

## 相関ルールを用いた誤答パターンの分析

Analysis of the Pattern of the Erroneous Answers Utilizing Association-rule Algorithm .

松河秀哉\*<sup>1</sup>・北村智\*<sup>2</sup>・山内祐平\*<sup>2</sup>・中野真依\*<sup>3</sup>・金森保智\*<sup>3</sup>・宮下直子\*<sup>3</sup>

Hideya MATSUKAWA\*<sup>1</sup>・Satoshi KITAMURA\*<sup>2</sup>・Mai NAKANO\*<sup>4</sup>・

Yasutomo KANAMORI\*<sup>4</sup>・Naoko MIYASHITA\*<sup>4</sup>

大阪大学大学教育実践センター\*<sup>1</sup>・東京大学大学院情報学環\*<sup>2</sup>・株式会社ベネッセコーポレーション\*<sup>3</sup>

Institute for Higher Education Research and Practice\*, Osaka University\*<sup>1</sup>, Interfaculty Initiative in Information Studies, The University of Tokyo\*<sup>2</sup>, Benesse Corporation\*<sup>3</sup>

〈あらまし〉本研究では、データマイニングの一手法である相関ルールを用いて、中学 1 年から高校 1 年までの数学の模擬試験の誤答データの時系列的な分析と、高校 1 年向けの英語の模擬試験の回答データの分析を行った。その結果、いずれの分析でも、生成されたルールから、誤答の原因をある程度推測可能であることが示唆された。

〈キーワード〉 相関ルール・誤答パターン・データマイニング

### 1. はじめに

教育工学の関連領域において、誤答分析はこれまで、算数などの教科を中心に、学習者が持っている誤りのモデルを明らかにしようとするバグ研究一環として盛んに行われてきた(Wenger 1987)。また、CAI 教材の作成の前段階として、教材が対象とする領域の典型的な誤答を収集することも数多く試みられている(洲脇ら 2001)。

こうした研究においては、誤答の分析は、例えば引き算の方法など、ごく限られた領域に焦点を当て、その範囲の誤答を研究者自身が分類して詳細に調べることが多く、労力の観点から、より広い領域間の関連はほとんど注目されることがなかったと言える。

一方、近年では、情報技術の急速な進歩に伴い、膨大な情報の中から有用なデータを発見するための、データマイニングの手法が発展してきた。

自動的に膨大な情報を処理するデータマイニングの手法を利用すれば、これまで困難であった、より広い領域にわたる、誤答の関係性を分析することが可能であると考えられる。

そこで、本研究では、中学 1 年の 1 学期から高校 1 年 1 学期まで 3 年間にわたる数学の学力テストの正誤データ、及び、高校 1 年時の英語の学力テストの回答データに対して、

相関ルールを用いた分析を行うことにより、中学段階と高校段階での誤答の関連性や、様々な問題のタイプを横断して現れる特徴的な誤答など、広範囲のデータにおける誤答のパターンを明らかにすることを目的とする。

### 2. 分析方法

#### 2.1. 相関ルールと評価の尺度

相関ルールは、A ならば B であるの様に、ある事象と別の事象の共起関係の起こりやすさを表すものである。例えば、ある店でパンとバターを同時に購入する人の 90 %は、ミルクも購入するといったものが、相関ルールのひとつの例である。マーケティング分野では、レジで収集される POS データを用いて、こうした分析が日常的に行われており、商品の効率的な陳列などに応用されている。

最近では、教育の分野においても、大量の学習履歴データの中から、「問題 A と B に正答すれば、90 %の確率で合格」といったように、最終成績を予測するルールを探索する研究(Minaei-Bidgoli 2003)などが行われるようになってきている。

相関ルールの「よさ」を計る尺度としては、支持度(support)と確信度(confidence)がよく用いられる。例えば、ある模擬試験の受講者数を N、その中で、問題 1 と問題 2 を間違った

人の数を  $N_a$ 、問題 3 を間違えた人の数を  $N_b$ 、問題 1, 2, 3 を同時に間違えた人の数を  $N_{ab}$  とすると、問題 1 と問題 2 に間違えなければ、問題 3 も間違えというルールは支持度は、 $N_{ab}/N$  と表すことができる。またそのルールの確信度は、 $N_{ab}/N_a$  と表すことができる。一般に、支持度が高いルールほどよく通用するルールであり、確信度が高いほどそのルールは信頼できる。

確信度は、有効なルールの判断基準となり得るが、誤答を分析する場合、結論部分に難問が含まれていると、確信度だけでは不十分となる場合もある。例えば上のルールで、問題 3 の誤答率が元々 90% 近くあったと仮定すると、「問題 1 と問題 2 に違えば、問題 3 も間違える」というルールは確信度が 90% あっても、「問題 1 と問題 2 に間違える」という条件部分はあまり特別な意味があるといえなくなってしまう。

そこで利用されるのが、「リフト」という指標である。リフトはルール全体の確信度を結論部分の確信度で割ったもので、 $(N_{ab}/N_a)/(N_b/n)$  と表すことができる。例えば、問題 3 単体の誤答率が 10% なのに、「問題 1 と問題 2 に違えば、問題 3 も間違える」というルール全体の確信度が 90% ならばリフトの値は 9 である。これは、「問題 1 と問題 2 に間違える」という条件がつくことで誤答が 9 倍起こりやすくなったことを意味する。反対に、問題 3 単体の誤答率が 89% で、ルール全体の確信度が 90% の場合は、リフトは  $90/89 \approx 1$  であり、確信度が同じでも、「問題 1 と問題 2 に間違える」という条件は、あまり重要でないことが判断できる。

本研究の分析においては、関連ルールの絞り込みに、上述した、確信度、支持度、リフトを尺度として用いた。

## 2.2. 分析対象

分析対象とするデータは、株式会社ベネッセコーポレーションの進研ゼミ中学講座と高校講座の全国学力模擬試験のデータである。

数学の誤答分析には、中学 1 年の 1 学期から中学 3 年 2 学期までの毎学期と、高校 1 年 1 学期の計 9 回分の模擬試験(2004 年度から 2007 年度)を全て受験している 4444 名の各設問に対する正誤データを用いた。分析対象と

なった変数の数は模試の設問の数と同じ 389 であった。

英語の誤答分析には、高校 1 年 1 学期のマークシート形式の模擬試験(2007 年度)の各設問に対する 25919 件の回答データを用いた。模試の設問数は 117 問であったが、複数の選択肢をダミー変数化したため、分析対象となった変数の数は 917 であった。

## 2.3. 分析ツール

### 2.3.1. 分析に用いたハードウェア

分析を行った PC のスペックは以下のとおりである。

CPU:Xeon 5365 3.0Ghz (4core)

Memory:DDR-2 677 FB-DIMM 16GByte

HDD:1TByte (500GByte × 4 Raid1+0)

OS:WindowsXP x64 Edition

### 2.3.2. 分析に用いたソフトウェア

関連ルールの分析には、Christian Borgelt によって開発された、apriori.exe を用いた。本ソフトは、関連ルールを求めるための代表的なアルゴリズムである、apriori アルゴリズムを実装したもので、フリーウェアとして、以下のサイトで公開されている。

<http://www.borgelt.net/apriori.htm>

このサイトで公開されている apriori.exe は 32bit 版であるが、本ソフトウェアを用いて WindowsXP 上で、試験的に分析を行ったところ、変数が多いため、プロセスがメモリ空間のに収まりきらず、分析が中断された。そこで、本ソフトウェアのソースコードを、VisualStudio2005 を用いて 64bit 環境(x64)向けにコンパイルして、分析に利用した。

分析の結果として生成された関連ルールの格納と検索を行うためのデータベースには MS-SQL Server2005 を用いた。

## 2.4. 分析手順

1. apriori.exe を用いて関連ルールを求めるにあたって、数学のデータに関しては、確信度 70%以上、ルールの支持度 0.1%以上、リフト 1.5 以上、条件の個数 3 以下という制約を与えた。英語のデータに関しては、確信度 90%以上、ルールの支持度 0.1%以上、リフト 1.5 以上、条件の個数 3 以下という制約を与えた。

2. apriori.exe の出力結果は、自作のスクリプトを用いてタブ区切りのデータに整形した

上で、DBにインポートした。

3. DBに格納したデータに対して、SQL文を発行し、確信度やリフト値によるソートをかけながら、誤答間の関係を解釈可能なデータを筆者が単独で探索した。

### 3. 結果

数学のデータに関しては、所与の条件で、420742255件のルールが作成された。このうち、ルールの結論部分に、高1の設問のみを含み、条件部分に中学の設問のみを含むルールは48641199件あった。ルールの件数が膨大であるため、そのすべてを精査することは不可能であった。したがって、ここでは、結論部の誤答率とルールの支持度が高く、最も明確に誤答間の関係を解釈可能であった4問の問題の組み合わせを図1から図4に示す。

図1から図3は中学時代の誤答を表し、図4は高校時代の誤答を表している。問題A、問題B、問題Cに同時に間違えたのは351人(全体の7.9%)であり、そのうち問題Dも間違えたのは250人(全体の5.6%)だった。つまり問題Aと問題Bと問題Cに間違えれば問題Dも間違えるというルールの確信度は71.2%であった。ルールのリフト値は2.694であり、問題Aと問題Bと問題Cを同時に間違えるという条件がつくことによって、問題Dは約

2.7倍間違えやすくなることが示されている。

次に英語の誤答に関して所与の条件で分析したところ、152352件のルールが生成された。このうち、誤答の原因として未回答を含まないものは123件あった。その中で、結論部の誤答率とルールの支持度が高く、最も明確に誤答間の関係を解釈可能であった4問の問題の組み合わせを図6に示す。

この結果のなかで、問題Aの回答が3,2かつ問題Bの回答が6,2かつ問題Cの回答が3であった学習者は37人であり、このうち、問題Dの回答が1なのは35人(全体の0.14%)であった。したがって「問題A,B,Cに間違えば問題Dも間違う」というルールの確信度94.6%であり、リフト値は2.34であることから、問題A,B,Cを同時に間違える学習者は、そうでない学習者より、2倍以上問題Dを間違えやすくなることが示されている。

### 4. 考察

#### 4.1. 数学の結果について

まず、数学の分析結果についての誤答間の関連性について考察する。はじめに、相関ルールの結論部分である高校時点の問題Dの性質からみると、この問題を解くには、図5のような図示を行い、次にその情報を用いて方程式を立式し、方程式を計算して解を出す

問11 aでわると商がbで余りがcとなる正の整数nがあります。整数nをa, b, cを用い

て表しなさい。

- ①  $n = a + b + c$
- ②  $n = a(b + c)$
- ③  $n = ab + c$
- ④  $n = \frac{b}{a} + c$

図1 中学時代の問題A(全体の誤答率29.5%)

問28 右の図で、直線ℓの式は $y = x + 1$ 、直線mの式は $y = -2x + 10$ です。直線ℓとmの交点をA(3, 4)とし、直線ℓ, mとx軸との交点をそれぞれB, Cとします。点Aを通り、△ABCの面積を2等分する直線の式を求めなさい。

- ①  $y = 2x - 5$
- ②  $y = 2x + 2$
- ③  $y = 4x - 8$
- ④  $y = -\frac{3}{2}x - 3$

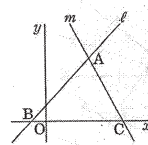


図2 中学時代の問題B(全体の誤答率21.9%)

次のそれぞれの場合について、yをxの式で表しなさい。

問17 周の長さが30cmである長方形の縦の長さがx cm、横の長さがy cm

- ①  $y = \frac{30}{x}$
- ②  $y = -x + 30$
- ③  $y = -x + 15$
- ④  $y = 15x$

図3 中学時代の問題C(全体の誤答率35.1%)

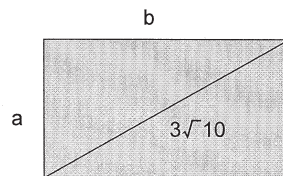


図5 問題Dを解くのに必要な図示

(6) 長さが24cmの針金を折り曲げて長方形を作ると、対角線の長さが $3\sqrt{10}$ cmになった。

この長方形の面積は、サシcm<sup>2</sup>である。

図4 高校時代の問題D(全体の誤答率26.4%だが、上記3題とも間違えると誤答率は71.2%)

問題A(正答率61.1%、このタイプの誤答率1.7%)

二つめの角を左に曲がると、白い建物がある。

・ You will find a white building ① — — ② — the second

1. and 2. at 3. if 4. left 5. turn 6. you

正答: You will find a white building if you turn left at the second  
誤答: You will find a white building if you turn at left the second

問題B(正答率47.2%、このタイプの誤答率7.2%)

父はたくさんの助言を私にしてくれた。

・ My father ① — — ② —

1. advice 2. a lot of 3. gave 4. me 5. to 6. was given

正答: My father gave a lot of advice to me.  
誤答: My father was given to me a lot of advice.

問題C(正答率14.6%、このタイプの誤答率41.9%)

・ We finally ① the top of the mountain.

1. reached 2. reached at 3. reached on 4. reached to

正答: We finally reached the top of the mountain.  
誤答: My father reached on the top of the mountain.

問題D(正答率47.8%、このタイプの誤答率38.8%)

→ABC誤答がそろると94.6%に高まる)

・ He showed ①.

1. me to his computer 2. his new computer me

3. me of his computer 4. his new computer to me

正答: He showed his new computer to me.  
誤答: He showed me to his new computer.

#### 図6 英語の誤答パターン例

いう手順を踏む必要があることが分かる。

次に、相関ルールの条件部分に含まれる中学時点の各設問の性質をみると、問題Aは文章から立式を行う問題、問題Bは文章の図示、図からの立式、式の計算を行う問題、問題Cは文章の図示と図からの立式を行う問題である。

このことから、相関ルールの条件部分に含まれている問題は、結論部分の問題を解くために必要なスキルの、一部、もしくは全部が必要な問題であることが分かる。即ち、このルールあてはまる学習者は、文章で書かれた問題から図示や立式をすること、その上で計算を行うことのどちらか、もしくは両方に以前からつまずいていた可能性が考えられる。

このように、相関ルールを用いることで、中学時代と高校時代のように時系列的に離れた誤答間の関連性の分析から、その学習者がある問題に誤答する原因を過去にさかのぼって探れる可能性が示唆され、その原因に応じて学習者に効果的な処遇を行うことなどが期待される。

#### 4.2. 英語の結果について

図6中の各設問について、その性質を見てみると、問題Aと問題Cは助詞を伴わない動詞の扱いに関連する誤答、問題Bと問題Dは二重目的語を取る動詞に関連する誤答、であり、一見したところ、全体としての関連は薄そうに思われる。しかし、日本語に直して考

えてみると、「左に曲がる」、「私にしてくれた」、「山頂に着いた」、「私に見せた」というように、全て「に+動詞」という訳が含まれる問題であることが分かる。即ち、このルールにあてはまる学習者は、日本語に引きずられて、動詞の扱いを誤っている可能性が考えられる。

このように、分析対象が1回の模擬テストの結果であっても、選択肢の内容まで含めて相関ルールを求めることで、比較的詳細に誤答の原因を推測できる可能性が示唆された。

#### 5. まとめ

本研究では、データマイニングの一手法である相関ルールを用いて、中学1年から高校1年までの数学の模擬試験の誤答データの時系列的な分析と、高校1年向けの英語の模擬試験の回答データの分析を行った。その結果、いずれの分析でも、生成されたルールから、誤答の原因をある程度推測可能であることが示唆された。

相関ルールから誤答の原因を推測するには、ルールの解釈が不可欠であるが、今後は、この作業に教科の専門家の協力を求めることで、より重要な誤答間の関連が発見されることが期待される。

#### 参考文献

Minaei-Bidgoli, B., Kashy, D.A., Kortemeyer G., Punch, W.F., (2003) Predicting Student Performance: An Application of Data Mining Methods with an educational Web-based System, IEEE conf. on Frontier In Education FIE 2003

岡本敏雄・溝口理一郎(監訳)(1990)知的CAIシステム. オーム社, 東京 (Etienne Wenger, (1987) Artificial Intelligence and Tutoring Systems. Los Altos, CA)

洲脇史朗・宮地功(2001)中学校における「一次方程式」の誤答分析. 電子情報通信学会技術研究報告. ET, 教育工学, Vol.101, No.397(20011019) pp. 1-8

\* 本研究は、東京大学大学院情報学環ベネッセ先端教育技術学講座のプロジェクトとして、(株)ベネッセコーポレーションと共同で行われている。

