1.28 という近い距離になってしまっていることが分かる。この傾向は時系列データが長いほど顕著で、行方向と列方向の最小値を 2 つの時系列データ間距離算出のもとにしている今回の手法ではできる限り避けたい現象である。そこでこの例以外でも固定窓幅 ω の値を変えて実験してみたところ $\omega = 15$ あたりでどの時系列データ間も安定してきたため、今回の解析では固定窓幅を 15

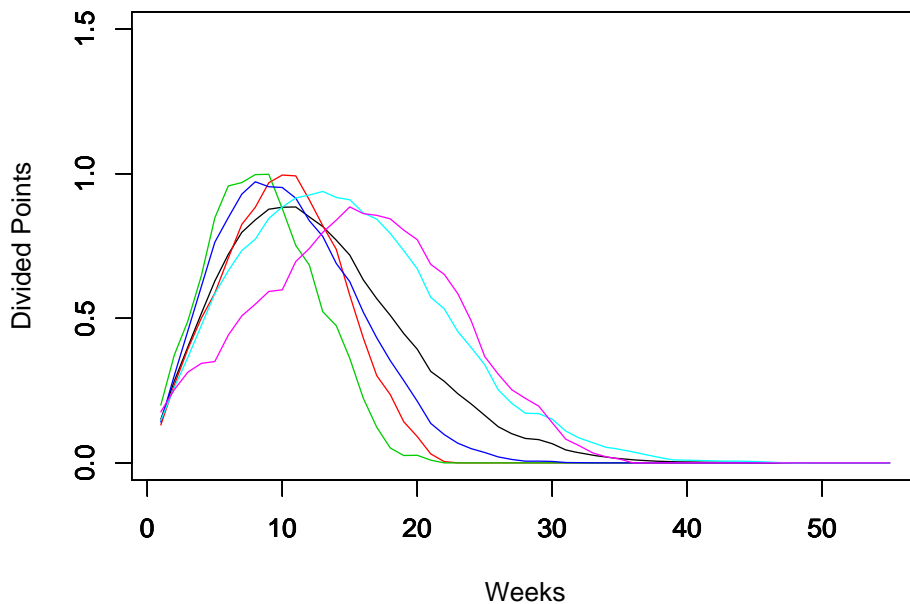


図 9: 2003 年クラスタ平均

- 第3クラスタ（赤線）
このクラスタも第1クラスタ同様に全ての曲が22週以下とランクイン期間が短い。しかし、ランクイン期間が短いとどうしても早くなりがちなヒットタイミングがこの第3クラスタに関しては特にそういった特徴はない。むしろこのクラスタはピーク後の下降スピードが急激で目立っている。
- 第4クラスタ（水色線）
このクラスタは20週から最高は46週まで幅広い曲が属すロングヒットのクラスタである。ロングヒットであるがどの曲も序盤に上昇し、昇降を繰り返すのではなくわりとしっかりしたヒットの山を持ち、緩やかに下降していく曲のクラスタである。
- 第5クラスタ（ピンク線）
このクラスタの大きな特徴は序盤のポイント上昇が他のどのクラスタよりも緩やかなことである。曲のピークが20週前後の曲が多く、それまではわりと激しく昇降を繰り返しながら上昇していく。上昇が緩やかな反面、ピーク後は一気に下降しランク外になる曲ばかりである。

9.1 まとめ

今回の例では大きく5つのクラスタに分けることができた。第1, 第2, 第3クラスタは良く似ていて違いが分かりづらいが、単にランクイン期間が短い第1クラスタの「短期間

ヒット群」に対し、第2クラスはランクイン期間が短いわけではなかったことから「ヒットタイミング早期群」と呼べる。また、第3クラスはランクイン期間が短い曲群にしてはヒットタイミングが少し遅めで、下降が急激なため「ヒットタイミング末期群」とそれぞれ現行の手法でとても細かい特徴まで読み取れていることが分かった。他にも第4、第5クラスのようにどちらも30週前後のロングヒット曲が多い中でヒットタイミングによるクラスタリングができていた。第4クラスはヒットがしっかりしていた「ロングヒット群」であり、第5クラスは上昇スピードにおいては緩やかであるが一気に下降する「下降急激群」であった。

この他の年についても同様に解析を行ったところ、いずれの年もまず「短期間ヒット群」のような変動スピードの速いクラスが大きく分かれた。逆にランクイン期間の長い曲についてはそれほどまとまったクラスにはならなかった。そのため、図9のようなクラス平均の曲線を代表と呼ぶ場合には注意を要する。すなわち、今回のように $\omega = 15$ の場合は15からその2倍の30くらいまでをランクイン期間でうまくクラスタリングする結果となり、序盤にピークを迎えるか、終盤にピークを迎えるかといった特性もそれぞれのクラスに現れていたが、長めの上昇や下降に対する特性にはばらつきが大きかった。したがって、クラス平均からピークタイミング付近の特徴とすその方の特徴を抽出する場合、その安定性に違いがあることを理解することが重要であり、個別のデータのプロットにも注意すべきである。

10 おわりに

本論文の目的はランキングデータにおける変動パターンをつかむクラスタリング手法の提案であった。変動パターンという点に関しては解析結果にも反映され良い手法ができていた。ランクイン回数やピークタイミングが違っていても似ている変動をしたものについてはまとめることができ、「視覚化した場合に似ているものをまとめる」という大きな目標は満たしていた。

参考文献

- [1] G. Das, K.-I. Lin, H. Mannila, G. Renganathan, and P. Smyth. (1998): Rule discovery from time series, Proc. the Fourth ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, pp.16-22.
- [2] 井手剛, 井上恵介 (2004) : 非線形変換を利用した時系列データからの知識発見, 第4回データマイニングワークショップ, 日本ソフトウェア科学会データマイニング研究会, 研究会資料シリーズ ISSN 1341-870X, No.29, pp.1-8.
(<http://spinglass.hp.infoseek.co.jp/>)
- [3] 井手剛 (2006) : 部分時系列クラスタリングの理論的基礎, 第20回人工知能学会全国大会予稿集, 2A1-2. (<http://spinglass.hp.infoseek.co.jp/>)
- [4] Top Hits Online : <http://www.tophitsonline.com/>.