# Exploring Common and Individual Characteristics of Students via Matrix Recovering

**Zhen Wang,**[1,2] **Ben Teng,** [1,3] **Yun Zhou** [1,3] **Hanshuang Tong,** [1,3] **Guangtong Liu** [1,3]

[1] Aixuexi Education Group
[2] wangzhenYSU@163.com
[3] {tengben0905,zhouyun.nudt,tonghanshuang.thu, feidieliuliuliu}@gmail.com

## Abstract

Balancing group teaching and individual mentoring is an important issue in education area. The nature behind this issue is to explore common characteristics shared by multiple students and individual characteristics for each student. Biclustering methods have been proved successful for detecting meaningful patterns with the goal of driving group instructions based on students' characteristics. However, these methods ignore the individual characteristics of students as they only focus on common characteristics of students. In this article, we propose a framework to detect both group characteristics and individual characteristics of students simultaneously. We assume that the characteristics matrix of students' is composed of two parts: one is a low-rank matrix representing the common characteristics of students; the other is a sparse matrix representing individual characteristics of students. Thus, we treat the balancing issue as a matrix recovering problem. The experiment results show the effectiveness of our method. Firstly, it can detect meaningful biclusters that are comparable with the state-of-the-art biclutering algorithms. Secondly, it can identify individual characteristics for each student simultaneously. Both the source code of our algorithm and the real datasets are available upon request.

## Introduction

A growing collection of educational data contributes to the research of modeling student characteristics. For instance, by exploiting the exercising records of students with knowledge tracing or knowledge diagnosis methods, researchers can track the change of each student's knowledge acquisition during their exercising activities, and output a student-knowledge mastery matrix, whose element presents the performance level of knowledge mastery for each student (Piech et al. 2015; Liu et al. 2019; Tong, Zhou, and Wang 2020).

Based on the characteristics of students, instructors can offer each student specific interventions to improve their performance (Lin-Siegler, Dweck, and Cohen 2016). However, it is too time-consuming for instructors to complete these tasks by hand, since a large number of students usually vary greatly in learning rates and knowledge levels. One feasible approach is conducting a clustering anaylsis on the student-knowledge mastery matrix to solve this challenge, which is

the practice of placing students of similar characteristics in the same group (M.A 2015; Romero and Ventura 2010). After discovering student groups according to students' common characteristics, personalized learning systems can be built and adaptive contents can be offered to promote effective group learning by instructors.

To date, serveral studies have been published to cluster students into meaningful groups based on their characteristics with a goal of driving group instructions (Amershi and CONATI 2009; Dutt, Ismail, and Herawan 2017; Mojarad et al. 2018; Henriques, Finamore, and Casanova 2019). These methods can be mainly divided into two major categories: traditional clustering methods and biclustering methods. Traditional clustering methods, such as hierarchical agglomerative clustering, K-means and model-based clustering, identify groups of students with similar characteristics in a global way, that is they simply group students according to all available values (all knowledge mastery levels in this article), thus being unable to identify local patterns. While biclustering algorithms, whose particularity is that partitioning is done in two dimensions of a matrix yielding to clustering according to students and their characteristics, allow the discovery of local patterns. (Mojarad et al. 2018) is one of typical algorithms based on traditional clustering methods to group students. Mean-shift clustering is used to select a number of clusters, and then k-mean clustering is applied to identify distinct student profiles. (Henriques, Finamore, and Casanova 2019) is one recent application of biclustering method in educational data. The authors use pattern-based biclustering approach to detect non-trivial, yet potentially relevant educational and statistically significant patterns from the performance of students' data. Their results confirm the unique role of biclustering in finding relevant patterns of students' performance.

Despite the advances of researches on student grouping, one important problem that students maybe do not have the same characteristics in some aspects even though they're in the same group is always ignored. Balancing group teaching and individual mentoring is not trivial in the education area. To our knowledge, there has been no existing computional method to solve this issue. In this article, we propose a novel method to solve this problem. Specifically, we use the method to identify both group characteristics and individual characteristics of students simultaneously in this

article. Under the assumption that common characteristics shared by multiple students form a low-rank matrix and individual characteristics for each student form a sparse matrix, we model the problem as a matrix recovering problem. Then we use an iterative algorithm to solve it.

To demonstrate the performance of the proposed method, we conduct comparison experiments using both synthesized datasets and real educational datasets. Simulation results show that the proposed method achieves comparable performances compared with existing methods in many settings. And experiment results on two real datasets demonstrate the effectiveness of our method.

Overall, the salient features of the method described in this article can be summarized as follows:

- it inherits the advantage of biclustering that local pattern can be discovered.

- the individual characteristics for each student can be detected meanwhile.

- statistical evaluation is employed to filter the spurious biclusters, which guarantees the statistical significant of the results.

## Method

In this paper, we propose a new framework for exploring common and individual characteristerics of students.
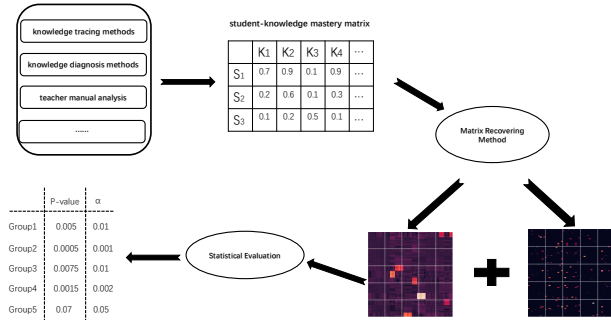


Figure 1: The framework of our method. It contains two steps: matrix recovering and statistical evaluation. Following above two steps, the method proposed can be very robust and accurate in the detection of common characteristics and individual characteristics for students.

Figure 1 gives the flow of the overall framework. There are two key steps in the framework. In the first step, we perform matrix recovering method on the students' characteristics from different sources, such as knowledge tracing methods and manual analysis by instructors. After that, statistical evaluation is used to detect robust biclusters. The result validation step is aimed to guarantee the statistic significance for each discovered clusters, as stated in (Henriques and Madeira 2018). In the following, we will explain these two steps in detail.
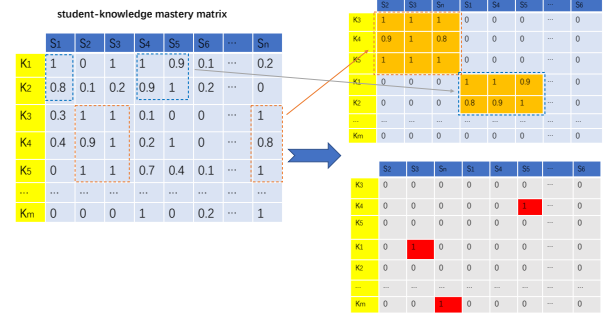


Figure 2: An example of matrix recovering for educational data. Base on the assumption that common characteristics form a low-rank matrix and individual characteristics form a sparse matrix, we can decompose the input matrix into two matrices.

## Matrix recovering

The goal of the first step is to recover low-rank component and sparse component from original input matrix, respectively. Figure 2 shows an example result of matrix recovering for educational data. In this subsection, we first present the mathematical formulation of matrix recovering problem. And then the solution to solve it is provided. Finally, we give some suggestions to select appropriate parameters.

**Formulation** Mathematically, we can express the characteristics of students as a matrix $\mathbf{D} \in \mathbb{R}^{n \times m}$, where each element $d_{ij}$ of the matrix is a value representing the mastery of knowledge for each student, and $n$ and $m$ are the numbers of students and knowledge topics in this article, respectively. Because our goal is to balance group teaching and individual mentoring, we need to be able to explore common characteristics shared by multiple students and individual characteristics for each student. As common characteristics of students can be represented as biclusters and individual characteristics for each student can be assumed as randomly distributed and sparse in the matrix, we can treat the balancing issue as a problem of recovering a low-rank matrix $\mathbf{X}$ and a sparse matrix $\mathbf{E}$ from the orginal characteristics matrix $\mathbf{D}$.

Naturally, the following matrix decomposition model is proposed to detect two types of characteristics from input with noise:

$$\mathbf{D} = \mathbf{X} + \mathbf{E} + \epsilon, \tag{1}$$

In Eq. (1), $\mathbf{X}$ refers to the common characteristeris component. $\mathbf{E}$ refers to the individual component and $\epsilon$ is a noise component.

We consider the following minimization problem to achieve the decomposition:

$$\min_{\mathbf{X},\mathbf{E},\epsilon} \frac{1}{2} \|\epsilon\|_F^2 + \alpha \text{rank}(\mathbf{X}) + \beta \|\mathbf{E}\|_0$$
$$\text{s.t. } \mathbf{D} = \mathbf{X} + \mathbf{E} + \epsilon, \tag{2}$$

where $\|\epsilon\|_F$ is the Frobenius norm, rank $(\mathbf{X})$ is the rank of matrix $\mathbf{X}$ and $\|\mathbf{E}\|_0$ is the $\ell_0$-norm. We can get a penalized

maximum likelihood estimate with respect to the variables $\mathbf{X}, \mathbf{E}, \epsilon$ through solving Eq. (2).

Since the model proposed in Eq.(2) is NP-hard, the convex relaxation approach is used to effectively recover $\mathbf{X}$ and $\mathbf{E}$. Specifically, the rank $(\cdot)$ is replaced by the nuclear norm and the $\ell_0$-norm is replaced by the $\ell_1$-norm. The nuclear norm is defined as the sum of the singular values of $\mathbf{X}$, whic is the tightest convex surrogate to the rank operator (Fazel 2002) and has been widely used for low-rank matrix recovery (Candès et al. 2011). The $\ell_1$-norm is defined as $\|\mathbf{X}\|_1 = \sum_{i,j} |X_{ij}|$. The $\ell_1$ relaxation has proven to be a powerful technique for sparse signal recovery (Tropp 2006).

Thus, we can solve the following problem, instead of directly solving Eq.(2):

$$\mathcal{F}(X, E) = \min_{\mathbf{X}, \mathbf{E}} \frac{1}{2}\|\mathbf{D} - \mathbf{X} - \mathbf{E}\|_F^2 + \alpha\|\mathbf{X}\|_* + \beta\|\mathbf{E}\|_1. \quad (3)$$

Eq.(3) is a convex problem so that the global optimal solution is unique. This means that the results of our method is stable.

**Solution**  We can solve the optimization problem of Eq.(3) by alternatively solving the following two sub-problems until convergence:

$$\hat{\mathbf{X}} \leftarrow \arg\min_{\mathbf{X}} \mathcal{F}(\mathbf{X}, \hat{\mathbf{E}}) \quad (4)$$

$$\hat{\mathbf{E}} \leftarrow \arg\min_{\mathbf{E}} \mathcal{F}(\hat{\mathbf{X}}, \mathbf{E}). \quad (5)$$

The theoretical proof for the convergence can be found in (Boyd 2010). The problem in Eq.(4) can be reduced to

$$\min_{\mathbf{X}} \frac{1}{2}\|\mathbf{D} - \hat{\mathbf{E}} - \mathbf{X}\|_F^2 + \alpha\|\mathbf{X}\|_*, \quad (6)$$

which becomes a nuclear-norm regularized least-squares problem and has the following closed-form solution (Cai, Candès, and Shen 2010),

$$\hat{\mathbf{X}} = \mathcal{D}_\alpha\left(\mathbf{D} - \hat{\mathbf{E}}\right), \quad (7)$$

where $\mathcal{D}_\lambda$ refers to the singular value thresholding (SVT)

$$\mathcal{D}_\lambda(\mathbf{M}) = \sum_{i=1}^{r} (\sigma_i - \lambda)_+ \mathbf{u}_i \mathbf{v}_i^T. \quad (8)$$

Here, $(x)_+ = \max(x, 0)$. $\{\mathbf{u}_i\}, \{\mathbf{v}_i\}$, and $\{\sigma_i\}$ are the left singular vectors, the right singular vectors, and the singular values of $\mathbf{M}$, respectively.

The problem in Eq.(5) can be rewritten as

$$\min_{\mathbf{E}} \frac{1}{2}\|\mathbf{D} - \hat{\mathbf{X}} - \mathbf{E}\|_F^2 + \beta\|\mathbf{E}\|_1. \quad (9)$$

It admits a closed-form solution

$$\hat{\mathbf{E}} = \mathcal{S}_\beta\left(\mathbf{D} - \hat{\mathbf{X}}\right), \quad (10)$$

where $\mathcal{S}_\beta(\mathbf{M})_{ij} = \text{sign}(M_{ij})(M_{ij} - \beta)_+$ refers to the elementwise soft-thresholding operator (Boyd 2010).

**Parameter selection**  Two parameters need to be estimated in the first step. In this article, we give some suggestions to select proper estimations via the analysis of the size of the input matrix $(n, m)$ and the standard variation of the noise $\sigma$ (Candès et al. 2011; Zhou et al. 2010).

Firstly, we estimate an intermediate variable $\sigma$ from the data by the median-absolute-deviation estimator (Meer et al. 1991)

$$\hat{\sigma} = 1.48 \text{ median} \{|\mathbf{D} - \text{median}(\mathbf{D})|\}. \quad (11)$$

For parameter $\alpha$, it serves as a threshold in the SVT step in Eq.(8) so that it should be large enough to threshold out the noise but not too large to over-shrink the signal (Zhou et al. 2010). A proper value is $\alpha = (\sqrt{n} + \sqrt{p})\sigma$, which is the expected $\ell_2$-norm of a $n \times p$ random matrix with entries sampled from $\mathcal{N}(0, \sigma^2)$. In practice, we can adjust $\alpha$ around this value to fit real data.

For parameter $\beta$, there is a relative weight $\lambda = \beta/\alpha$ that balances the two terms in $\alpha\|\mathbf{X}\|_* + \beta\|\mathbf{E}\|_1$ and consequently controls the rank of $\mathbf{X}$ and the sparsity of $\mathbf{E}$. It has been proved that $\lambda = 1/\sqrt{m}$ gives a large probability of recovering $\mathbf{X}$ and $\mathbf{E}$ under their assumed conditions (Candès et al. 2011). In practice, we should set different values of $\lambda$ to keep sufficient characteristics in $\mathbf{X}$ in specific applications.

## Statistical evaluation

When the low-rank matrix $X$ is recovered, we can get the biclusters in two ways. The first way is to perform biclustering methods on the low-rank matrix, rather than on the original matrix. The other one is similar to Spectra Biclustering (Kluger et al. 2003), by applying traditional clustering methods on the low-rank matrix to cluster students and knowledge topcis, respectively. In order to compare our framework with biclustering methods in the experiments section, we use the latter one to detect biclusters in this article. However, as small biclusters can have high levels of homogeneity by chance, some ones in the detected biclusters are spurious and insignificant. In order to filter false positive biclusters and control the false discovery rate, statistical evaluation are needed.

In this article, we adopt the method BSig proposed in (Henriques and Madeira 2018) to evaluate the statistical significance of detected biclusters and reject those biclusters with high $P$-values. The BSig method provides the unprecedented possibility to minimize the number of false positive biclusters without incurring on false negatives. It first approximates a null model of the target educational data and then appropriately tests each bicluster in accordance with its underlying coherence. Finally, we reject those biclusters with $P-$values higher than Bonferroni correction thresholds.

## Experiments

To test the method proposed in this article comprehensively, we conduct several comparison experiments both on synthetic and real datasets with those biclustering algorithms drawn from important studies in the biclustering literature. In detail, we compare the following four typical biclustering methods on synthetic datasets:

1. **FABIA (Hochreiter et al. 2010).** This algorithm is based on factor analysis and it is a multiplicative model that assumes non-Gaussian signal distributions with heavy tails.

2. **LAS (Shabalin et al. 2009).** This method searches for patterns from input matrix by locally maximizing a Bonferroni based significance score iteratively.

3. **ISA (Ihmels, Bergmann, and Barkai 2004).** Iterative Signature Algorithm first randomly selects columns and rows, and then evaluates and updates them through iterative steps until convergence.

4. **Spectra Biclustering (Kluger et al. 2003).** Spectral biclustering method assumes that the input data matrix has a hidden checkerboard structure and uses singular value decomposition to find it.

For these methods, we use the default settings of them in biclustlib (Padilha and Campello 2017), a python library of biclustering algorithms. Meanwhile, we use six evaluation metrics in the experiments to compare the biclustering results: liu wang match score (Liu and Wang 2006), prelic recover score (Prelić et al. 2006), prelic relevance score (Prelić et al. 2006), csi (Campello 2010), cluster error (Patrikainen and Meila 2006) and fabia consensus score (Hochreiter et al. 2010). They give scores from different views by comparing the predicted biclusters against the ground-truth ones. Larger scores assigned by these metrics indicate better performances of the biclustering methods. The details of these metrics are referred to the corresponding papers.

Additionally, as our method can identify sparse signals at the same time, we validate the performance of identifying sparse signals using precision, recall and F1-score.

All the experiments are tested on the ThinkPad T480 Laptop computer with 1.80GHz CPU and 16G main memory.

### Experiments with synthetic data

Firstly, we test the effectiveness of our method on four synthetic biclustering type: constant, shift, scale and shift & scale. The biclustering datasets are generated based on the procedure proposed by (Eren et al. 2013). The details of the settings for each data type are listed in Table 1. After generating biclustering data, we add sparse signals to the matrix. Specifically, a sparse matrix $E$ is generated through the following way: each element of this sparse matrix independently takes on value 0 with probability 1 - $p_s$, and value 6 with probability $p_s = 0.01$. Finally, we shuffle the resulted matrix.

To be fair in the performance comparison, for each synthetic data type, we perform 5 independent runs to obtain an average test result for each method. The detailed results are presented in Figure 3. As Figure 3 shows, there is no perfect algorithm that performs best on all the synthetic datasets. That is because each biclustering model is always designed to fit a specific data type.

Specifically, for LAS method, it performs best on the constant biclusters, while it can not find shift and scale biclusters effectively. FABIA method can perform competitive with our method on shift and scale datasets. Howerver, it performs worse on constant datasets. Overall, our method can detect well biclusters on each dataset.

Table 1: **The detail of settings for each synthetic data.**

| type | setting |
| --- | --- |
| Constant biclusters | nrows=300, ncols=50, nclusts=5, bicluster_signals=5, nclustcols=8, nclustrows=5, noise=1, shuffle=True |
| Shift biclusters | nrows=300, ncols=50, nclusts=5, bicluster_noise=[0.01]*5, noise=1, base_loc=1, shift_loc=1, shift_scale=3, shuffle=True |
| Scale biclusters | nrows=300, ncols=50, nclusts=5, bicluster_noise=[0.01]*5, base_loc=1, scale_scale=3, scale_loc=0, shuffle=True |
| Shift-scale biclusters | nrows=300, ncols=50, nclusts=5, bicluster_noise=[0.01] * 5, base_scale=1, scale_scale=2, shift_scale=3, shuffle=True |

We also test the ability of our method to identify sparse signals. Figure 4 shows the performance of our method on 4 sythetic datasets. Our method can achieve great F1-scores. Note that the recall is almost 1 no matter what type the data is. It guarantees that every sparse signal will not be missed, meaning that each student's indivial characteristics will not be ignored in practice.

In Figure 5, we present illustrations of the orignal matrix, low-rank matrix and sparse matrix recovered by our method for each dataset.

### Experiments with real data

Since synthetic data can only reflect certain aspects of reality, we also perform our framework on two real datasets to test the effectiveness. In this subsection, we will first descripe the detail of the two datasets and then present the performances of our method.

**Datasets and Preprocessing** We collected two datasets of students' performance: ADS data (Henriques, Finamore, and Casanova 2019) and MATH data. The ADS dataset is provided by the Department of Informatics of the Pontifical Cathaolic University of Rio de Janero (PUC-Rio), which captures the performance of students along the topics of the Advanced Data Structures (ADS) course. The MATH dataset is collected from AIXUEXI Education Group LTD, which contains 76 students in Junior High School Grade 1 and 54 knowledge topics.

We do some preprocessings on the two datasets. For ADS data, as the goal is to detect the weaknesses of students, we first use 100 minus the performance scores and divide the resulted scores into 10 levels where 0 denotes an excelling grade and 9 a low grade. For MATH data, we first use the DKT method (Piech et al. 2015) to get the mastery matrix and then use 1 minus the elements in the matrix as the input data.

**Results** The results of our method on ADS dataset is shown in Figure 6. For the resulted sparse component,
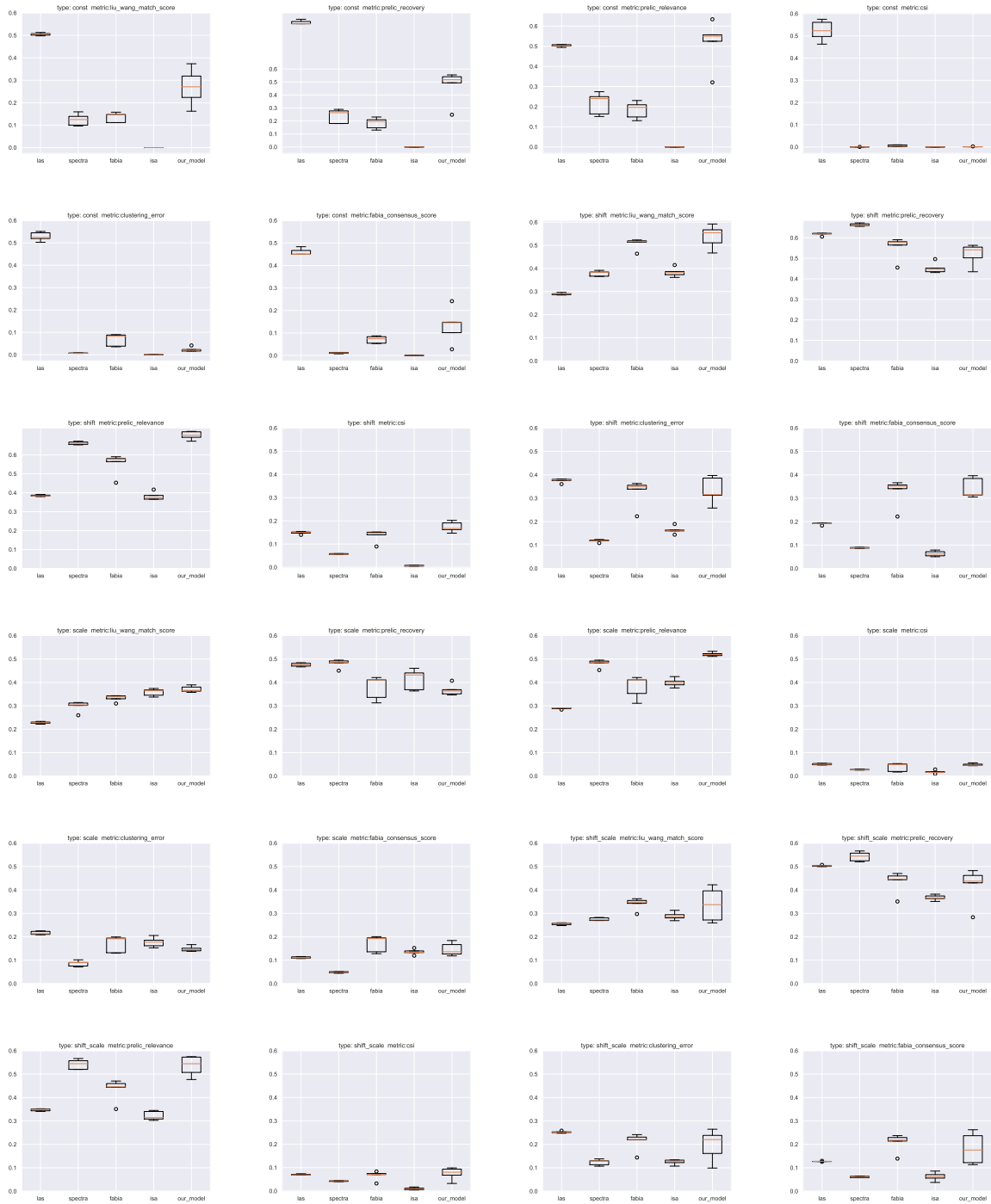
Figure 3: The results of synthetic data experiments. We compare our framework with 4 typical biclustering algorithms on 4 synthetic data type using 6 metrics.
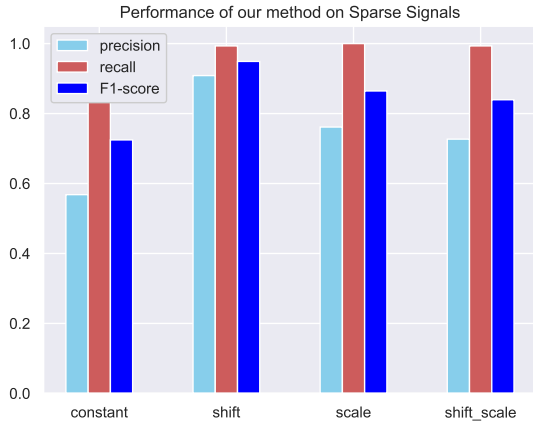
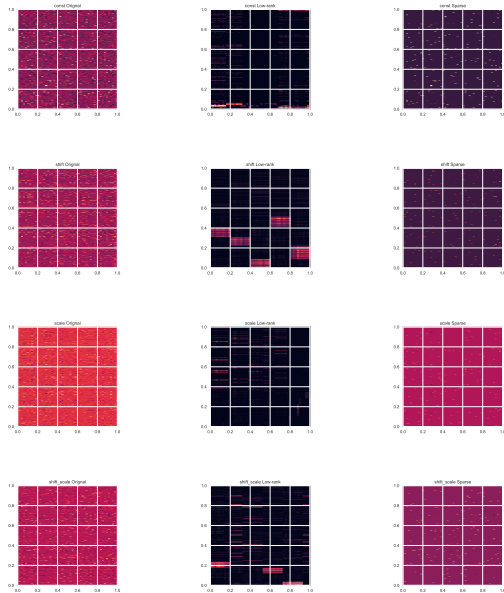Figure 4: Performance of our framework on Sparse Signals.



Figure 5: Illustrations of the results on synthetic data. The left ones are the original input matrix. The middle and the right ones are the results recovered by our framework.

the sparse rate is approximately 0.05. For the low-rank part, there are approximately 20 students and 5 knowledge topics per bicluster on average. We can find that the framework proposed in this article can efficiently and comprehensively detect student-topic biclusters. Table 2 gives the $P-$values of some biclusters found in ADS dataset by our method. The filter step of the framework ensures the statistic significant of the biclusters detected. Figure 7 presents the patterns of four biclusters recovered in the low-rank matrix. Here, four patterns are divided into 2 knowledge topic sets: {GraphTraversal, Gridfile, Complexity} and {Árvorebinária, TiposEstruturados(incluivetordeestruturas), AVLTrees}. We can find that different biclusters follow different patterns. For the left two biclusters, students in the upper one maybe have difficulties in learning these three topics, while students in the lower one may be not good at {Árvorebinária, AVLTrees}, but for TiposEstruturados(incluivetordeestruturas), they are not the same as students in the upper one. This phenomenon is more obvious in the right two biclusters. Thus, based on the different patterns, instructors can capture students' weak knowledge and give specific interventions to each group. We can get similar results from MATH dataset, as the patterns shown in Figure 8.

Furture more, we do some analysis on the recovered low-rank matrix. Inspired by Latent Semantic Analysis, we do Singular Value Decomposition on the recovered low-rank matrix. The right singular vectors can be regard as latent features for each topic. As shown in Figure 9, it is clear to see that some knowledge topics are clustered based on these features. This is in accord with intuition and practical. First, some topics may have the similar difficulties so that most students cannot master them easily. Secondly, for students, to master some knowledge topcis is highly relied on the mastery of same prior knowledge. Low proficiency on prior knowledge results in the fact that the following topics show a similar degree of mastery. This result gives us a new perspective to find the relationships among topics and enlightens us that well-defined relationships among knowledge topics could be useful in the detection of meaningful groups.

## Conclusion

Balancing group teaching and individual mentoring is an important topic in education area. In this article, we propose a matrix recovery based method for detecting common characteristics for a group of students and identifying individual ones for each student. In addition, statistical evaluation method is applied to filter the spurious biclusters. The experiment results show that the method can provide more stable biclusters and accurate sparse signals. As shown in our experimental results, some relationships among knowledge topics are detected. Inspired by this, we can integrate information from knowledge graph to identify more meaningful groups in the future work.

Table 2: $P-$values of some biclusters found in ADS dataset.

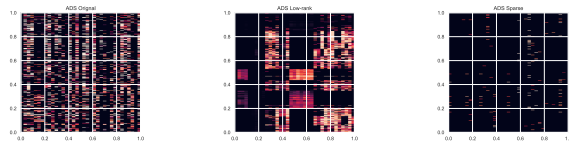| Students' ID | Knowledge Topics | P-value |
|---|---|---|
| 16 students: [34 74 82 ... 168 174 175] | ListasEncadeadas, Bitvector, GenericTrees, B-TreesInsertion | 7.1E-28\|3.8E-5 |
| 34 students: [12 20 31 ... 251 254 261] | Árvorebinária, TiposEstruturados, AVLTrees | 4.8E-5\|4.0E-4 |
| 14 students: [3 19 30 ... 117 124 144] | GraphTraversal, Gridfile, Complexity | 1.1E-13\|2.8E-5 |
| 12 students: [9 86 128 ... 244 245 281] | ListasEncadeadas, Bitvector, GenericTrees, B-TreesInsertion | 1.6E-22\|3.8E-5 |



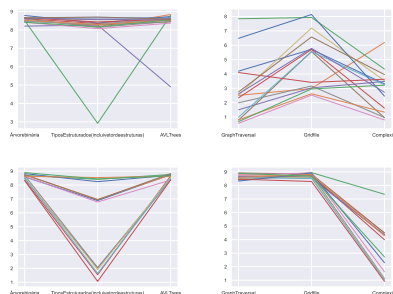Figure 6: The results of ADS data experiments.



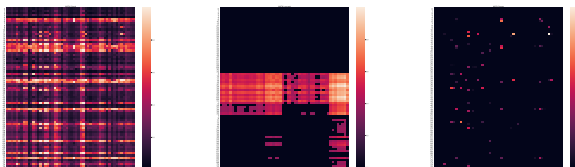Figure 7: The patterns of four biclusters recovered in ADS dataset



Figure 8: The results of MATH data experiments.
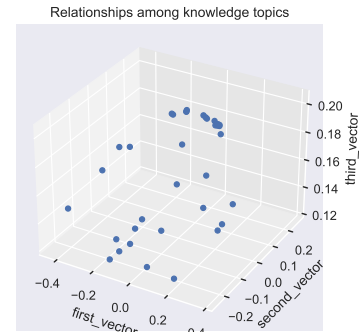


Relationships among knowledge topics

Figure 9: Analysis on the first three of the right singular vectors. We can regard the elements in the first three vectors as the important features of the corresponding topics. Clearly, some topics are grouped together based on these features.

## References

Amershi, S.; and CONATI, C. 2009. Combining Unsupervised and Supervised Classification to Build User Models for Exploratory Learning Environments. *Journal of Educational Data Mining* .

Boyd, S. 2010. Distributed Optimization and Statistical Learning via the Alternating Direction Method of Multipliers. *Foundations and Trends® in Machine Learning* 3(1): 1–122.

Cai, J.; Candès, E.; and Shen, Z. 2010. A Singular Value Thresholding Algorithm for Matrix Completion. *SIAM Journal on Optimization* 20: 1956.

Campello, R. J. G. B. 2010. Generalized external indexes for comparing data partitions with overlapping categories. *Pattern Recognition Letters* 31(9): 966–975.

Candès, E.; Li, X.; Ma, Y.; and Wright, J. 2011. Robust principal component analysis? *Journal of the ACM* 58(3): 11.

Dutt, A.; Ismail, M. A.; and Herawan, T. 2017. A Systematic Review on Educational Data Mining. *IEEE Access* 5: 15991–16005.

Eren, K.; Deveci, M.; Küçüktunç, O.; and Çatalyürek, Ü. V. 2013. A comparative analysis of biclustering algorithms for gene expression data. *Briefings Bioinform.* 14(3): 279–292.

Fazel, M. 2002. *Matrix rank minimization with applications*. Ph.D. thesis, Stanford University.

Henriques, R.; Finamore, A. C.; and Casanova, M. A. 2019. On the Discovery of Educational Patterns using Biclustering. In *Intelligent Tutoring Systems - 15th International Conference, ITS 2019, Kingston, Jamaica, June 3-7, 2019, Proceedings*, volume 11528 of *Lecture Notes in Computer Science*, 133–144. Springer.

Henriques, R.; and Madeira, S. C. 2018. BSig: evaluating the statistical significance of biclustering solutions. *Data Min. Knowl. Discov.* 32(1): 124–161.

Hochreiter, S.; Bodenhofer, U.; Heusel, M.; Mayr, A.; Mitterecker, A.; Kasim, A.; Khamiakova, T.; Sanden, S. V.; Lin, D.; Talloen, W.; Bijnens, L.; Göhlmann, H. W. H.; Shkedy, Z.; and Clevert, D. 2010. FABIA: factor analysis for bicluster acquisition. *Bioinform.* 26(12): 1520–1527.

Ihmels, J.; Bergmann, S.; and Barkai, N. 2004. Defining transcription modules using large-scale gene expression data. *Bioinformatics* 20(13): 1993–2003. ISSN 1367-4803.

Kluger, Y.; Basri, R.; Chang, J. T.; and Gerstein, M. 2003. Spectral Biclustering of Microarray Data: Coclustering Genes and Conditions. *Genome Research* 13(4): 703–716.

Lin-Siegler, X.; Dweck, C. S.; and Cohen, G. L. 2016. Instructional interventions that motivate classroom learning. *Journal of Educational Psychology* 108(3): 295–299.

Liu, Q.; Huang, Z.; Yin, Y.; Chen, E.; Xiong, H.; Su, Y.; and Hu, G. 2019. EKT: Exercise-aware Knowledge Tracing for Student Performance Prediction. *IEEE Transactions on Knowledge and Data Engineering* 1–1.

Liu, X.; and Wang, L. 2006. Computing the maximum similarity bi-clusters of gene expression data. *Bioinformatics* 23(1): 50–56. ISSN 1367-4803.

M.A, I. 2015. Clustering Algorithms Applied in Educational Data Mining. *International Journal of Information Engineering and Electronic Business* .

Meer, P.; Mintz, D.; Rosenfeld, A.; and Kim, D. 1991. Robust regression methods for computer vision: A review. *International Journal of Computer Vision* 6(1): 59–70.

Mojarad, S.; Essa, A.; Mojarad, S.; and Baker, R. S. 2018. Data-Driven Learner Profiling Based on Clustering Student Behaviors: Learning Consistency, Pace and Effort. In *Intelligent Tutoring Systems - 14th International Conference, ITS 2018, Montreal, QC, Canada, June 11-15, 2018, Proceedings*, volume 10858 of *Lecture Notes in Computer Science*, 130–139. Springer.

Padilha, V. A.; and Campello, R. J. G. B. 2017. A systematic comparative evaluation of biclustering techniques. *BMC Bioinformatics* 18(1): 55.

Patrikainen, A.; and Meila, M. 2006. Comparing subspace clusterings. *IEEE Transactions on Knowledge and Data Engineering* 18(7): 902–916.

Piech, C.; Bassen, J.; Huang, J.; Ganguli, S.; Sahami, M.; Guibas, L. J.; and Sohl-Dickstein, J. 2015. Deep Knowledge Tracing. In *Advances in Neural Information Processing Systems 28*, 505–513. Curran Associates, Inc.

Prelić, A.; Bleuler, S.; Zimmermann, P.; Wille, A.; Bühlmann, P.; Gruissem, W.; Hennig, L.; Thiele, L.; and Zitzler, E. 2006. A systematic comparison and evaluation of biclustering methods for gene expression data. *Bioinformatics* 22(9): 1122–1129.

Romero, C.; and Ventura, S. 2010. Educational Data Mining: A Review of the State of the Art. *IEEE Trans. Syst. Man Cybern. Part C* 40(6): 601–618.

Shabalin, A. A.; Weigman, V. J.; Perou, C. M.; and Nobel, A. B. 2009. FINDING LARGE AVERAGE SUBMATRICES IN HIGH DIMENSIONAL DATA. *The Annals of Applied Statistics* 3(3): 985–1012.

Tong, H.; Zhou, Y.; and Wang, Z. 2020. Exercise Hierarchical Feature Enhanced Knowledge Tracing. In *International Conference on Artificial Intelligence in Education*, 324–328. Springer.

Tropp, J. A. 2006. Just relax: Convex programming methods for identifying sparse signals in noise. *IEEE Transactions on Information Theory* 52(3): 1030–1051.

Zhou, Z.; Li, X.; Wright, J.; Candes, E.; and Ma, Y. 2010. Stable principal component pursuit. In *Proceedings of the IEEE International Symposium on Information Theory*.