

発見ルールフィルタリングシステムの予備評価

A Preliminary Evaluation of Discovered Rule Filtering System

飯田 暁¹, 北村 泰彦², 朴 勤植³, 辰巳 昭治¹
Akira Iida, Yasuhiko Kitamura, Keunsik Park, Shoji Tatsumi

¹ 大阪市立大学大学院工学研究科 ² 関西学院大学理工学部 ³ 大阪市立大学大学院医学研究科

¹ Graduate School of Engineering, Osaka City University

² School of Science and Technology, Kwansai Gakuin University

³ Graduate School of Medicine, Osaka City University

Abstract: A data mining system semi-automatically discover knowledge by mining a large volume of data, but the discovered knowledge is not always novel to the user. We discuss a discovered rule filtering method to filter rules discovered by a data mining system into ones that are novel to the user by using information retrieval results from the Internet. We have two approaches; the micro view and the macro view approaches, to achieve discovered rule filtering. In the micro view approach, we extract keywords from a discovered rule and rank the rule referring to the number of hits when we submit the keywords to the MEDLINE database. Unfortunately, this approach does not work well because the retrieved result contains a number of noisy documents and because the number of hits depends much on the number of keywords but on the novelty of rule. In the macro view approach, we first retrieve documents by submitting every pair of the extracted keywords and then make the results gather in clusters. As a preliminary evaluation of the macro view approach, we made a survey on a relation between a set of outputs of the system and judgments made by medical experts. We found a tendency that keywords in a known rule form a single cluster.

1 はじめに

アクティブマイニングは情報収集、データマイニング、ユーザリアクションの技術を融合することにより、ユーザの目的にあった質の高い知識の効率的な発見を目指すデータマイニングの新しいアプローチである[1]。本稿ではインターネット上からの文献情報検索結果に基づきデータマイニング結果をフィルタリングする発見ルールフィルタリング手法[4,5]とその予備評価について述べる。

データマイニングは大量のデータを機械処理することにより、ユーザにとって有用な知識を自動的に発見しようとする手法である。一般的には与えられたデータに含まれる属性間の関係から統計的に意味のある関係を発見する手法がとられている。しかしながら単に統計的な特徴だけでデータマイニングを行うなら、(1)ユーザが扱いきれないほどの大量の知識が得られる、(2)ユーザにとって既知の知識が得られる、といった問題が生じる。そこで本研究ではルール形式で得られる大量の発見知識の中から新規なもののみをフィルタリングする発見ルールフィルタリングの開発を行っている。この発見ルールフィルタリングを実現するには発見されたルールが新規かどうかを判定する必要があるが、

これは当然のことながらマイニングの対象となる情報源以外の外部情報源を参照する必要がある。そこで、本研究ではインターネット上の情報源を用いて、発見ルールの新規性判定を行おうとしている。

本論文では以下、2章において発見ルールフィルタリングの概要について述べる。発見ルールフィルタリングにはマイクロビューアプローチとマクロビューアプローチの二つが考えられるが、3章ではマイクロビューアプローチとその予備評価について述べる。4章ではマクロビューアプローチとその予備評価について述べる。最後に、5章でまとめと今後の展望について述べる。

2 発見ルールフィルタリング

発見ルールフィルタリング[4,5]ではデータマイニングシステムにより発見されたルール形式の知識に対して、それに関連する情報をインターネット上から検索し、その結果に基づき、発見ルールのフィルタリングを行う。ここでの議論をより具体的なものとするために、肝炎データからのマイニングを知識発見の対象領域として議論を進める。発見ルールフィルタリングの具体的な手順は以下の通りである。

2.1 発見ルールの獲得

データマイニングシステムを利用してルール形式の知識を得る。静岡大学の山口グループでは、千葉大学医学部より提供された肝炎データから肝炎の進行具合を示す血液データ(GPT)と他の検査データとの時系列的な相関関係に着目し、様々なルール形式の知識発見を行っている[2]。発見されたルールの一例を図1に示す。

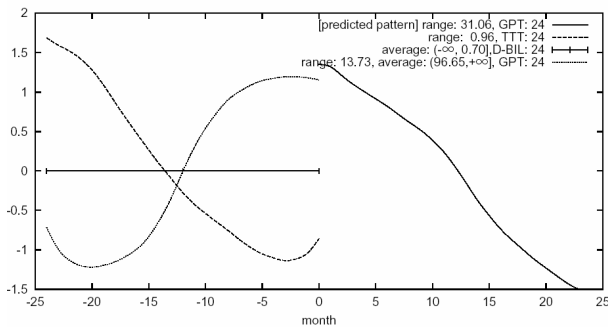


図1：発見ルールの一例

これはGPT(グルタミン酸ピルビン酸トランスアミナーゼ)とTTT(チモール混濁試験), D-BIL(直接ビリルビン)の関係を図的に表すものであるが、「24ヶ月間、D-BILが一定で、TTTが減少し、GPTが増加しているなら、その後24ヶ月間のGPTは減少する」というルール表記を行うことも可能である。

2.2 発見ルール駆動情報検索

発見されたルールに関連する情報を、インターネットを利用して収集する。肝炎に関連する一般のWeb情報源はあまりにも雑多な情報が含まれているので、本研究では医学・生物学関係の文献データベースであるMEDLINEを用いている。MEDLINE(MEDlars on LINE)は、米国をはじめ70カ国で出版された4000誌を超える医学・生物学系学術雑誌からのアブストラクトを含む書誌情報データベースであり、1966年以降の1100万件以上のデータが蓄積されている。PubMed [3] (<http://www.ncbi.nlm.nih.gov/entrez/query.fcgi>)はNCBI(National Center for Biotechnology Information)によりインターネット上に無料提供されているMEDLINEの検索サービスであり、一般の検索エンジン同様、キーワードを入力することにより、MEDLINE文献の検索を行うことができる。

発見ルールに関連する情報をMEDLINEデータベースから検索するためには、それに関連するキーワードを獲得する必要がある。このキーワードは発見ルール、マイニング領域から導出している。例えば、図1の例ではGPT, TTT, D-BILが発見ルールより抽出される。

またマイニング領域は肝炎であるので、hepatitisが領域キーワードとなる。これらのキーワードを組み合わせたものがMEDLINEデータベースの検索に用いられる。

図1で示すルール自体は各検査値の変動に関する情報も含んでいるが、現時点ではそのような情報を文献検索に反映することは困難であり、これは今後の課題となる。

2.3 発見ルールフィルタリング

MEDLINE文献検索の結果に応じて発見ルールのフィルタリングを行う。その手法としては、文献検索結果に応じてルールのランキングを行ったり、ある閾値を設定して、それ以下のものを排除したりすることである。文献検索結果とランキングをどのように対応付けるかは本研究の中核となる重要な研究課題であるが、マイクロビューとマクロビューの二つのアプローチが考えられる[4]。次章以下ではそれぞれのアプローチの詳細と予備評価について述べる。

3 ミクロビューアプローチとその評価

3.1 ミクロビューアプローチの仮説

マイクロビューアプローチは文献検索のヒット数に応じて発見ルールのランク付けをするアプローチである。マイクロビューアプローチでは以下の仮定をおいた。

【仮定】(マイクロビューアプローチ)

1. 既知のルールのヒット数は多い。
2. 未知のルールのヒット数は少ない。
3. ゴミのルールのヒット数はゼロである。

この仮定では既知のルールはよく研究がなされており、それに関連する文献は多く出版されていることを前提としている。ヒット数が少なければ、それは十分な研究がなされているとは限らないので、未知ルールに分類することができると考えられる。一方で、ルールが無意味なものに関しては誰も研究しないのでゼロであると仮定している。

3.2 仮説の検証と考察

以上の仮説の妥当性を検討するために、医学専門家へのアンケートに基づく予備評価を行った。

アンケートの質問文を図2に示す。問1の質問項目には、静岡大学のグループにより得られた30個のルールの中からランダムに20個選択してそこに含まれる属性名を使用した。アンケートの対象は、医師国家試験を直前に控えた医学部の6回生、47人である。国家試験を直前にひかえた医学生であることから、医学的

問1：下のキーワードの組み合わせでMEDLINEを検索するとしたときに、検索結果がどのようになると思いますか？ 次の三つの中から選んでください。

1. 当然（どのような内容の文献が見つかるか、ある程度予想できる）
2. おもしろい（自分が知らないことが載っている文献が見つかるかもしれない）
3. 無意味（ほとんどヒットしないだろう）

- | | | | |
|---------|----------|------|------------------------|
| (1) ALT | TTT | | (1.当然、2. おもしろい、3. 無意味) |
| (2) TTT | 直接ビリルビン | ALT | (1.当然、2. おもしろい、3. 無意味) |
| (3) ALT | 総コレステロール | C型肝炎 | (1.当然、2. おもしろい、3. 無意味) |
| ... | | | |

問2：次の検査項目同士の関係について、次の四つの中から適切だと思うものを選んでください。

1. 強い関係がある.
2. 関係がある.
3. あまり関係は無い.
4. 関係ない.

- | | | | |
|-----|-----|----------|----------------|
| (1) | ALT | 総ビリルビン | (1、 2、 3、 4) |
| (2) | ALT | 総コレステロール | (1、 2、 3、 4) |
| (3) | ALB | 総コレステロール | (1、 2、 3、 4) |
| ... | | | |

図2：アンケートの質問文

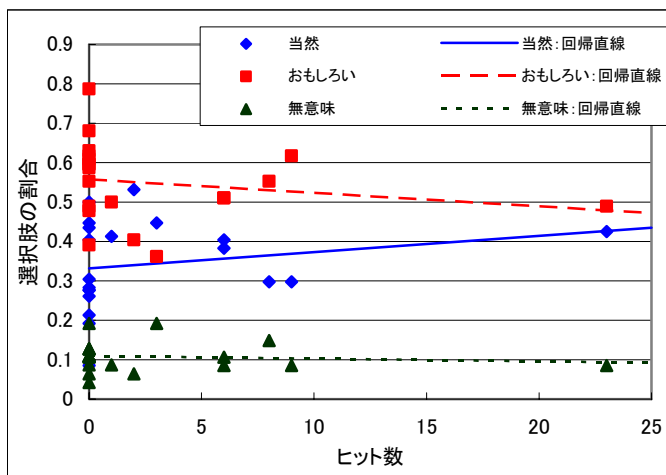


図3：選択肢の割合とヒット数の関係

な常識とも言える教科的な知識についてはかなり詳しいと考えられる。

図3に問1のアンケート結果と文献検索結果をまとめたものを示す。それぞれの発見ルールに対して、三つの選択肢が選ばれた割合と関連文献検索したときのヒット数の関係をプロットしたものである。またそれぞれの選択肢について求めた回帰直線も表示している。それぞれの選択肢について、その回帰直線の傾きを有

意水準95%でt検定[6]した。自由度18（データ数20から2を引いたもの）での2.5%点の値（2.101）を、それぞれのt値の絶対値が上回っていれば、ヒット数は選択肢の割合に影響を与えていると言える。

それぞれの選択肢についてt値を求めると、「当然」では0.887、「おもしろい」では-0.812、「無意味」では-0.441となり、アンケート結果とヒット数の間には有意な相関を見出すことはできなかった。

マイクロビューアアプローチによる検索結果と専門家との評価の間に相関が得られなかった理由としては、検索ヒット数は検索に用いるキーワード数に依存することが考えられる。図4はキーワード数とヒット数の関係と、そのキーワード数ごとの平均値をプロットしたものである。この図より、キーワード数が増えると、ヒット数は減少する傾向にあることがわかる。よって、文献検索ヒット数は、キーワード数に依存しているために、抽出されたキーワードをそのまま文献検索に用いる手法は有効でないことが明らかになった。

ただし、キーワード数が少ない場合は専門家の評価とヒット数に相関があることが予想できる。そこで、図2の問2に示されるアンケート調査を実施した。問2は、発見ルールに含まれるキーワード対の組み合わせ

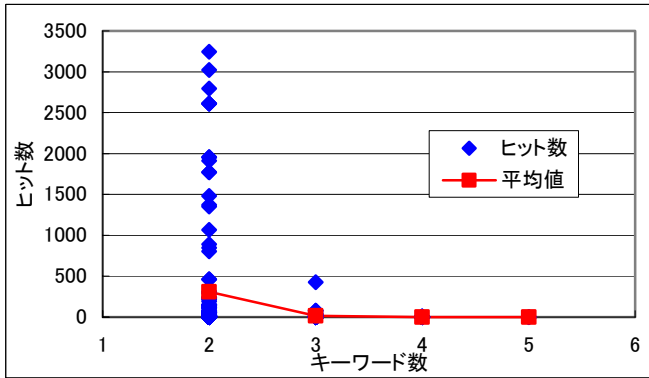


図 4 : キーワード数とヒット数の関係

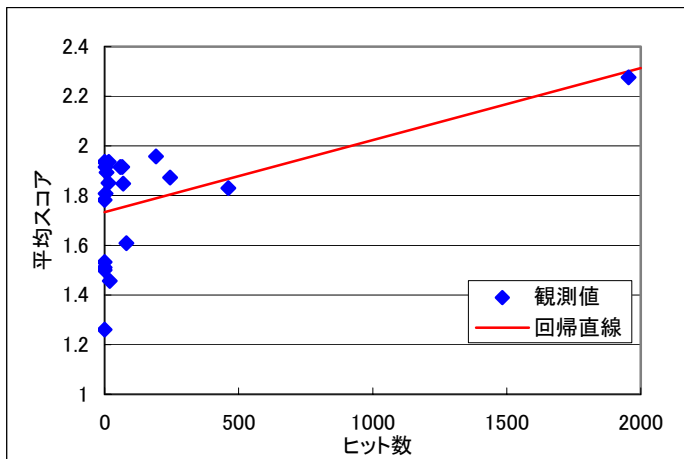


図 5 : ヒット数と平均スコアの関係

の中から 20 組を選択したものを用いた。

図 5 に問 2 の結果を示す。「強い関係がある」、「関係がある」、「あまり関係はない」、「無関係」のそれぞれに 3, 2, 1, 0 のスコアをつけ、その平均と MEDLINE 検索によるヒット数をプロットしたものである。平均スコアとヒット数との間の相関係数は 0.54 であり、ヒット数と専門家の評価には相関があることがわかった。したがって文献検索に用いるキーワード数を 2 に限定すれば、検索結果と専門家の評価には相関がある。この結果は次章のマクロレビューアプローチへつながってゆく。

4 マクロレビューアプローチとその評価

4.1 マクロレビューアプローチの仮説

マイクロレビューアプローチはキーワード検索を用いて発見ルールに直接関連する文献を検索しようとするアプローチであったが、現状では十分なフィルタリング能力を発揮することはできなかった。そこでマクロレビューアプローチでは発見ルールから抽出されるキーワードの全ての対に関する共起文献数を参照することで、発見ルールに関連する研究動向を推測しようとするものである。

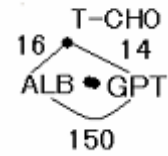


図 6 : GPT, ALB, T-CHO 間の共起関係

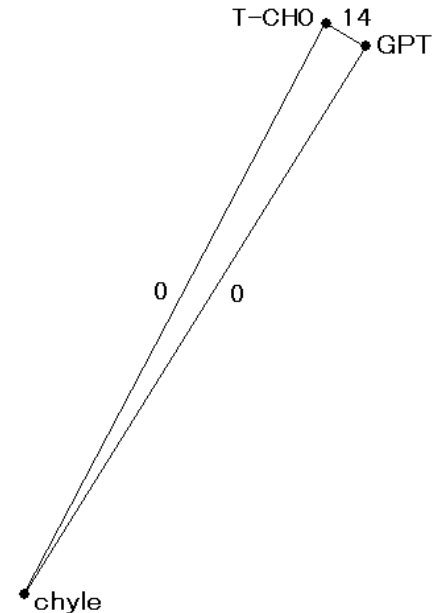


図 7 : GPT, T-CHO, chyle 間の共起関係

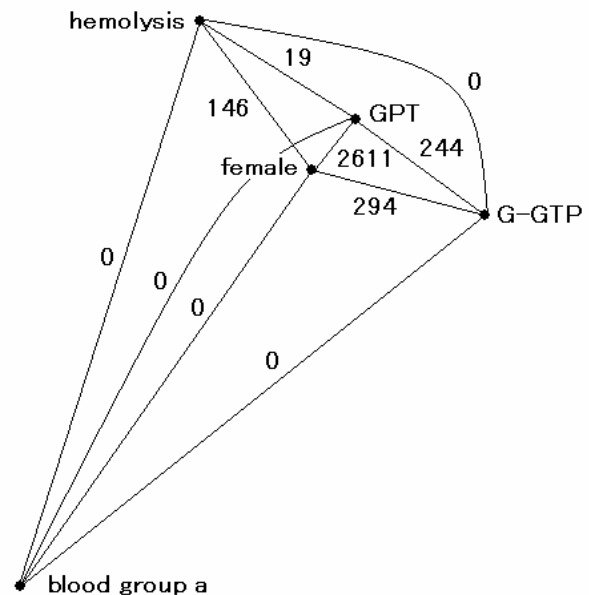


図 8 : GPT, G-GTP, hemolysis, blood group a, female 間の共起関係

図 6, 7, 8 は 3 つのルールでの共起文献数をグラフとして表したものである。グラフ中でのノードはそれぞれのキーワードを表し、エッジの長さは二つのキーワードでの共起数の逆数、添字は共起数を表す。

図 6 のルールでは、ALB, GPT, T-CHO ともにそれ

ぞれ関係があるといえる。図7のルールでは、T-CHO、GPTの間には関係があるが chyle との間には関係がないことが分かる。図8のルールでは、GPT、female、G-GTP間には相互に関係があるが、hemolysisとG-GTP間と blood group a とその他のキーワードとの間には関係がないことが分かる。

それぞれのルールを明確に特徴付けるために、階層クラスタ分析法[7]を用いてキーワードのクラスタリングを行った。階層クラスタ分析とは、グラフ中のノードをいくつかのグループに分けることであるが、そのクラスタを生成する手法としては共起数の逆数を距離として、最長距離法 (CLINK: complete linkage clustering method) を用いた。この手法はまず、最も距離の近いノード同士を統合しクラスタを生成し、新たに統合されたクラスタと他のクラスタとの距離にはそれぞれのクラスタに属するノードのうちで最も遠いもの同士の距離として定め、ある閾値以内にあるノードはクラスタに加えてゆく。こうしてできたそれぞれのクラスタに含まれるキーワード同士は関係が深く、それらに関する知識はすでに知られていると見なすことができる。そこから、生成されたクラスタの数により、ルールを以下のように分類できると仮定する。

【仮定】(マクロビューアアプローチ)

1. 既知のルールのクラスタ数は1である。
2. 未知のルールのクラスタ数は2である。
3. ゴミのルールのクラスタ数は3以上である。

クラスタ数が1の場合は、すべてのキーワード対に関する研究が活発な場合であるので、そのようなキーワードを含むルールは既知であると考えられる。クラスタ数が2の場合は、キーワード集合が二つのクラスタにグループ化される場合である。それぞれのグループに関する研究は行われていても、グループ間の研究は十分ではないと考えられるので、それに関するルールは未知であると判定できる。しかし、クラスタ数が3以上の場合は、ゴミルールとした。これは未知の関係があるグループが複数存在する場合には、ルールとして複雑すぎると考えられるからである。

例えば、クラスタリングの閾値を距離1(キーワード間の共起数が1件のもの)とすると、図6のルールでは、全てのキーワードが1つのクラスタに統合さる。図7のルールでは、GPTとT-CHOは統合されるが chyle は統合されず、クラスタ数は2となる。図8のルールではGPT、G-GTPとfemaleが一つのクラスタに統合され、hemolysisと blood group a は統合されないで、クラスタ数は3となる。

4.2 仮説の検証と考察

以上の仮説の妥当性を検討するために、3章で行ったアンケート調査(問1)の結果を用いた。表1にその結果を集計したものを示し、図9にクラスタ数と選択肢の割合の関係をグラフにして示す。それぞれの選択肢が選ばれた割合と、階層クラスタ分析の結果のクラスタ数との関係をプロットしたものである。なお、クラスタリングの閾値は、距離1とした。

仮説の検証を行うために、 χ^2 検定[8]を用いて、クラスタ数の変化に伴い、それぞれの選択肢の割合に有意な変化があるかを検証した。有意水準としては95%を用いた。

その結果、「当然」の割合は、クラスタ数1と2、1と3の間に有意な差があることが認められた。「おもしろい」の割合にも、クラスタ数1と2、1と3の間に有意な差があることが認められた。「無意味」の割合は、クラスタ数の変化による差は認められなかった。よって、クラスタ数が1から複数になると「当然」の割合は減少し、「おもしろい」の割合は増加する、と言える。また、「無意味」の割合はクラスタ数の変化と関係があるとは言えない。

このことから、クラスタ数が1のルールは「当然」なものであり、複数個あるルールは「おもしろい」ものであると考えられる。

また、与えられたルールに「無意味」と回答する割合は少ない。したがって、ルールをフィルタリングしてしまうよりは、クラスタ数の多い順にランキングする方が適切であると考えられる。

表1：問1の結果

クラスタ数	当然	おもしろい	無意味	無回答	計
1(10例)	189	233	46	2	470
2(8例)	110	219	43	4	376
3(2例)	22	60	12	0	94
計	321	512	101	6	940

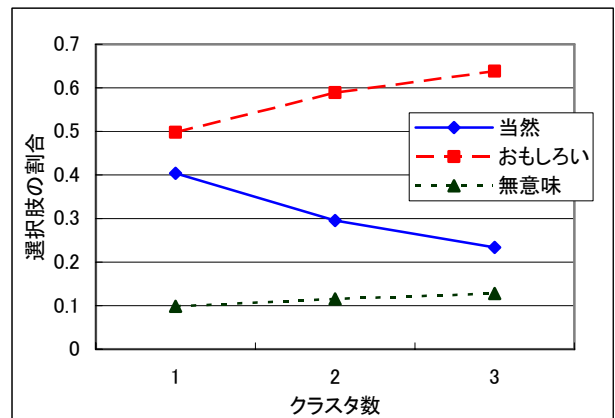


図9：選択肢の割合とクラスタ数の関係

なお、今回の調査ではクラスタ数は最大で3であるので、クラスタ数が4以上のルールに対する検証は今後の課題である。

5 まとめ

MEDLINE 情報検索の結果を用いてデータマイニングにより得られた発見ルールを、利用者にとって未知なものにフィルタリングする発見ルールフィルタリングについて述べ、二つのアプローチとその予備評価について述べた。マイクロビューアプローチで行った文献検索は単純なキーワード検索に基づく手法であり、十分な精度を得るようなフィルタリングを行うことが困難であった。一方、共起関係を用いたキーワードクラスタリングによるマクロビューアプローチによる結果は医学専門家との判断とも相関があることがわかった。ただし、医学専門家は発見ルールに対して「無意味」のラベルを貼ることに慎重であるので、発見ルールをフィルタリングしてしまうことは問題があるかもしれない。むしろ、発見ルールをその重要性に応じてランキングする手法をとることが望ましい。

今後の課題としては以下のものが挙げられる。

(1)情報検索の精度の向上。マクロビューアプローチにおいても、キーワードの共起関係を求めるために検索を行っているが、その精度は必ずしも高いとはいえない。文献検索の精度の向上を行えば、さらに正確なフィルタリングが可能になると考えられる。一つの手法は得られた文献のアブストラクトを自然言語処理の手法を用いて解析することである[9]。例えば、発見ルールに含まれる属性が文献アブストラクト中で離れた場所に存在するとするならばそれらの属性の関連を述べた文献である可能性は低いかもしれない。したがってアブストラクトの構文解析を行うことにより、属性キーワードが同一文中に現れるかどうかを確認できればフィルタリングの精度は向上すると考えられる。さらに属性キーワードが結果や結論に関連する文の中に現れるかどうか、属性キーワード間の関係を修飾する文節は何か、などということが明らかになればフィルタリングの精度はさらに向上すると考えられる。

また先にも述べたが本研究での文献検索は単純なキーワード検索によっているため、属性間の関係や時系列的な変化を反映したものではない。これに関しても自然言語処理の手法を応用することで対処できるかもしれない。

(2)発見ルールフィルタリングの効果の実証。開発した手法を肝炎データマイニング分野に応用し、利用者の支援にどの程度役に立つかを明らかにする。これに関しては静岡大学などで行われている肝炎データマイニングの結果を用いてその検証を行う予定である。

謝辞:本研究にあたり、肝炎データベースからの発見ルールを提供していただいた静岡大学山口高平教授、自然言語処理に関して議論いただいた奈良先端大学院大学の松本裕二教授に感謝の意を表します。

参 考 文 献

- [1] H. Motoda (Ed.), *Active Mining: New Directions of Data Mining*, IOS Press, Amsterdam, 2002.
- [2] 畑澤寛光, 佐藤芳紀, 山口高平, 波形パターンを分類クラスとするルールの発見支援システムの構成法—慢性肝炎データセットを対象にして—, 人工知能学会研究会資料, SIG-KBS-A201-11, 2002.
- [3] 懸俊彦, *PubMed 活用マニュアル*, 南江堂, 東京, 2000.
- [4] Y. Kitamura, K. Park, A. Iida, and S. Tatsumi. *Discovered Rule Filtering Using Information Retrieval Technique*. *Proceedings of International Workshop on Active Mining*, pp. 80-84, 2002.
- [5] Y. Kitamura, A. Iida, K. Park, and S. Tatsumi. *Micro View and Macro View Approaches to Discovered Rule Filtering*. *Proceedings of 2nd International Workshop on Active Mining*, pp.14-21, 2003.
- [6] 岸野洋久, *社会現象の統計学*, 朝倉出版, 1992.
- [7] 岡太彬訓, 今泉忠, *パソコン多次元尺度構成法*, 共立出版, 1994.
- [8] 津村善郎, 淵脇学, 築村昭明, *社会統計入門[第2版]—経済学を学ぶ人のために—*, 東京大学出版会, 1988.
- [9] T. Yamasaki, M. Shimbo, and Y. Matsumoto: *A MEDLINE document search system using section information*, 人工知能学会研究会資料, SIG-KBS-A301-05, 2003.