LMS 蓄積データを用いた学習特徴の抽出における 変数の粒度の検討

小川賀代^{†1} ピトヨ ハルトノ^{†2}

近年,LMS の普及に伴い,アクセスログ,成績,アンケートなどの学習に関する様々なデータの取得が可能になり,膨大なデータとなりつつある。これらの履歴データをデータマイニングすることで学生の理解度や学習パターンの把握ができ、個人に適した学習の提供が期待されている。本報告では、定量的な複数の指標を用いてクラスター数を決定し、各クラスターの学習特徴を解析、可視化する方法を述べる。また、解析に使用する変数の粒度についても検討を行い、粒度の違いにより、異なる解析結果が得られることが確認でき、解析者の見たい情報に応じて粒度を変化させる必要性があることがわかった。

Evaluation of the granularity of variables for extracting learning characteristics using data accumulated by LMS

KAYO OGAWA^{†1} PITOYO HARTONO^{†2}

The popularization of Learning Management Systems recently made it possible to obtain educational data such as students' access logs and grades. Now, a vast amount of data is becoming available to researchers. Data-mining of these historical data is expected to contribute to the understanding of the level of comprehension and learning characteristics of students, enabling the provision of education suitable to individual students. This paper reports a method for determining the number of clusters using multiple quantitative indices, and analyzing and visualizing the learning characteristics of each cluster. The paper also examines the granularity of variables used for the analysis. We confirmed that the analyses yields different results depending on the granularity, and found that it is necessary to change granularities according to the information an analyst is interested to obtain.

1. はじめに

近年、学習を裏付け深めるための情報として、学習管理システム(LMS:Learning Management System)から取得できる学習履歴データを活用することが期待されている。特に膨大なデータの活用のためには、機械的に学習者を類似した学習特徴を持つグループに分類し、意味付けを行うことが必要となる。これを可能とする解析手法としてクラスター分析がある。しかし、分類数(クラスター数)は経験的に決定されているため、経験者の判断に左右される。そのため、データの自然な分布を必ずしも反映しているとはいえない。また、ビッグデータの取り扱いにおいては、できるだけ機械的に行える方が望ましい。

これまでの研究において、評価観点の異なる3つの指標を用いてクラスター数を決定し、学習者の特徴を抽出する解析手法を提案してきた. [1] また、分類結果の直感的な理解に向け、自己組織化マップ(SOM:Self-organizing map)による可視化の検討を行った。ただし、解析を行うにあたってデータ変数の粒度の影響やどの程度の粒度で用いるべきであるかについての検討はなされていない。そこで本研究では、本研究では学習履歴データの変数の粒度を変化させたときの学習特徴の抽出を行い、粒度の違いによる結果

2. 学習履歴データの解析手法

本研究で提案する解析手法の概要を図1に示し,以下に 説明する.

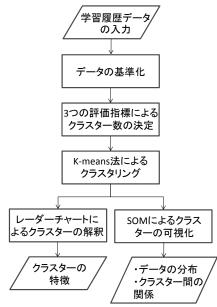


図1 解析手法の概要

2.1 3 つの指標を用いたクラスター数の決定

学習履歴データを定量的に分類するために, クラスター 分析を用いる. クラスター分析は, サンプル間の距離に基

の比較検討を行った.

^{†1} 日本女子大学

Japan Women's University

^{†2} 中京大学

Chukyo University

づき類似したサンプルを1つのクラスターに分類する手法 である.

本論文では、クラスター分析で最も用いられる手法である K-means 法で行う。K-means 法は、予め想定したクラスター数を用いて、式(1)に示す評価関数 を最小化することでデータを分類する。式(1)では、 $Z_j \in \{Z_{1j}, Z_{2j}, \cdot \cdot \cdot, Z_{nj}\}^{\mathcal{E}}$ 学習者jの基準化した学習サンプルとし、 C_k はクラスターkの重心であり、式(2)で計算する.

$$E = \sum_{k=1}^{K} \sum_{j=1}^{m} \alpha_{kj} \left\| C_k - Z_j \right\|^2$$
 (1)

$$C_k = \frac{1}{m_k} \sum_{z_j \in C_k} z_j \tag{2}$$

ただし、n:履歴項目数、m:学習者数、K:クラスター数、 $\alpha_{k_j} = \begin{cases} 1 & Z_j \in C_k \\ 0 & Z_j \notin C_k \end{cases}$

示す.

従来, クラスター数は解析者の経験と主観に大きく依存 して決定することが多いが、そのクラスター数がデータの 自然な分布を反映しているとは限らない. また, 想定した クラスター数を用いてクラスタリングの実行後に、その妥 当性 (Cluster Validity) を特定の評価指標を用いて評価する が、データには様々な分布特徴が存在するため、1 つの評 価指標のみを適用すると,異なる特徴を持つデータに対し, クラスター評価の妥当性を保証することができない. そこ で本解析手法では、評価観点の異なる3つの評価指標を用 意し, データに潜む自然な分布を反映したクラスター数を 決定する. ここでは、先ず、予め定めたクラスター数の範 囲で K-means 法を実施し、3 つの評価指標を適応する. そ して, 最も顕著な結果を示す評価指標を選び, その評価指 標の中で顕著に結果を示すクラスター数を決定する. この ように、クラスター数の決定は、グラフ化された定量的な データに基づいて、総合的に判断を行う. 本解析手法で用 いたクラスター数に対する評価指標を以下に示す.

(a) Dunn index[2]

クラスターの最大直径Dとクラスター重心値間の最小距離Mの比

$$DI = M/D \tag{3}$$

が最大となるクラスター数を選択する.

(b) DB index[3]

各クラスターにおいて、クラスター内分散 S_i とクラスター

重心値間の距離 M_{ii} の比

$$R_{ii} = \left(S_i + S_i\right) / M_{ii} \tag{4}$$

が最大になる値を R_i として、式(3)に示す類似度

$$\overline{R} = 1/K \sum_{i=1}^{K} R_i \tag{5}$$

を計算し、この値が最小となるクラスター数Kを選択する. (c) ベイズ情報量基準 (BIC:Bayesian information criterion)[4]

統計学における情報量基準の1つで、最大対数尤度 $\hat{l}_{M}(D)$

とモデルのパラメータ数 $P_{M}/2 \cdot \log R$ の差を計算し、この値が最大となるクラスター数を選択する.

2.2 可視化によるクラスターの分析

2. 1 で得られた自然なクラスター数に従い、データを K-means 法によりクラスタリングした後、形成されたクラ スターを可視化することで学習者の特性を抽出する.まず, 各クラスターの重心を項目毎に基準化し, クラスター別に レーダーチャートを描いて、クラスターの特徴を把握する. レーダーチャートによる視覚的な把握を行うことで, 学習 者の特性を容易に解釈することができる. これと同時に, Self-Organizing Map (SOM) を用いたクラスターの可視化 を行う. 1 個人が有する履歴項目は LMS に残された学習ロ グや成績など多次元に亘るため、個々人を比較することは 難しい. また, 多次元データのクラスタリングを行ってい るため、クラスター間の関係を理解するのは困難である. そこで本解析手法では、高次元データの位相関係を保ちな がら、低次元空間にそれらのデータを配置できる SOM を データの可視化手法として用いることにした. SOM は非線 形関係を捉えることができ、新しいデータをマップ上に配 置することができる. また, SOM の結果の上に, クラスタ ーごとに異なる色で結果を表示することで, クラスター間 の関係を把握することができる.

3. データ解析における変数の粒度の検討

3.1 使用する学習履歴データの概要

企業が提供している e ラーニングコンテンツの学習履歴 データを用いて評価実験を行った. 受講者数が 238 名, 講座は 10 章からなっており, 3~4 章ごとに全 3 つの章末テストと全 4 つの学習後アンケート, また最後に修了試験, 模擬試験が組まれている. 講義, テスト及びアンケートの受講に制限はない. 変数の詳細を以下に示す. 括弧内の数値は、変数の数を示している.

・講座の受講回数(10)

・修了テストの初回点数(1)

・アンケートの回答回数(4)

・修了テストの最高点数(1)

・各章末テストの初回点数(3)

・修了テストの受験回数(1)

・各章末テストの最高点数(3)

・模擬試験の受験回数(1)

・各章末テストの受験回数(3)

今回は講座の受講回数とアンケートの回答回数について、合計値を用いた粒度の粗い場合と、講座ごと、アンケートごとに分けた値を用いた粒度の細かい場合の2通りについて解析を行った.よって、粒度の粗い場合は15変数、粒度の細かい場合は27変数で解析を行った.また、これら全てのデータに対し基準化を行い解析を行った.

3.2 評価実験の結果

3.2.1 クラスター数の決定

データのクラスター分析を行うために、まずクラスター数を決定する。 今回のデータについては解析条件をクラスター数範囲 $2\sim10$ 、繰り返し回数 100 回として指標の計算を行った。その結果を図 2、3 に示す。

粒度の粗い場合は、指標値結果から極大・極小値は見当たらないが、指標値の変化量を見ると全ての指標でクラスター数4以降で値が小さくなっている。よって、クラスター数を4と決定した.一方、粒度の細かい場合は、Dunn IndexとBICにおいてクラスター数8で極大値を示している.よって、クラスター数を8と決定した.

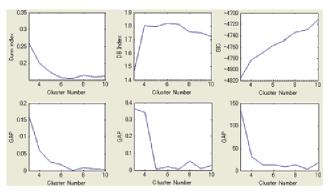


図2 指標計算結果(上)と指標値変化量(下) (粒度の粗い場合)

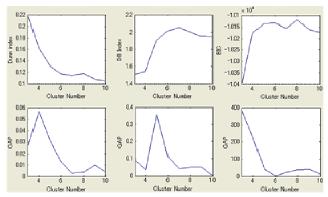
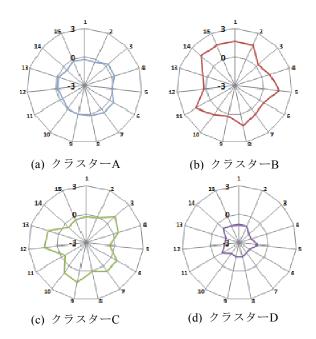


図3 指標計算結果(上)と指標値変化量(下) (粒度の細かい場合)

3.2.2 クラスターの特徴の把握

クラスターの学習特徴を把握するため,クラスターごとにレーダーチャートを描き、学習特徴の解釈を行った.レーダーチャートを図 4,5 に、レーダーチャートから読み取れる学習特徴と所属人数を表 1,2 に示す.なお、クラスタ

一名は A-H とした. また, 図 4, 5 の軸は解析に用いた変数を表しており, 軸の値の 0 はデータの平均値を示している.



1: 10章からなる講座の合計受講回数

2: 4つのアンケートの合計回答回数

3-5: 章末テスト1の初回点数,最高点数,受験回数

6-8: 章末テスト2の初回点数,最高点数,受験回数

9-11: 章末テスト3の初回点数,最高点数,受験回数

12-14: 修了テストの初回点数、最高点数, 受験回数

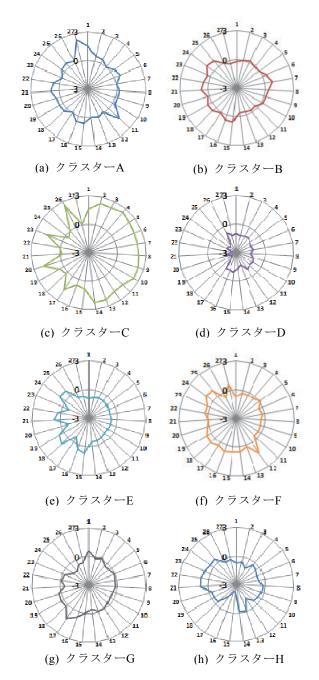
15: 模擬試験受験回数

図4 各クラスターのレーダーチャート (粒度の粗い場合)

表1 各クラスターの学習特徴と人数(粒度の粗い場合)

| Cluster | 学習特徴 | 人数 (人) | 割合 (%) |
|---------|---------------------------------------|-----------|-----------|
| A | 平均的なクラスター | 61 | 25.6 |
| В | 講義・テスト共に繰り返し取 り組む | 19 | 8.0 |
| С | 講義・テスト共に繰り返し取 り組まないが, テスト点数は 高い | 145 | 60.9 |
| D | 全体的にコンテンツにしっ かり取り組まない | 13 | 5.5 |

変数の粒度を細かくすることにより、クラスター数が増え、さらに細かい学習特徴を読み取ることができた。例えば、粒度の細かい場合のクラスターGとHからは、それぞれ最初のコンテンツのみよく取り組む、後半のコンテンツのみよく取り組む、といった粒度が粗かった時には見えなかった結果が得られた。これにより、学習者のコンテンツに取り組む姿勢の変移を把握することができた。



1-10: 10章からなる講座の各受講回数

11-14: 4つのアンケートの各回答回数

15-17: 章末テスト1の初回点数,最高点数,受験回数

18-20: 章末テスト2の初回点数,最高点数,受験回数

21-23: 章末テスト3の初回点数,最高点数,受験回数

24-26: 終了テストの初回点数、最高点数,受験回数

27: 模擬試験受験回数

図 5 各クラスターのレーダーチャート (粒度の細かい場合)

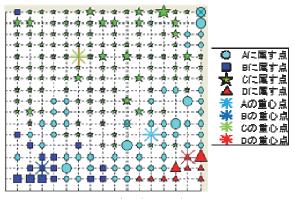
また、特に章末テストによく取り組む、講義よりテスト に力を入れて取り組む、といった特徴は粒度の粗い場合に も得られてもよい特徴であるが、粒度の粗い解析結果から は、このような特徴は見られなかった。よって、変数の粒 度を上げたことでより明確に特徴抽出を行うことができる といえる.

表 2 各クラスターの学習特徴と人数(粒度の細かい場合)

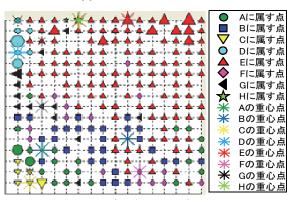
| Cluster | 学習特徴 | 人数 (人) | 割合 (%) |
|---------|---|-----------|-----------|
| A | 全体的によく取り組むが,特に模擬試験を繰り返し取り 組む | 24 | 10.1 |
| В | 特に章末テストに取り組む | 29 | 12.2 |
| C | 全講義を繰り返し受講し, テ ストも何回も取り組む | 7 | 2.9 |
| D | 全体的にコンテンツにしっ かり取り組まない | 11 | 4.6 |
| Е | 講義・テスト共に繰り返し取 り組まないが, テスト点数は 高い | 134 | 56.3 |
| F | 講義よりテストに力を入れ て取り組む | 15 | 6.3 |
| G | 講義・テスト共に最初のコン テンツは頑張るが、後半にな るにつれて取り組まない | 14 | 5.9 |
| Н | 全体的にあまり取り組まないが、章末テストに関しては 部分的によく取り組む | 4 | 1.7 |

3.2.3 SOM を用いたクラスターの可視化

続いて、SOM を用いてデータ分布の様子を把握する. マップを生成するにあたり、マップサイズを 16×16 、繰り返し回数を 1500 回、学習率を 0.1 として解析を行った. また、



(a) 粒度の粗い場合



(b) 粒度の細かい場合

図6 SOMによる学習履歴データの可視化

3.2.1 の結果で得られた各学習者の所属クラスターが、ど こにマッピングされたかわかるように、SOM の結果の上に、 クラスターごとに異なるマークと色で示し,各クラスタの 重心を「*」で示した. 更に、同じユニットに同クラスタ ーに属するデータが重なる場合はプロット点を大きくし, 異なるクラスターに分類されるにもかかわらず、データが 同じユニットに重なる場合は「×」印を付けた. 結果を図 6 に示す. 図 6(a)を見ると, クラスターA については, 性 質が最も異なるとされる(2次元空間において最も距離が 離れている)対角線の位置にも分布していることがわかる. 一方,図 6(b)では、クラスター数が増えたことも要因だと 考えられるが、全体的にデータがクラスターごとにまとま って分布する結果が得られた. このことから,変数の粒度 を細かくして解析を行ったことで、類似した特徴を持つデ ータがよりまとまるようにクラスタリングされたことが, SOM の結果より把握することができる.

4. まとめ

本報告では、企業が提供しているeラーニングコンテンツから取得した LMS に蓄積されている学習履歴データを用いて学習者の特徴を抽出するために、定量的な複数の指標を用いてクラスター数を決定し、各クラスターの学習特徴を抽出、可視化する方法を用いて解析を行った.更に、蓄積学習履歴データの変数の粒度を変化させて解析したときの学習特徴の抽出を行い、比較を行った.この結果より、変数の粒度の細かい方が、より明確な学習特徴を得られることが確認できた.また、SOM による可視化の結果においても、粒度が細かい方が、よりクラスターごとにまとまって分布する結果が得られた.これらの結果より、変数の粒度は解析結果の精度に影響を及ぼすことがわかった.しかし、学習履歴データを蓄積するサーバーの容量には制限があるため、どの変数をどの粒度で用いればよいかは、目的に応じて粒度を変化させて解析を行う必要がある.

謝辞 本研究で使用した解析データは、株式会社プロシーズからご提供頂いたものである.ここに感謝申し上げる.また、本研究は、科学研究費補助金基盤研究(C)(課題研究番号 25330419)の助成を受けたものである.

参考文献

- 1) 石川晶子,小川賀代,ピトヨ ハルトノ:学習履歴データを活用した学習者の特徴抽出方法の検討,教育システム情報学会誌, Vol. 31, No. 2, pp. 185-196 (2014).
- 2) Dunn, J. C.: A fuzzy relative of the ISODATA process and its use in detecting compact well-separated clusters, Cybernetics and Systems, Vol. 3, No. 3, pp. 32-57(1973).
- 3) Davies, D. and Bouldin, D.: A cluster separation measure, IEEE Trans. Pattern Analysis and Machine Intelligence, Vol. PAMI-1, No. 2, pp. 224-227(1979).
- 4) Schwarz, G. :Estimating the dimension of a model, Annals of Statistics, Vol. 6, No. 2, pp. 461-464(1978).