

# Better Prerequisite Discovery From Data: Implementation and Evaluation Issues

Leave Authors  
Anonymous  
for Submission  
City, Country  
e-mail address

Leave Authors  
Anonymous  
for Submission  
City, Country  
e-mail address

Leave Authors  
Anonymous  
for Submission  
City, Country  
e-mail address

## ABSTRACT

Understanding the prerequisites of skills is useful for designing curriculum, assessing mastery and for student modeling. For example, in student modeling, prerequisites can be used for remediation of student knowledge. Prerequisites are often encoded as graphs designed by subject matter expert in a costly and time-consuming process. In this paper, we introduce *Combined student Modelling and prerequisite Discovery* (COMMAND), a novel algorithm for inferring a prerequisite graph and a student model from data. Unlike prior work that evaluates prerequisite graphs using synthetic data, we develop a methodology to evaluate prerequisite graphs with data collected from students interacting with a digital learning environment. Our experiments suggest that COMMAND improves on previously published approaches.

## Keywords

Prerequisite discovery, Bayesian network, student modeling

## 1. INTRODUCTION

Students usually acquire skills in a meaningful sequence which starts from relatively simple concepts and gradually approaches more complex ones. Among these skills, some are preliminary of others such that they must be mastered before the subsequent concepts can be learned. For instance, students have to know how to do addition before they learn to do multiplication. In this paper, we use prerequisite structure to refer to the relationships among skills that place strict constraints on the order in which these skills can be acquired. Determining the prerequisite relations among skills is crucial for designing curriculum and for assessing mastery.

Prerequisite structures of skills are often hand-engineered by subject matter experts in a costly and time-consuming process. Moreover, the prerequisite structures specified by the experts are seldom tested and might be unreliable in the sense that experts may hold “blind spots”. Nowadays, large volume of educational data has been cumulated through the digital tutoring systems. Thus, learning prerequisite structures from educational data has drawn substantial interest from both education and data mining communities [7, 20, 3, 16,

5, 15]. However, inferring the prerequisite relationships between skills is still challenging since a student’s knowledge of a skill is an unobserved latent variable in the data. In this paper, we introduce *Combined student Modelling and prerequisite Discovery* (COMMAND), a novel pipeline for discovering prerequisite structure of skills from data. COMMAND models the prerequisite structure of the skills with a statistical model called Bayesian network. It then learns the Bayesian network through a two-stage process.

## 2. PREVIOUS WORK

Prior work have investigated the discovery of prerequisite structures from data. Many of these approaches discover the relationship of items without using latent variables. That is, the prerequisites do not use skill mappings, and only find dependencies between items [7, 20, 15]. **What’s wrong with finding dependencies between items?**

To account for latent skills, a mapping of items to skills (often called a  $Q$ -matrix) is used. In this setting, students are tested on items-questions, problems, parts of questions- each of which requires one or multiple background skills. To estimate the relationship between two background skills, Brunskill proposed a method to compute and compare the log likelihoods between the prerequisite model and the flat model (two skills are independent) to determine which model better fits the data [3]. Chen et al. [5] used DINA (Deterministic Input Noisy AND) model to infer the knowledge states of students from their response data on test items and then used probabilistic association rules mining to determine the prerequisite relationships between two skills. However, these methods are only capable of estimating the pairwise relationships, instead of optimizing the full structure of the skills.

Bayesian network is a promising tool to model skill topologies as demonstrated in [13, 11]. Bayesian network allows the modeling of full structure of skills and can encode conditional independence relationships between the skills. This enables more accurate modeling of the probabilistic dependence between skills. By modeling prerequisite relations as a Bayesian network, Scheines et al. converted the problem of prerequisite structure discovery to a Bayesian network structural learning problem, and proposed an algorithm to find Bayesian network structures that fit the data [16]. However, the output from their algorithm is a class of Bayesian networks, which are equivalent in representing a probability distribution. These *equivalent* Bayesian networks have the same skeleton but differ in the orientation of some edges, thus represent very different prerequisite structures. Scheines and his colleagues did not offer any solution to discriminate between these Bayesian networks.

In this paper, we propose the *Combined student Modelling and prerequisite Discovery* (COMMAND) pipeline to infer the prerequisite structures of skills and a student model from data. The pipeline also uses the promising Bayesian network model to represent the prerequisite relationships. The pipeline first uses a popular Bayesian network structure learning algorithm, called Structural EM to select a class of Bayesian network models based on the distribution fitting of the data. It then uses the domain knowledge that it is unlikely for a student to master a skill before mastering its prerequisite skills to select one Bayesian network, which gives a unique prerequisite structure. Different from [16], Structural EM converts the Bayesian network learning to an optimization problem. It is more scalable and can be applied to data that contain missing values. Further, the approach in [16] only outputs the topology of the Bayesian network, while Structural EM outputs both the structure and the probability distribution (a set of conditional probabilities), which is very useful for student modeling. We demonstrate the robustness of our proposed approach on noisy data that contains unaccounted latent variables or missing values. Finally, we show the discovered prerequisite models are useful for student modeling.

### 3. THE COMMAND PIPELINE

COMMAND learns the prerequisite structure of the skills from data with a statistical model called Bayesian network [14, 17]. Bayesian networks are also called probabilistic graphical models because they can be represented visually and algebraically as a collection of nodes and edges. A tutorial description of Bayesian networks in education can be found elsewhere [1], but for now we say that they are often described with two components: the nodes represent the random variables, which we describe using *conditional probability tables* (CPTs), and the set of edges that form a *directed acyclic graph* (DAG) represent the conditional dependencies between the variables. Bayesian networks are a flexible tool that can be used to model an entire curriculum.

Figure 1 illustrates an example of a prerequisite structure modeled with a Bayesian network. Here, we relate four test items with the skills of addition and multiplication. Addition is a prerequisite of multiplication thus there is an arrow from addition to multiplication. Modeling prerequisites as edges in a Bayesian network allows us to frame the discovery of the prerequisite relationships as the well-studied machine learning problem of learning a Bayesian network from data (with the presence of latent variables).

Algorithm 1 describes the COMMAND pipeline. Suppose we collect data from  $n$  students, answering  $p$  items. Then, the input to COMMAND is a matrix  $\mathbf{D}$  with  $n \times p$  dimensions, and an item to skill mapping  $Q$ -matrix. Each entry in  $\mathbf{D}$  encodes the performance of a student (see Table 1 for an example).

Table 1: Example data matrix to use with COMMAND. The performance of a student is encoded with 1 if the student answered correctly the item, and 0 otherwise.

User	Item 1	Item 2	Item 3	Item $p$
Alice	0	1		0
Bob	1	1	...	1
Carol	0	0		1
		...		

<sup>1</sup> $P(S_i = a | S_j = b)$  can be computed using any Bayesian network inference algorithm such as Junction tree algorithm[12].

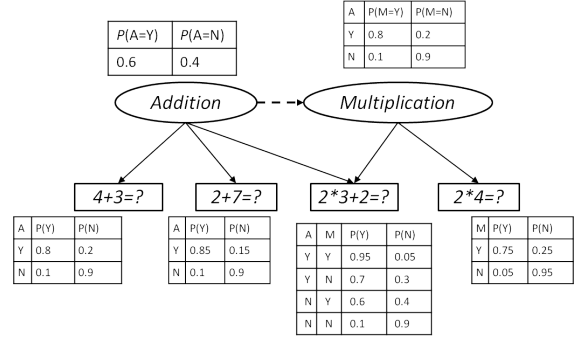


Figure 1: A hypothetical Bayesian network. Solid edges are given by item to skill mapping, dashed edges between skill variables are to be discovered from data. The conditional probability tables are to be learned.

#### Algorithm 1 The COMMAND algorithm

**Require:** A matrix  $\mathbf{D}$  of student performance on a set of test items, skill-to-item mapping  $Q$  (containing a set of skills  $\mathbf{S}$ ).

- 1:  $G_0 \leftarrow \text{Initialize}(\mathbf{S}, Q)$
- 2:  $i \leftarrow 0$
- 3: **do**
- 4:   *E*-step:
- 5:    $\theta_i^* \leftarrow \text{ParametricEM}(G_i, \mathbf{D})$
- 6:    $\mathbf{D}_i^* \leftarrow \text{Inference}(G_i, \theta_i^*, \mathbf{D})$
- 7:   *M*-step:
- 8:    $\langle G_{i+1}, \theta_{i+1} \rangle \leftarrow \text{BNLearning}(G_i, \mathbf{D}_i^*)$
- 9:    $i \leftarrow i + 1$
- 10: **while** stop criterion is not met
- 11:  $RE \leftarrow \text{FindReversibleEdges}(G_i)$
- 12:  $EC \leftarrow \text{EnumEquivalentDAGs}(G_i)$
- 13:  $DE \leftarrow \{\}$
- 14: **for** every reversible edge  $S_i - S_j$  in  $RE$  **do**
- 15:    $ratio \leftarrow \frac{P(S_i=1|S_j=1) \cdot P(S_j=0|S_i=0)}{P(S_j=1|S_i=1) \cdot P(S_i=0|S_j=0)}$
- 16:   **if**  $ratio \geq 1$  **then**
- 17:      $ratio^* = ratio$
- 18:      $DE \leftarrow DE \cup S_i \rightarrow S_j$
- 19:   **else**
- 20:      $ratio^* = \frac{1}{ratio}$
- 21:      $DE \leftarrow DE \cup S_i \leftarrow S_j$
- 22:   **end if**
- 23: **end for**
- 24:  $sort(DE)$  by  $ratio^*$  in descending order
- 25: **while**  $DE$  is not empty **do**
- 26:    $e \leftarrow dequeue(DE)$
- 27:   **if**  $\exists G \in EC$   $e \in G$  **then**
- 28:      $\forall G \in EC$ , remove  $G$  from  $EC$  if  $e \notin G$
- 29:   **end if**
- 30: **end while**
- 31: **return**  $EC$

COMMAND relies on a popular machine learning algorithm called Structural Expectation Maximization (EM), which has not been used in educational applications. A secondary contribution of our work is introducing Structural EM for learning Bayesian network structures from educational data. We now describe the steps of COMMAND in detail.

### 3.1 Initial Bayesian Network

COMMAND represents the prerequisite structure using Bayesian networks that use latent binary variables to represent the student knowledge of a skill, and observed binary variables that represent the student performance answering items (e.g., correct or incorrect). We first create an initial Bayesian network that complies to the skill-to-item  $Q$ -matrix. That is, we create an arc to each item from each of its required skills and leave all skill variables disconnected. With the created Bayesian network as an initial network, we learn the arcs between the skill variables using Structural EM.

### 3.2 Structural EM

A common solution to learning a Bayesian network from data is the score-and-search approach [6, 10]. This approach uses a scoring function<sup>2</sup> to measure the fitness of a Bayesian network structure to the observed data, and attempts to find the optimal model in the space of all possible Bayesian network structures. However, the conventional score-and-search approaches rely on efficient computation of the scoring function, which is only feasible for problems where data contains observations for all variables in the Bayesian network. Unfortunately, our domain has skill variables that are not directly observed. An intuitive work-around is to use the Expectation Maximization (EM) to estimate the scoring function. However, EM in this case takes a large number (hundreds) of iterations to converge and each iteration requires Bayesian network inference, which is computationally prohibitive. Further, we need run EM for each candidate structure. However, the number of possible Bayesian network structures is super-exponential with respect to the number of nodes. The Structural Expectation Maximization algorithm [8, 9] is an efficient alternative.

Structural EM is an iterative algorithm that inputs a matrix  $\mathbf{D}$  of student performance (see example Table 1). Figure 2 illustrates one iteration of the Structural EM algorithm. The relevant steps are also sketched in Algorithm 1. Each iteration consists of an Expectation step ( $E$ -step) and a Maximization step ( $M$ -step). In  $E$ -step, we first find the maximum likelihood estimate  $\theta^*$  of the parameters for the current structure  $G$  calculated from previous iteration using parametric EM.<sup>3</sup> We then do Bayesian inference to compute the expected values for the hidden variables using the current model  $(G, \theta^*)$  and use the values to complete the data. In the  $M$ -step, we use the conventional score-and-search approach to optimize the structure according to the completed data. Since the space of possible Bayesian network structures is super-exponential, exhaustive search is intractable and local search algorithms, such as greedy hill-climbing search, are often used. The  $E$ -step and  $M$ -step interleave and iterate until some stop criterion is met, e.g., the scoring function does not change significantly. Contrast to the conventional score-and-search algorithm, Structural EM runs EM only on one structure in each iteration, thus is computationally more efficient.

COMMAND's initialization step fixes the arcs from skills to items according to the  $Q$ -matrix. In the  $M$ -step of our Structural EM implementation we only consider the candidate structures that comply with the  $Q$ -matrix.

An advantage of using Structural EM to discover the prerequisite

<sup>2</sup>A commonly used scoring function is Bayesian information criterion (BIC), which is composed of a likelihood term and a term to penalize the model complexity

<sup>3</sup>In the first iteration, the current network is created from the initialization step.

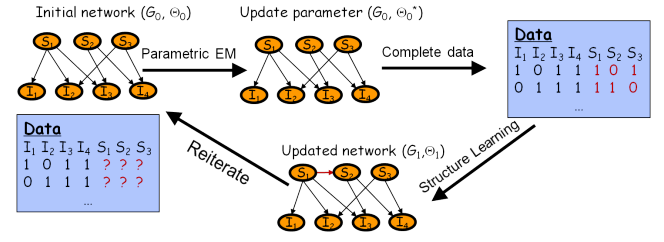


Figure 2: An illustration of the Structure EM algorithm to discover the structure of the latent variables.

relationship of skills, is that it is easily extensible to incorporate domain knowledge. For example, we can place constraints on the output structure to force or to disallow a skill to be a prerequisite of another other skill.

Another advantage of Structural EM is that it can be applied when there are missing data in the student performance matrix  $\mathbf{D}$ . That is, some students do not answer all the items thus some responses are missing. This "missing values" problem is very common in real-world situations. Structural EM can be applied to this type of data since it was originally developed to solve the "missing values" problem [8]. The general idea is, in the  $E$ -step, the algorithm also computes the expected values for missing data points, in addition for latent variables.

### 3.3 Discriminate Between Equivalent BNs

Structural EM selects Bayesian network model based on how well it explains the distribution of the data. Bayesian network theory states that some Bayesian networks are statistically equivalent in representing the data. Thus, the output from Structure EM is an equivalence class that may contain many Bayesian network structures. These equivalent Bayesian networks have the same skeleton and the same  $v$ -structures<sup>4</sup>. For instance, Figure 3 gives an example of two simple Bayesian networks that are not distinguishable by Structural EM algorithm. They share the skeleton but differ in the orientation of the edge (we will call such an edge a reversible edge). They apparently represent two different prerequisite structures.

#### 3.3.1 Use Domain Knowledge

To determine a unique structure, we shall determine the orientation of each reversible edge. For this purpose, we can use the following domain knowledge.

**Knowledge 1.** If  $S_1$  is a prerequisite of  $S_2$ , i.e.,  $S_1 \rightarrow S_2$ , then  $P(S_1 = 1 | S_2 = 1) = 1$  and  $P(S_2 = 0 | S_1 = 0) = 1$ .

**Knowledge 1** says if a skill  $S_1$  is the prerequisite of a skill  $S_2$ , a student must master skill  $S_1$  before he masters  $S_2$  and it is impossible for him to master  $S_2$  if he has not mastered  $S_1$ . In other words,  $S_1$  is not a prerequisite of  $S_2$  if at least one of the conditions does not hold. This puts a constraint on the joint distribution encoded by the Bayesian network to be learned. We can check for any reversible edge  $S_1 - S_2$  in the Bayesian network, either  $P(S_1 = 1 | S_2 = 1) = 1$  and  $P(S_2 = 0 | S_1 = 0) = 1$ , or  $P(S_2 = 1 | S_1 = 1) = 1$  and  $P(S_1 = 0 | S_2 = 0) = 1$ . If the former holds,  $S_1 \rightarrow S_2$ ; otherwise,  $S_1 \leftarrow S_2$ .

<sup>4</sup>A  $v$ -structure in a Bayesian network  $G$  is an ordered triple of nodes  $(u, v, w)$  such that  $G$  contains the directed edges  $u \rightarrow v$  and  $w \rightarrow v$  and  $u$  and  $w$  are not adjacent in  $G$ . [19].

However, in real situations, it is possible for a student to master the post-requisite even he does not know the prerequisite. Further, there is always noise in the data. Thus, the two conditional probabilities  $P(S_1 = 1|S_2 = 1)$  and  $P(S_2 = 0|S_1 = 0)$  are not exactly but close to 1. Thus, we use the following empirical rule: if  $P(S_1 = 1|S_2 = 1) \cdot P(S_2 = 0|S_1 = 0) \geq P(S_2 = 1|S_1 = 1) \cdot P(S_1 = 0|S_2 = 0)$ , we determine  $S_1 \rightarrow S_2$ ; otherwise, we determine  $S_1 \leftarrow S_2$ . Note that these conditional probabilities can be computed easily from the Bayesian network model output from Structural EM. The intuition behind this is that the two probabilities  $P(S_1 = 1|S_2 = 1)$  and  $P(S_2 = 0|S_1 = 0)$  are certificates of the prerequisite structure  $S_1 \rightarrow S_2$ . The larger of these two probabilities, the more likely the relationship  $S_1 \rightarrow S_2$  holds. Since here we are concerned with which direction the edge goes, we simply compare the products of two probabilities and select the direction that is more probable.



Figure 3: Two equivalent Bayesian networks representing different prerequisite structures.

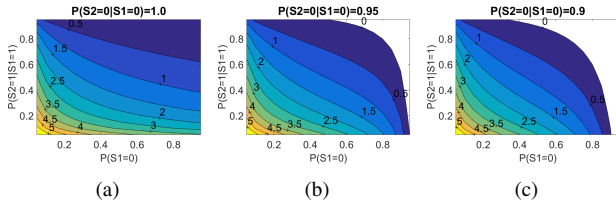


Figure 4: Contour plots of  $\log(\text{ratio})$  against  $P(S_1 = 0)$  and  $P(S_2 = 1|S_1 = 1)$  for various values of  $P(S_2 = 0|S_1 = 0)$ .

### 3.3.2 Theoretical Justification

The empirical rule proposed in § 3.3.1 needs theoretical justification. We consider the simple Bayesian network in Figure 3a, which has only one reversible edge  $S_1 - S_2$ . This model has three free conditional probability parameters:  $P(S_1 = 0) = p$ ,  $P(S_2 = 0|S_1 = 0) = q$ ,  $P(S_2 = 1|S_1 = 1) = r$ . If we observe data generated from this model, can we use the rule to determine  $S_1 \rightarrow S_2$ ? Equivalently, we use rule  $\text{ratio} = \frac{P(S_1=1|S_2=1) \cdot P(S_2=0|S_1=0)}{P(S_2=1|S_1=1) \cdot P(S_1=0|S_2=0)} \geq 1$ . Using Bayes rule and rules of probability, the rule becomes

$$\text{ratio} = \frac{(1-p)(qp + 1 - r - p + pr)}{p[r(1-p) + p(1-q)]} \geq 1, \quad (1)$$

which can be simplified to

$$\frac{(1-p)(1-r) - p(1-q)}{p[r(1-p) + p(1-q)]} \geq 0. \quad (2)$$

The rule holds only if  $(1-p)(1-r) - p(1-q) \geq 0$ . Thus, we can use the rule to determine  $S_1 \rightarrow S_2$  if  $(1-p)(1-r) - p(1-q) \geq 0$ . Since  $\text{ratio}$  depends on  $p$ ,  $q$  and  $r$ , we study how  $\text{ratio}$  changes with these parameters. Figure 4 shows the contour plots of  $\log(\text{ratio})$  against  $P(S_1 = 0)$  and  $P(S_2 = 1|S_1 = 1)$  for three different values of  $P(S_2 = 0|S_1 = 0)$ . The white region in each contour plot is the region where the rule is not applicable. That is, in this region,  $\text{ratio} < 1$  even  $S_1 \rightarrow S_2$  is the true model. As it is showed, if  $P(S_2 = 0|S_1 = 0) = q = 1$ , the rule is always applicable no matter what values  $P(S_1 = 0)$  and  $P(S_2 = 1|S_1 = 1)$  are. With  $P(S_2 = 0|S_1 = 0)$  decreasing, the white region becomes larger and the rule

becomes more inapplicable.  $P(S_2 = 0|S_1 = 0)$  can be interpreted as the strength of the prerequisite relationship. When the prerequisite relationship is weaker, the rule is less applicable and using the rule to determine the prerequisite relationship is less accurate.

### 3.3.3 Orient All Reversible Edges

Using the empirical rule, we can orient every reversible edge in the network structure. However, orienting each reversible edge is not independent and may conflict with each other. Having oriented one edge would constrain the orientation of other reversible edges because we have to ensure the graph is a DAG and the equivalence property is not violated. For example, in Figure 5a, if we have determined  $S_1 \rightarrow S_2$ , the edge  $S_2 \rightarrow S_3$  is enforced. In this paper, we take an ad-hoc strategy to determine the orientation for all reversible edges. For each reversible edge  $S_i - S_j$ , we let  $\text{ratio}^* = \text{ratio}$  if  $\text{ratio} > 1$  and  $\text{ratio}^* = \frac{1}{\text{ratio}}$  otherwise. The larger the  $\text{ratio}^*$  is, the more confidently when we decide the orientation. We sort the list of reversible edges by  $\text{ratio}^*$  in descending order. We then orient the edges by this ordering. When one edge is oriented, the constraint is propagated to other reversible edges. In practical implementation, we use the following strategy: we first enumerate all equivalent Bayesian networks and make them a list of candidates; when an edge is oriented to  $S_i \rightarrow S_j$ , we remove all contradicting Bayesian networks from the list. Eventually only one Bayesian network structure stands. This procedure is detailed in the *Discriminate between equivalent BNs* section of Algorithm 1. The *EnumEquivalentDAGs( $G_i$ )* implements the algorithm of enumerating equivalent DAGs in [4].

## 4. EVALUATION

In § 4.1, we evaluate COMMAND with simulated data to assess the quality of the discovered prerequisite structures. Then, in § 4.2 we use data collected from real students. In all our experiments, we use the Bayesian information criterion (BIC) as the scoring function in Structural EM.

### 4.1 Simulated Data

Synthetic data allow us to study how COMMAND compares to ground truth. For this, we engineered three prerequisite structures (DAGs), shown in Figure 5. Here, each figure represents different causal relations between the simulated latent skill variables.

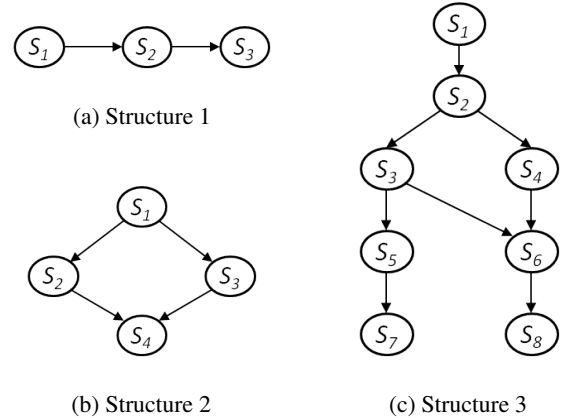


Figure 5: Three different DAGs between latent skill variables. Item nodes are omitted.

For clarity, Figure 5 omits the item nodes; but each skill node is parent of six item variables and each item variable has 1-3 skill nodes

as parents. All of these nodes are modeled using binary random variables. More precisely, the latent nodes represent whether the student achieves mastery of the skill, and the observed nodes indicate if the student answers the item correctly. Notice that these Bayesian networks include the prerequisite structures as well as the skill-item mapping.

We consider simulated data with different number of observations ( $n = 150, 500, 1000, 2000$ ). For each sample size and each DAG, we generate ten different sets of conditional probability tables randomly, with two constraints. First, we enforce that achieving mastery of the prerequisites of a skill will increase the likelihood of mastering the skill. Second, mastery of a skill increases the probability of student correctly answering the test item. Thus, in total we generated 120 synthetic datasets (3 DAGs x 4 sample sizes x 10 CPTs), and report the average results.

We evaluate how COMMAND can discover the true prerequisite structure using metrics designed to evaluate Bayesian networks structure discovery. In particular, we use the  $F_1$  adjacency score and the  $F_1$  orientation score. The adjacency score measure how well we can recover connections between nodes. It is a weighted average of the true positive adjacency rate and the true discovery adjacency rate. On the other hand, the orientation score measures how well we can recover the direction of the edges. It is calculated as a weighted average of the true positive orientation rate and true discovery orientation rate. In both cases, the  $F_1$  score reaches its best value at 1 and worst at 0. Moreover, for comparison, we compute the  $F_1$  adjacency score for Bayesian network structures whose skill nodes are fully connected with each other. These fully connected DAGs will serve as baselines for evaluating the adjacency discovery.<sup>5</sup> For completeness, we list these formulas in tables 2 and 3, respectively.

Table 2: Formulas for measuring adjacency rate (AR)

Metric	Formula
True positive ( $TPAR$ )	$\frac{\# \text{ of correct adjacencies in learned model}}{\# \text{ of adjacencies in true model}}$
True discovery ( $TDAR$ )	$\frac{\# \text{ of correct adjacencies in learned model}}{\# \text{ of adjacencies in learned model}}$
$F_1\text{-}AR$	$\frac{2 \cdot TPAR \cdot TDAR}{TPAR + TDAR}$

Table 3: Formulas for measuring orientation rate (OR)

Metric	Formula
True positive ( $TPOR$ )	$\frac{\# \text{ of correctly directed edges in learned model}}{\# \text{ of directed edges in true model}}$
True discovery ( $TDOR$ )	$\frac{\# \text{ of correctly directed edges in learned model}}{\# \text{ of directed edges in learned model}}$
$F_1\text{-}OR$	$\frac{2 \cdot TPOR \cdot TDOR}{TPOR + TDOR}$

We use these metrics to evaluate the effect of varying the number of observations of the training set (sample size) on the quality of learning the prerequisite structure. We designed experiments to specifically answer the following four questions:

1. How does the type of items affect COMMAND’s ability to recover the prerequisite structure? We consider the situation where in the model each item requires only one skill and the situation where each item requires multiple skills.

<sup>5</sup>We do not compute  $F_1$  orientation score for fully connected DAGs because all edges in a fully connected DAG are reversible.

2. How well does COMMAND perform when there is noise in the data? We focus on studying noise due to the presence of unaccounted latent variables.
3. How well does COMMAND perform when the student performance data have missing values?
4. How is COMMAND compared with other prerequisite discovery methods? In particular, we compare COMMAND to the Probabilistic Association Rules Mining (PARM) method [5].

We now investigate these questions.

#### 4.1.1 Single-skill vs Multi-skill Items

We consider two situations where different types of  $Q$ -matrix are used. In the first situation, each item node maps to only one skill node. In the second one, each item maps to 1-3 skills. Figure 6 compares the  $F_1$  of adjacency discovery and edge orientation result under two types of  $Q$ -matrix. We observe that the accuracy for both the adjacency and the edge orientation improves with the amount of data. With just 2000 observations, the algorithm can recover the true structures almost perfectly. Additionally, when the model contains multiple-skill items, the accuracy is slightly lower than that where all items in the model are single-skilled items, probably because there are more parameters to be estimated in the case of multi-skill items.

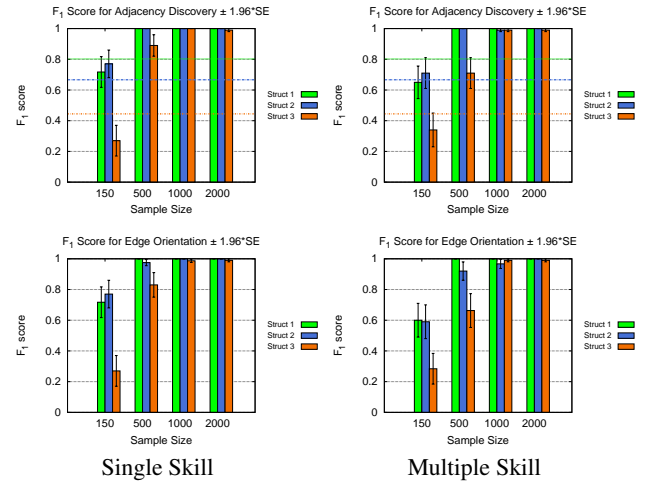


Figure 6: Comparison of  $F_1$  scores for adjacency discovery (top row) and for edge orientation (bottom row). Horizontal lines are baseline  $F_1$  scores computed for fully connected (complete) Bayesian networks. The error bars show the 95% confidence intervals, i.e.,  $\pm 1.96 \cdot SE$ .

#### 4.1.2 Sensitivity to Noise

Real-world data sets often contain various types of noise. For example, noise may occur due to latent variables that are not explicitly modeled. To evaluate the sensitivity of COMMAND to noise, we synthesize the three Bayesian networks in Figure 5 to include a *StudentAbility* node that takes three possible states (low/med/high). In these Bayesian networks, students’ performance depends not only on whether they have mastered the skill, but also on their individual ability. For simplicity, all items in the setting are single-skilled items. We first simulated data from Bayesian networks that have a *StudentAbility* variable to generate “noisy” data samples, and then



use this data to recover the prerequisite structure. Figure ?? illustrates the procedure of this sensitivity analysis experiment for Structure 1.

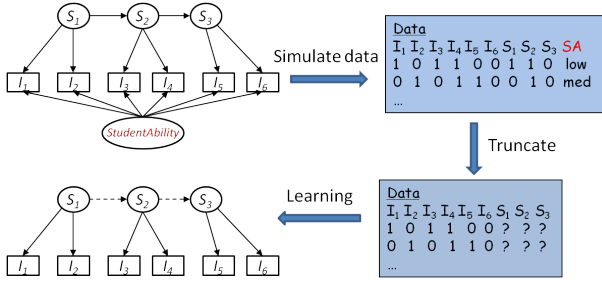


Figure 7: Evaluation of COMMAND with noisy data

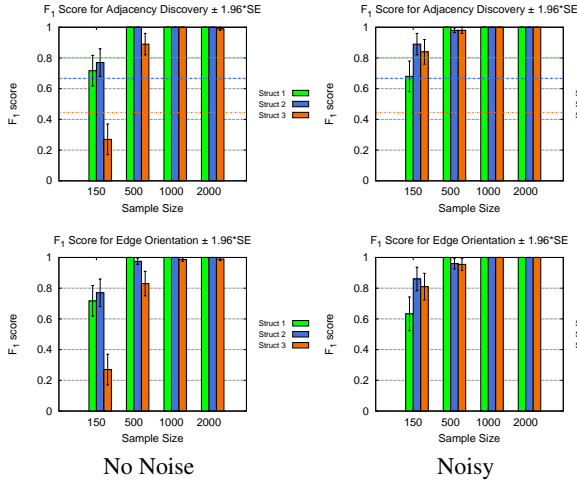


Figure 8: Results of adding systematic noise. Top: Comparison of  $F_1$  scores for adjacency discovery. Horizontal lines are baseline  $F_1$  scores computed for fully connected Bayesian networks. Bottom: Comparison of  $F_1$  scores for edge orientation.

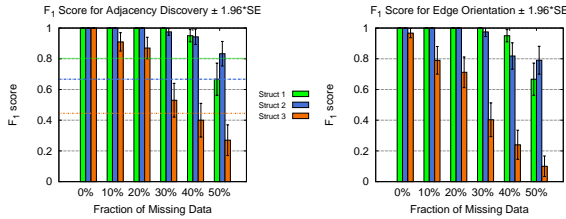


Figure 9: Results of learning with missing data. Left: Comparison of  $F_1$  scores for adjacency discovery. Horizontal lines are baseline  $F_1$  scores computed for fully connected Bayesian networks. Right: Comparison of  $F_1$  scores for edge orientation.

Figure 7 compares the results where noise was introduced or not. Interestingly, the noise does not harm the learning accuracy at all, and actually improves the accuracy. This improvement is more evident when the sample size is small (see  $n = 150$ ). We hypothesize that the correlations caused by *StudentAbility* node would be compensated by adding more edges between skill nodes. For small sample size, this might help because BIC score tends to give sparse structures.

#### 4.1.3 In Presence of Missing Values

Real-world student data often have missing values, i.e., students do not respond to all items. To evaluate how COMMAND performs on data with missing values, we generated data sets of size 1000 with varying fraction of randomly missing values (10%, 20%, 30%, 40%, 50%) and used COMMAND to recover the structures from these data sets. Again, the models only contain single-skilled items. The results are presented in Figure 8. It is seen that the accuracy decreases when the fraction of missing values increases. However, COMMAND is still able to recover the true structures for Structure 1 and 2 even when the data contain up to 30% missing values.

#### 4.1.4 Comparison With PARM

Chen and his colleagues proposed to use Probabilistic Association Rules Mining (PARM) for discovering the prerequisite relationships between skills [5]. Since PARM discovers pair-wise prerequisite relationship, instead of constructing the full structure, we only compare COMMAND with PARM for discovering the pair-wise relationships. That is, we derived pair-wise prerequisite relationships from the Bayesian network structure and see how the two approaches discover these relationships. Here we simulated data from Structure 3 (Figure 5(c)) (with single-skilled items), which has 21 pair-wise prerequisite relationships. When experimenting with PARM, the following cutoff values were used:  $minsup = 0.125$ ,  $minconf = 0.76$ ,  $minprob = 0.9$ . These values have been used in the simulation experiments in [5]. The evaluation metric used is  $F_1 = \frac{2 \cdot TPR \cdot TDR}{TPR + TDR}$ , where  $TPR = \frac{\text{\# of correct relationships learned}}{\text{\# of relationships in true model}}$  and  $TDR = \frac{\text{\# of correct relationships learned}}{\text{\# of relationships in learned model}}$ . The result is presented in Figure 9. It shows that COMMAND significantly outperforms PARM for sample size  $n \geq 500$ . The low  $F_1$  score by PARM is caused by low  $TPR$ , i.e., PARM failed to discover many prerequisite relationships (data not shown). It is speculated that selecting a different set of cutoff values for PARM might improve the results. However, determining these thresholds is not trivial and requires experts' intervention. By contrast, COMMAND does not require tuning of many parameters. Another disadvantage of PARM is that it does not respect transitivity. That is, PARM may not identify  $S_1 \rightarrow S_3$  even when it has discovered  $S_1 \rightarrow S_2$  and  $S_2 \rightarrow S_3$ .

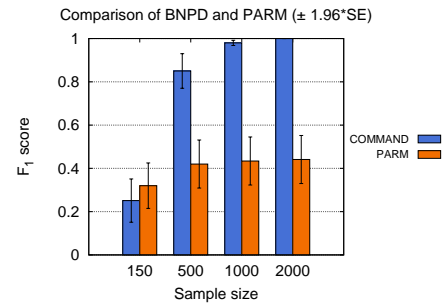


Figure 10: Comparison of COMMAND and PARM for discovering prerequisite relationships in Structure 3. The cutoff values used for PARM experiments are:  $minsup = 0.125$ ,  $minconf = 0.76$ ,  $minprob = 0.9$ .

## 4.2 Real Student Performance Data

We now evaluate COMMAND using two real-world student data sets.

#### 4.2.1 ECPE Data Set

The ECPE (Examination for the Certification of Proficiency in English) data set was collected from a test by the English language Institute of the University of Michigan to examine individual cognitive skills required for understanding English language grammar [18]. It contains a sample of 2922 examinees who are tested by 28 items on 3 skills, i.e., *morphosyntactic rules* ( $S_1$ ), *cohesive rules* ( $S_2$ ) and *lexical rules* ( $S_3$ ). Each item requires either one or two of the three skills. The prerequisite structure output from COMMAND is depicted in Figure 10. The estimated conditional probability tables (CPT) are shown correspondingly. The discovered structure says *lexical rules* is a prerequisite of *cohesive rules* and *morphosyntactic rules*, and *cohesive rules* is a prerequisite of *morphosyntactic rules*. This totally agrees with the findings in [18] and that by the PARM method in [5]. Further, COMMAND also outputs the conditional probabilities associated with each skill and its direct prerequisite. We clearly see that the probability of student mastering a skill increases when the student has acquired more prerequisites of the skill.

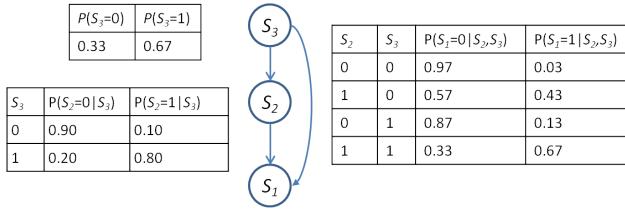


Figure 11: Result on ECPE data set.  $S_1$ : Morphosyntactic rules;  $S_2$ : Cohesive rules;  $S_3$ : Lexical rules. The estimated conditional probability tables (CPTs) are shown correspondingly.

#### 4.2.2 Math Data Set

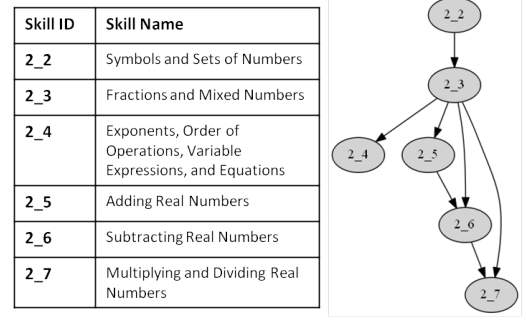
We now evaluate COMMAND using data collected from a commercial non-adaptive tutoring system. The textbook items are classified in chapters, sections, and objectives; but for this paper, we only use the data from Chapter 2 and Chapter 3 of the book. We use performance data while students solve test items. That is, students are tested on the items after they have been taught all relevant skills.

**Q-matrix and preprocessing.** We define skills as book sections. We use an item-to-skill mapping (*Q*-matrix) that assigns each exercise to a skill solely as the book section in which the item appears.<sup>6</sup> For each chapter, We process the data set to find a subset of items and students that does not have missing values. This is, the dataset we use in COMMAND has students responding to *all* of the items.

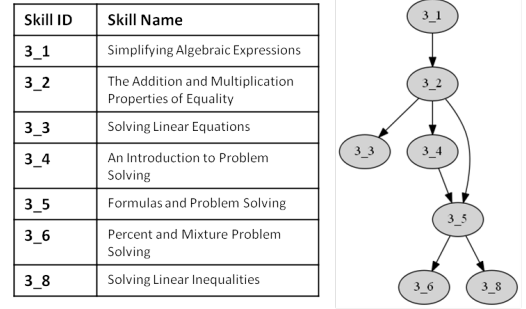
After filtering, two data sets, Math-chap2 and Math-chap3, were obtained for Chapter 2 and 3 respectively. In Math-chap2, six skills are included and each skill is tested on three to eight items, for a total of 30 items. In Math-chap3, seven skills are included and each skill has three to seven items, for a total of 33 items. Math-chap2 includes student test results for 1720 students, while the Math-chap3 has test results for 1245 students.

For simplicity we use binary variables to encode performance data (i.e., correct or incorrect) and skill variables (i.e., mastery or not mastery).

<sup>6</sup>Here we assume the items are single-skilled despite that they might be multi-skilled. In future a more accurate *Q*-matrix should be used.



(a) Prerequisite structure learned for Math-chap2.



(b) Prerequisite structure learned for Math-chap3.

Figure 12: Prerequisite structures constructed by COMMAND for HED data sets.

**Prerequisite Structure Discovery.** The Bayesian networks generated with the COMMAND algorithm are illustrated in Figure 11. Our observation is that the topological order of the sections in both structures are fully consistent with the book ordering heuristic. This shows an agreement between our fully data-driven method and human experts. We also ran PARM approach to learn pair-wise prerequisite relationships from these data sets. Given  $minsup = 0.125$ ,  $minconf = 0.76$  and  $minprob = 0.9$ ,  $2_5 \rightarrow 2_6$ ,  $2_5 \rightarrow 2_7$  and  $2_6 \rightarrow 2_7$  are discovered for Math-chap2,  $3_1 \rightarrow 3_3$  and  $3_2 \rightarrow 3_3$  are discovered for Math-chap3. These relationships are small subset of the set of relationships discovered by COMMAND.

**Predictive Performance of Prerequisite Models.** COMMAND not only discovers the topological relationships among skills, but also outputs a statistical model in the form of Bayesian network. This statistical model can be used for inference and predictive modeling. For example, given a student's response to a set of items, we can infer the student's knowledge status of any skill. This is very useful for remedial intervention, i.e., discovering when a student is lacking background knowledge. In this paper, we evaluate these models by evaluating how well the generated Bayesian networks predict student performance on a test item given response on other items. In particular, we compute the posterior probability of a student's response to an item  $I_i$  given his performance on all other items  $\mathbf{I}_{-i} = \mathbf{I} \setminus \{I_i\}$ , by marginalizing over the set of latent variables  $\mathbf{S}$ :

$$P(I_i|\mathbf{I}_{-i} = \mathbf{i}_{-i}) = \sum_{\mathbf{S}} P(I_i, \mathbf{S}|\mathbf{I}_{-i} = \mathbf{i}_{-i}). \quad (3)$$

This can be computed efficiently using the Junction tree algorithm [12]. We then do binary classification based on the posterior probability. We compare the Bayesian network models generated from COMMAND with five baseline predictors:

- A *majority* classifier which always classifies an instance to the majority class. For example, if majority of the students get an item wrong, other students would likely get it wrong.
- A Bayesian network model in which the skill variables are *disconnected*. This model assumes that the skill variables are marginally independent of each other. Most existing knowledge tracing approaches make this assumption.
- A Bayesian network model in which the skill variables are connected in a *chain* structure, i.e.,  $2-2 \rightarrow 2-3 \rightarrow 2-4 \rightarrow \dots$ . This assumes that a section (skill) only depends on the previous section. In other words, a first-order Markov chain dependency structure.
- A Bayesian network model constructed using the pairwise relationships output from *PARM*. That is, we create an edge  $S_i \rightarrow S_j$  if *PARM* says  $S_i$  is the prerequisite of  $S_j$ .
- A *fully connected* Bayesian network where skill variables are fully connected with each other. This model assumes no conditional independence between skill variables and can encode any joint distribution over the skill variables. However, it has exponential number of free parameters and thus can easily overfit the data.

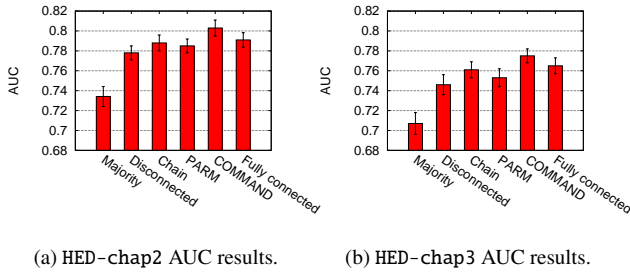


Figure 13: Comparison of AUC of six models to predict student performance on test items. The results are from 10 fold cross-validation.

The parameters of these baseline Bayesian network predictors are estimated from the data. The model predictions were evaluated using the *Area Under the Curve* (AUC) of the Receiver Operating Characteristic (ROC) curve metric calculated from 10-fold cross-validation. Results are presented in Figure 12. The error bars show the 95% confidence intervals calculated from the cross-validation. On both *Math-chap2* and *Math-chap3* data sets, the *COMMAND* models outperform the other five models. The *fully connected* models are the second best performing models. On *Math-chap2*, *COMMAND* model has an AUC of  $0.803 \pm 0.008$  and the *fully connected* model has an AUC of  $0.791 \pm 0.007$  (Figure 12a). A paired *t*-test reveals that the AUC difference of two models are statistically significant with a *p*-value of 0.0022. On *Math-chap3*, *COMMAND* model has an AUC of  $0.775 \pm 0.007$  and the *fully connected* model has an AUC of  $0.765 \pm 0.008$  (Figure 12b). The AUC difference of two models are also statistically significant with a *p*-value of 0.01. The *fully connected* models are outperformed by the much simpler prerequisite models, suggesting overfitting.

Using prerequisite structure to improve student models has been demonstrated in previous works [11, 2]. In their works, the prerequisite structure are hand-engineered by human experts. But here, we demonstrate with fully data-driven prerequisite structures.

## 5. CONCLUSION AND DISCUSSION

Discovering prerequisite structures of skills from data is challenging since skills are latent variables in the data. In this paper, we proposed *COMMAND*, a novel pipeline for learning the prerequisite dependencies of latent skill variables from students' response data on test items. We demonstrated our approach using both synthetic and real-world data under different conditions such as data contain missing values or noises. We show the discovered prerequisite models are useful for student modeling and remedial intervention by demonstrating its predictive performance on unseen data.

Studying prerequisite structure discovery in the context of Bayesian network structure learning is not new [16]. However, the approach in [16] does not discriminate between equivalent Bayesian network structures which represent different prerequisite structures. In our approach, we address this problem by using an ad-hoc strategy to refine the prerequisite model. In future, it would be wiser to refine the model within the Structural EM learning. For example, in Structural EM search, we can use a scoring function with an additional term that penalizes the structures that violate the CPT constraint. In such way, learning may be more accurate.

## 6. REFERENCES

- [1] Russell G Almond, Robert J Mislevy, Linda Steinberg, Duanli Yan, and David Williamson. 2015. *Bayesian networks in educational assessment*. Springer.
- [2] Anthony Botelho, Hao Wan, and Neil Heffernan. 2015. The prediction of student first response using prerequisite skills. In *Proceedings of the Second (2015) ACM Conference on Learning@ Scale*. ACM, 39–45.
- [3] Emma Brunskill. 2010. Estimating prerequisite structure from noisy data. In *Educational Data Mining 2011*.
- [4] Yetian Chen and Jin Tian. 2014. Finding the k-best Equivalence Classes of Bayesian Network Structures for Model Averaging.. In *AAAI*. 2431–2438.
- [5] Yang Chen, Pierre-Henri Wuillemin, and Jean-Marc Labat. 2015. Discovering Prerequisite Structure of Skills through Probabilistic Association Rules Mining. In *Proceedings of the 8th International Conference on Educational Data Mining*. 117–124.
- [6] Gregory F Cooper and Edward Herskovits. 1992. A Bayesian method for the induction of probabilistic networks from data. *Machine learning* 9, 4 (1992), 309–347.
- [7] Michel C Desmarais, Peyman Meshkinfam, and Michel Gagnon. 2006. Learned student models with item to item knowledge structures. *User Modeling and User-Adapted Interaction* 16, 5 (2006), 403–434.
- [8] Nir Friedman. 1997. Learning belief networks in the presence of missing values and hidden variables. In *ICML*, Vol. 97. 125–133.
- [9] Nir Friedman. 1998. The Bayesian structural EM algorithm. In *Proceedings of the Fourteenth conference on Uncertainty in artificial intelligence*. 129–138.
- [10] David Heckerman, Christopher Meek, and Gregory Cooper. 1997. *A Bayesian approach to causal discovery*. Technical Report. MSR-TR-97-05, Microsoft Research.
- [11] Tanja Käser, Severin Klingler, Alexander Gerhard Schwing, and Markus Gross. 2014. Beyond knowledge tracing: Modeling skill topologies with bayesian networks. In *Intelligent Tutoring Systems*. Springer, 188–198.
- [12] Daphne Koller and Nir Friedman. 2009. *Probabilistic graphical models: principles and techniques*. MIT press.
- [13] RJ Mislevy, RG Almond, D Yan, and LS Steinberg. 2000. *Bayes Nets in Educational Assessment: Where do the Numbers Come From, CRESST/Educational Testing Service*. Technical Report. CSE Technical Reports 518.
- [14] Judea Pearl. 2000. *Causality: models, reasoning and inference*. Vol. 29. Cambridge Univ Press.
- [15] Chris Piech, Jonathan Spencer, Jonathan Huang, Surya Ganguli, Mehran Sahami, Leonidas Guibas, and Jascha Sohl-Dickstein. 2015. Deep Knowledge Tracing. *arXiv preprint arXiv:1506.05908* (2015).
- [16] Richard Scheines, Elizabeth Silver, and Ilya Goldin. 2014. Discovering prerequisite relationships among knowledge components.



In *Educational Data Mining 2014*.

- [17] Peter Spirtes, Clark Glymour, and Richard Scheines. 2001. *Causation, prediction, and search*. MIT Press.
- [18] Jonathan Templin and Laine Bradshaw. 2014. Hierarchical diagnostic classification models: A family of models for estimating and testing attribute hierarchies. *Psychometrika* 79, 2 (2014), 317–339.
- [19] Thomas Verma and Judea Pearl. 1990. Equivalence and synthesis of causal models. In *Proceedings of the Sixth Annual Conference on Uncertainty in Artificial Intelligence*. 255–270.
- [20] Annalies Vuong, Tristan Nixon, and Brendon Towle. 2010. A method for finding prerequisites within a curriculum. In *Educational Data Mining 2011*.