

1. はじめに

教材のデジタル化に伴い、電子教材上での学生の学習行動ログを大規模に収集することが可能となった。この学習行動ログデータを分析することで、学生一人一人に合わせた学習サポートの実現が期待されている。先行研究では、行動ごとの出現回数をヒストグラム特徴とした各学生の成績予測手法が提案されている。成績予測は最終的な成績を予測するものであり、学習行動の順序に起因する異常行動や特異性を捉えることは困難である。そこで本研究では、多くの学生が行う標準的な学習行動から逸脱した行動系列を異常と定義し、学習行動の文脈情報を捉える埋め込み表現による教師なし異常検知手法を提案する。

2. 先行研究

小濱らは、電子教材における行動ごとの出現回数をヒストグラム特徴として表現して成績予測を行う手法を提案した[1]。この手法は成績予測に有効である一方、行動の出現回数に依存した予測のため、学習行動における行動の違いや文脈的特異性を捉えることが難しく、学習プロセスの分析には課題がある。

3. 提案手法

本研究では、学習行動における文脈的な異常を捉えるため、行動系列の文脈情報を埋め込みベクトル化し、Isolation Forest[2]を用いて異常行動を検知する手法を提案する。本手法は、Masked Language Model (MLM) によるモデル学習、ベクトル生成、異常検知の3段階で構成される。提案手法の概要図を図1に示す。

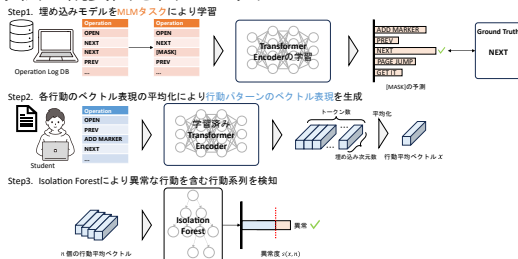


図1: 提案手法の概要図

3.1. 行動系列の埋め込みベクトル生成

Step1では、系列内の離れた行動間の関係性を捉えるためTransformer Encoderを採用し、一部をマスクした行動系列を入力としてMLMにより学習を行う。Step2では、学習したモデルに行動系列を入力し、各行動に対応する文脈情報を保持した埋め込みベクトルを得る。このベクトルを平均化することで、行動平均ベクトルを生成する。

3.2. Isolation Forest による異常行動の検知

Step3では、生成したベクトルをIsolation Forestに入力し、正解ラベルを用いずに特徴空間上で孤立した異常な行動パターンを持つ行動系列を検知する。Isolation Forestは、決定木に基づく教師なし異常検知手法である。学習時は、ランダムに選択された特徴量と、その最大値・最小値の間からランダムに決定された分割点を用いてデータを再帰的に二分割し、多数の決定木を構築する。各データが葉ノードに孤立するまでの平均パス長に基づき、式(1)で異常度を算出する。

$$s(x, n) = 2^{-\frac{E(h(x))}{c(n)}} \quad (1)$$

ここで、 $h(x)$ はデータ x が葉ノードに到達するまでのパス長、 $E(h(x))$ はその平均値、 $c(n)$ はデータ数 n における平均パス長である。正常データは孤立までに多くの分割を要する一方、異常データは少ない分割で孤立するためパス長が短くなる傾向がある。この性質を利用して異常行動を検知する。

4. 評価実験

本実験では、提案手法を用いて異常行動パターンを検知する。また、検知したデータの行動履歴を用いてその行動の分析を行う。

4.1. 実験条件

本実験では、九州大学で収集された学習行動ログデータセットを用いる。訓練データに2019年から2021年（1,209名）までのデータを、評価データに2022年（237名）のデータを使用し、いずれも1から8週目までの講義時間内のデータに限定して実験を行う。モデルは、BERT-Baseの構成（入力512トークン、次元数768、12層）を基にしたTransformer Encoderを使用する。学習には、学習率を $1e-5$ 、バッチサイズを32、エポック数を300に設定した。損失関数には、各行動の希少性に基づく重みであるIDFを適用したWeighted Cross Entropy Lossを用いる。また、異常判定の閾値は、異常度の上位5%に設定する。

4.2. 実験結果

異常検知の結果を表1に示す。表1より、全データ数1,463件に対し、74件が異常と判定された。正常データの平均行動回数は約152.2回であるのに対し、異常データは平均30.6回と、系列長が短い傾向が見られた。また、成績分布を確認すると、異常データ群では成績FやDの割合が高い傾向が見られる。一方、成績Aの学生も異常データの約3割（28.4%）を占めている。これは、標準的な学習行動とは異なる特異な行動パターンを持つ週や、遅刻や早退によって学習時間が極端に短くなり、行動回数が少なくなった週が検知されたためと考えられる。

表1: 異常度（平均パス長）に基づく検知結果

判定結果	データ数	割合 [%]	平均行動回数	成績分布 [%]				
				A	B	C	D	F
正常	1,389	94.9	152.2	38.7	27.1	19.8	9.8	4.6
異常	74	5.1	30.6	28.4	16.2	23.0	21.6	10.8

4.3. 定性的評価

行動平均ベクトルの異常度の分布を図2に、各クラスターを色で表示した結果を図3に示す。図2より、異常度が高いデータが空間の左端に密集し、独立した領域を形成していることがわかる。この領域を図3のクラスタリング結果と照合すると、行動回数が少なく、“OPEN”（教材を開く）や“CLOSE”（教材を閉じる）といった行動が主体であるクラスターに対応している。これは、システムへのアクセス後に内容を読み進めることなく即座に離脱した行動を示しており、多くの学生に見られるページ遷移の文脈を持たないため、特徴空間上で孤立し異常と判定されたと考えられる。一方、正常データに対応するクラスターは、行動回数が多く、“NEXT”（次のページに移動）を主体としたページ遷移を多く伴う行動系列であることが確認された。

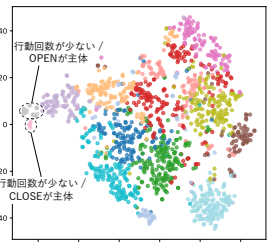
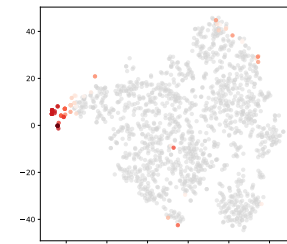


図2: 異常度の空間分布

図3: クラスタ分布

5. おわりに

本研究では、文脈を考慮した行動特徴とIsolation Forestを用いた異常行動の検知および分析を行った。実験の結果、異常データは即座離脱による文脈欠如が多く、成績F・Dの割合が高かったのに対し、正常データはページ遷移を多く伴い、成績Aの割合が高かった。このことから、ページ遷移を伴う継続的な学習行動が成績向上に寄与すると考えられる。今後は、ページ滞在時間などの時間情報を加え、閲覧行動の質的な差異を考慮した異常検知を目指す。

参考文献

- [1] H. Kohama, et al., “Recommending Learning Actions Using Neural Network”, ICCE, 2023.
- [2] F. Liu, et al., “Isolation Forest”, ICDM, 2008.