

平成 22 年 5 月 20 日現在

研究種目：若手研究 (B)  
研究期間：2008～2009  
課題番号：20700639  
研究課題名 (和文) データ・テキストマイニングを活用した授業評価アンケートの分析

研究課題名 (英文)  
Analysis of Course Evaluation Questionnaire Utilizing Data and Text Mining Technology  
研究代表者  
松河 秀哉 (HIDEYA MATSUKAWA)  
大阪大学・大学教育実践センター・助教  
研究者番号：50379111

## 研究成果の概要 (和文)：

大阪大学の授業評価アンケートのデータを用いて、(1) 授業の満足度や理解度に影響を与える具体的な要因、及び (2) 数値で回答されるアンケート項目と、自由記述中に含まれるキーワードの関係を検証した。その結果、講師の話し方や授業の難易度など、授業の理解度や満足度に影響を与える要因やその程度が明らかになった。また、「速い」という単語が自由記述に含まれる場合、学生は授業が難しいと感じるなど、字義どおりではない意味をもつキーワードが存在することが明らかになった。

## 研究成果の概要 (英文)：

I examined (1) the concrete factor that affects the satisfaction to and understanding of the lecture, and (2) the relationship between items included in the questionnaire and keywords included in the free description, utilizing the data collected as the course evaluation questionnaire at Osaka University. As a result, it became clear that Speaking skill of the lecturer or difficulty level of the lecture etc. would affect the satisfaction and understanding level of the lecture. In addition, I found that some keywords did not represent the literal meaning: When the keyword “rapid” was included in the free description, learners tended to feel that the lecture was “difficult”.

## 交付決定額

(金額単位：円)

	直接経費	間接経費	合計
2008 年度	2,400,000	720,000	3,120,000
2009 年度	600,000	180,000	780,000
年度			
年度			
年度			
総計	3,000,000	900,000	3,900,000

研究分野：教育工学

科研費の分科・細目：科学教育・教育工学・教育工学

キーワード：データマイニング・テキストマイニング、授業評価アンケート、フィードバック

## 1. 研究開始当初の背景

日本の各大学においては全学規模の授業評価アンケートが実施されるようになってきたが、その結果は極めて簡単に集計されるだけのことが多く、そこで蓄積されているであろう大量の数値データやテキストデータは有効に活用されているとは言い難い。

従来、業評価アンケートのデータを扱った研究では、主に因子分析や共分散構造分析などを用いて、授業に対する満足度や理解度との関係を考察する研究などが行われてきたが、こうした手法では複数のアンケート項目が一つの因子にまとめられることによって必然的に抽象化される。

そうして得られる結果の解釈には専門的な知識が必要であり、その結果授業を担当している一般の教員にそのままフィードバックすることは困難であった。

そこで本研究ではデータ・テキストマイニングの手法に注目するデータマイニングとは多くの変数を含む大量のデータセットの中から価値のある情報を引き出す解析的手法の総称であり、テキストマイニングはその考え方をテキストデータに応用したものである。

こうした手法を授業評価アンケートデータに応用すれば、Aという教育手法とBという教育手法を同時に採用している教員の授業は、学生の理解度が高い、もしくは、自由記述に「うれしい」というキーワードが入っている場合は学生の満足度が高い、といったことが分かる可能性がある。こうして得られる結果は、アンケートに採用されている具体的な質問項目の組み合わせとなり、従来の因子分析等を用いた場合よりも、遙かに解釈が容易である。このため、分析を通して得られた結果を、授業を行う教員に、より直接的に伝えることが可能である。

今回の研究では、データマイニング的な手法を用いることによって、従来の研究のように授業の理解度や満足度に影響を与える具体的な要因を因子化・抽象化することなく同定し、一般の教員に対して理解しやすい形で分析結果を提供できるものと考えられる。

また、テキストマイニングによってアンケート項目とキーワードの関係を調べることで、満足度や理解度、授業の難易度などに関係が深いキーワードが発見されることが期待できる。

## 2. 研究の目的

以上のような背景から、本研究では次の2点を目的とする。

(1) 授業の満足度や理解度に影響を与える具体的な要因をアンケート項目の中から明らかにする。

(2) 数値で回答されるアンケート項目と、

自由記述中に含まれるキーワードの関係を明らかにする。

## 3. 研究の方法

上記の目的を達成するため、まず、2009年までの授業評価アンケートのデータを格納したデータベースを作成する。

自由記述に関しては、形態素分析を施した上で、各単語とその単語が含まれていたアンケートとの対応関係を保持した上でデータベース化する。

(1) 上記のデータベースを利用し、目的1については次のような方法で、満足度や理解度に影響を与える具体的な要因を検証する。

大阪大学の共通教育において2007年度以降に行われた授業評価アンケートは、表1に示した質問項目から構成されており、(2)~(12)は順序尺度もしくは間隔尺度として取り扱い可能な1から5の5段階で回答するようになっている。

そこで、2007年1学期~2009年2学期までの83461件の回答データを利用し、設問(2)~(10)までの設問にどのように回答するかによって、理解度を問う設問(11)や満足度を問う設問(12)の値がどのように変化するかを回帰二進木分析を用いて分析する。

回帰二進木分析は、目的となる変数をもつ集合を、いくつかの基準となる変数の特定の値で次々に2分割していく手法で、分割後の集合の分散が分割前の分散に比べてできるだけ小さくなるように、基準となる変数とその値が決定される点に特徴がある。

表1 アンケートの質問項目

(1)あなたの所属学部・学科を選択して下さい。
(2)授業の難易度は
(3)授業内容の分量は
(4)理解度を深めるための配慮(小テスト、中間レポート、ノート提出等)が払われていた。
(5)一学期を通して授業は体系的に組み立てられ、適切な時間配分をもって行なわれた。
(6)教員の話し方・説明の仕方は、わかりやすかった。
(7)黒板・OHP・PowerPointのスライド・ハンドアウト等による情報提示は、わかりやすかった。
(8)学生とのコミュニケーション(質問を促す、ディスカッションの機会を設ける、Web-CTを利用して質問に答える等)にたいして教員は、熱心だった。
(9)あなたがこの授業に時間通りに出席した割合は、
(10)この授業の予習・復習に当てた平均時間(1週当たり)は、
(11)全体として授業内容を、よく理解することができた。
(12)総合的に見て、この授業に私は満足している

(2) 次に、目的2については以下の方法で、数値で回答されるアンケート項目と、自由記述中に含まれるキーワードの関係を明らかにする。

前述したように、2007年度以降のアンケートは、1から5の5段階で回答する11の設問

と自由記述中から構成されている。

過去の自由記述の全集合から、 $x$  個のキーワードを抽出し、各キーワードが 1 件 1 件の自由記述に含まれているか/いないかという 2 値情報を得ることができれば、数値の回答情報を持つ 11 個の変数と、2 値の情報を持つ、 $x$  個の変数をあわせた、 $11+x$  変数間で相関ルールを求めることができる。

キーワードの選定にあたっては、形態素分析素ソフトである chasen を用いて、2004 年度以降の前自由記述 42913 件の自由記述に形態素分析した結果データベースに格納されている、14519 種類、1327054 件の単語データを利用し、 $tf \cdot idf$  法を用いて、以下の式で重みづけを行う。

$$tf \cdot idf = f_{ti} \cdot \log(N/N_i)$$

( $f_{ti}$ : 単語  $t_i$  の総出現頻度、 $N$ : 自由記述件数、 $N_i$ : 単語  $t_i$  が含まれる自由記述件数)

その結果、重みづけの値が高いものから、名詞 500 個、動詞 200 個、形容詞 100 個、副詞 100 個をキーワードとして選択する。

その上で、アンケート内容が異なる健康スポーツ科目を除いた 2007 年度以降の全のアンケート結果の中で、数値での回答と自由記述が共に存在する 27882 件のデータを用いて、単語 900 変数 + 数値での回答 11 変数の、計 911 変数の間で相関ルールを求める。その後の解釈を簡単にするため、5 段階の回答は、1, 2 を 1, 3 を 2, 4, 5 を 3 の 3 段階に変換してから分析を行う。

相関ルールの算出にあたっては、hristian Borgelt によって開発された、フリーソフトである、apriori.exe を用いる (<http://www.borgelt.net/apriori.html>)。

実際のルールの抽出にあたっては、確信度が 70% 以上、リフト値が 1.5 以上で、条件部分の数が 5 個以下のルールのみを抽出する。

#### 4. 研究成果

(1) 授業の満足度や理解度に影響を与える具体的な要因については、図 1 の様な結果が得られた。

この図では、頂点のノードが授業の理解度に関する全回答の集合を表しており、5 層にわたって、授業の理解度の差を最大化させるように、アンケートの含まれるいずれかの設問項目によって、集団が 2 分割されていった結果を表している。分割の結果、授業の理解度(5 段階)の平均値が 1.2 の集団から、4.6 の集団まで 31 の集団に分割された。

この結果、例えば理解度の平均値が 1.2 の集団は、教員の話し方や説明の仕方が分かりやすいかについては、5 段階で 1 (ほとんどそう思わない)、授業の難易度については 5 (難しすぎる)、学生とのコミュニケーションのが熱

心だったかについては 1 (ほとんどそう思わない)、授業が体系的に組み立てられていたと思うかについては 1 (ほとんどそう思わない)、と回答した集団であることが読み取れる。

一方、理解度の平均値が 4.6 の集団は、教員の話し方や説明の仕方が分かりやすいかについては、5 段階で 5 (かなりそう思う)、授業の難易度については 1~3 (やさしすぎる~ちょうどよい)、授業が体系的に組み立てられていたと思うかについては 5 (かなりそう思う)、板書やパワーポイントの情報提示が分かりやすいかについては 5 (かなりそう思う)と回答した集団であることが読み取れる。

このように、本研究の結果から、アンケートの各項目にどのように回答した回答者の理解度が高い/低いかを全体的な傾向として把握することが可能である。

また上記の 2 例で示したアンケート項目については、はじめに示された項目の方が、後に示された項目よりも、理解度を左右する力が強い項目である。つまり、最初の方に出てきた項目について、改善策を考える方が、より効果的といえる。

この結果を利用すれば、どの項目をどのような順序で改善していけば、学生の理解度が向上する可能性があるかを、根拠をもって示すことができるため、これまで経験則に頼りがちであった FD (授業改善活動) の改善に繋がると考えられる。

なお、同様の結果が、満足度についても得られている。また、授業全体のデータから計算した結果だけではなく、より範囲を限定した科目群ごととの結果も、計算可能であるため、授業改善活動の規模や目的に応じて、本研究の結果は利用することが考えられる。

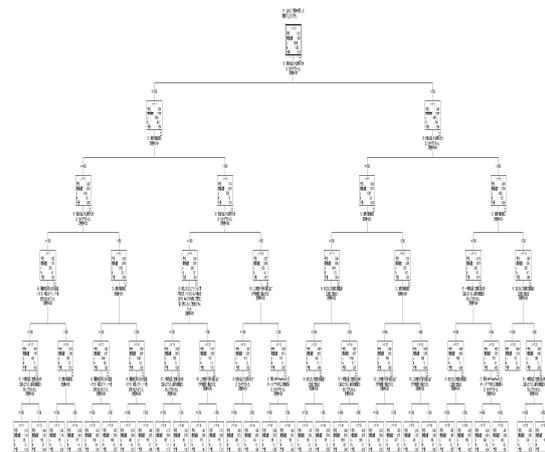


図 1 理解度についての回帰二進木

(2) 数値で回答されるアンケート項目と、自由記述中に含まれるキーワードの関係について、前述の方法で相関ルールを抽出したところ、約 2900 万件のルールが得られた。

最も解釈を容易にするため、その中で、条件部の数が一つだけで、「ある単語Xが含まれるならば、ある設問の3段階の解答がYである」という形になっているルールのみを抽出したところ、112個のルールが得られた。

そのルールの中で、支持度・確信度・リフトの積が大きいものから10個を選び、解釈しやすいように日本語に直したものが表2である。

これを見ると、例えば「速い」という単語が出現した場合、学生は授業内容の分量が多いと感じると同時に、授業を難しいと感じている可能性が高いという、字義通りではない2通りの意味を、具体的な数字を伴って読み取ることが可能である。

従来、テキストマイニングは、クラスター分析やコレスポネンズ分析(塚本 2006)など、単語間の関係を様々な手法で視覚化することが多く、分析結果の意味を解釈するのが難しいという問題があった。しかし、相関ルールを用いるこの方法は、単語の出現の有無と、アンケートの数値の回答データの関係が、確信度、リフト、支持度といった、はっきりとした意味を持つパラメータを伴ったルールという形で示されるため、意味の解釈が非常に簡単であるという利点を持っていると思われる。

こうした情報をルールが該当する自由記述と合わせて提示することで、より深い自由記述の解釈が可能になると考えられる。

以上の結果を踏まえたうえで、本研究ではさらに、分析結果を教員にフィードバックするためのシステムを試作した。

このシステムは、授業評価アンケートの数

- ・「難しい」が自由記述に含まれると、「(2) 授業の難易度は」という質問に対して、75.8%の確率で、学生から「むずかしい」と評価されます。その確率は、その単語が含まれない場合に比べて1.9倍に高くなります。このルールに該当するのは、自由記述がある回答の7.6%です。
- ・「楽しい」が自由記述に含まれると、「(8) 学生とのコミュニケーション(質問を促す、ディスカッションの機会を設ける、Web-CTを利用して質問に答える等)にたいして教員は、熱心だったかどうか」という質問に対して、73.8%の確率で、学生から「そう思う」と評価されます。その確率は、その単語が含まれない場合に比べて1.58倍に高くなります。このルールに該当するのは、自由記述がある回答の4.9%です。
- ・「英語」が自由記述に含まれると、「(8) 学生とのコミュニケーション(質問を促す、ディスカッションの機会を設ける、Web-CTを利用して質問に答える等)にたいして教員は、熱心だったかどうか」という質問に対して、71.3%の確率で、学生から「そう思う」と評価されます。その確率は、その単語が含まれない場合に比べて1.53倍に高くなります。このルールに該当するのは、自由記述がある回答の4.2%です。
- ・「質問」が自由記述に含まれると、「(8) 学生とのコミュニケーション(質問を促す、ディスカッションの機会を設ける、Web-CTを利用して質問に答える等)にたいして教員は、熱心だったかどうか」という質問に対して、71.4%の確率で、学生から「そう思う」と評価されます。その確率は、その単語が含まれない場合に比べて1.53倍に高くなります。このルールに該当するのは、自由記述がある回答の2.1%です。
- ・「ついていけない」が自由記述に含まれると、「(2) 授業の難易度は」という質問に対して、75.5%の確率で、学生から「むずかしい」と評価されます。その確率は、その単語が含まれない場合に比べて1.89倍に高くなります。このルールに該当するのは、自由記述がある回答の0.9%です。
- ・「速い」が自由記述に含まれると、「(2) 授業の難易度は」という質問に対して、70.3%の確率で、学生から「むずかしい」と評価されます。その確率は、その単語が含まれない場合に比べて1.76倍に高くなります。このルールに該当するのは、自由記述がある回答の1%です。
- ・「会話」が自由記述に含まれると、「(8) 学生とのコミュニケーション(質問を促す、ディスカッションの機会を設ける、Web-CTを利用して質問に答える等)にたいして教員は、熱心だったかどうか」という質問に対して、78.4%の確率で、学生から「そう思う」と評価されます。その確率は、その単語が含まれない場合に比べて1.68倍に高くなります。このルールに該当するのは、自由記述がある回答の0.9%です。
- ・「楽しむ」が自由記述に含まれると、「(12) 総合的に見て、この授業に私は満足しているかどうか」という質問に対して、94.3%の確率で、学生から「そう思う」と評価されます。その確率は、その単語が含まれない場合に比べて1.53倍に高くなります。このルールに該当するのは、自由記述がある回答の0.7%です。
- ・「グループ」が自由記述に含まれると、「(8) 学生とのコミュニケーション(質問を促す、ディスカッションの機会を設ける、Web-CTを利用して質問に答える等)にたいして教員は、熱心だったかどうか」という質問に対して、84.5%の確率で、学生から「そう思う」と評価されます。その確率は、その単語が含まれない場合に比べて1.82倍に高くなります。このルールに該当するのは、自由記述がある回答の0.6%です。
- ・「コミュニケーション」が自由記述に含まれると、「(8) 学生とのコミュニケーション(質問を促す、ディスカッションの機会を設ける、Web-CTを利用して質問に答える等)にたいして教員は、熱心だったかどうか」という質問に対して、76.9%の確率で、学生から「そう思う」と評価されます。その確率は、その単語が含まれない場合に比べて1.65倍に高くなります。このルールに該当するのは、自由記述がある回答の0.7%です。

値データを集計して、表やグラフで視覚的に提示する(図2)とともに、講義の理解度や満足度と、授業の難易度など、アンケートで問われる項目との間にどのような関係があるかを、回帰二進木分析で得られた樹形図によって提示する。

また、自由記述については、a.相関ルールを用いたアンケート項目と自由記述に含まれる単語の関係の表示、b.係り受け解析による自由記述の要点の抽出機能、c. tf-idf 法を用いたキーワード抽出機能、d.選択したキーワードや、自由記述の件数、平均情報量の変化などを時系列で表示する機能を備えている(図3)。分析は、各講義ごと、いくつかの科目を集めた科目群ごと、全ての講義の3段階のレベルについて行うことができる。

今後はこのシステムを発展させ、実際にFD活動に活用することを検討していきたいと考えている。

表2 生成されたルールの例

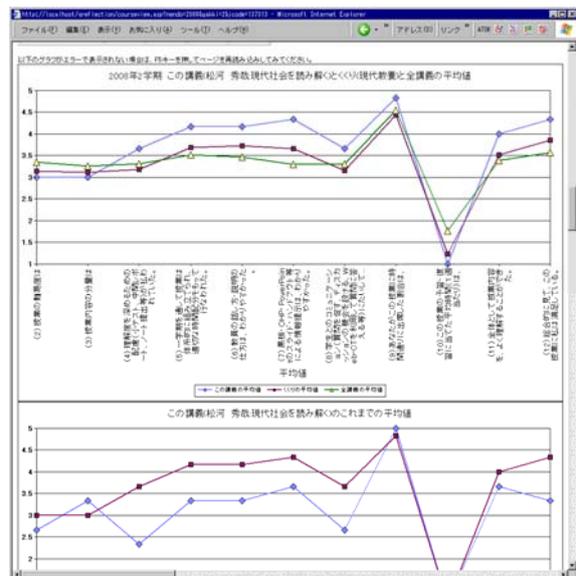


図2 集計機能の一例

図3 自由記述の分析機能の一例

## 5. 主な発表論文等

〔学会発表〕(計3件)

①松河秀哉、テキストマイニングを応用した自由記述に対する分析システムの紹介、平成21年度 第5回 大学コンソーシアム石川FD研修会テーマ:「IR ( Institutional

松河 秀哉:現代社会を読み解くに対する自由記述

自由記述	相関ルール分析による、自由記述に含まれる単語と各設問の解答との関係
他の一般教養の授業とは比べ物にならないくらいいろいろな工夫をこらした授業をして下さって、とてもおもしろかったです。この授業をとらなかつたら一生経験できなかったらと思うようなことに触れる機会がたくさんありました。	「経験」が自由記述に含まれると、「(8)学生とのコミュニケーション(質問を促す、ディスカッションの機会を設ける、Web-CTを利用して質問に答える等)」について教員は「熱心だったかどうか」という質問に対して、73.4%の確率で、学生から「そう思う」と評価されます。その確率は、その単語が含まれない場合に比べて1.82倍に高まります。このルールに該当するのは、自由記述がある回答の0.5%です。

5件が空欄でした。

松河 秀哉:現代社会を読み解くの要素

[元ページへリンクする](#)

係り元	係り先	元の自由記述
授業とは	おもしろかったです	<a href="#">元の自由記述を表示</a>
比べ物に	ならない	<a href="#">元の自由記述を表示</a>
工夫を	こらした	<a href="#">元の自由記述を表示</a>
授業を	して下さい	<a href="#">元の自由記述を表示</a>
授業を	とらなかつたら	<a href="#">元の自由記述を表示</a>
経験できなかったら	思うような	<a href="#">元の自由記述を表示</a>
ことに	触れる	<a href="#">元の自由記述を表示</a>
機会が	ありました	<a href="#">元の自由記述を表示</a>

松河 秀哉:現代社会を読み解くのキーワード

名詞

\*は未知語を表す

くくりの中の重要度		全体の中での重要度	
形態素	頻度 tfidf	形態素	頻度 tfidf
一 生	1 9.7	比 べ 物	1 13.49
比 べ 物	1 9.7	一 生	1 11.07
経 験	1 7.37	一 般	1 8.61
一 般	1 6.89	教 養	1 8.23

Research) と授業評価」、大学コンソーシアム石川 シティカレッジ教室 1、金沢、2010年 2月 26日、招待講演

②松河秀哉、データ・テキストマイニングを活用した授業評価アンケートの分析とフィードバックシステムの開発、日本教育工学会研究報告集 Jset09-5、pp65-70、2009年 12月 19日

③松河秀哉、データ・テキストマイニングを応用した授業評価アンケートフィードバックシステムの開発、教育システム情報学会第34回全国大会論文集、pp486-487、2009年 8月 20日

## 6. 研究組織

### (1) 研究代表者

松河 秀哉 (HIDEYA MATSUKAWA)

大阪大学・大学教育実践センター・助教

研究者番号：50379111