# R-VGAE: Relational-variational Graph Autoencoder for Unsupervised Prerequisite Chain Learning

**Irene Li[1], Alexander Fabbri[1], Swapnil Hingmire[2] and Dragomir Radev[1]**
[1]Yale University, USA
[2]Tata Consultancy Services Limited (TCS), India

## Abstract

The task of concept prerequisite chain learning is to automatically determine the existence of prerequisite relationships among concept pairs. In this paper, we frame learning prerequisite relationships among concepts as an unsupervised task with no access to labeled concept pairs during training. We propose a model called the **R**elational-**v**ariational **G**raph **A**uto**E**ncoder (R-VGAE) to predict concept relations within a graph consisting of concept and resource nodes. Results show that our unsupervised approach outperforms graph-based semi-supervised methods and other baseline methods by up to 9.77% and 10.47% in terms of prerequisite relation prediction accuracy and F1 score. Our method is notably the first graph-based model that attempts to make use of deep learning representations for the task of unsupervised prerequisite learning. We also expand an existing corpus which totals $1,717$ English Natural Language Processing (NLP)-related lecture slide files and manual concept pair annotations over 322 topics.

## 1 Introduction

With the increasing amount of information available online, there is a rising need for structuring how one should process that information and learn knowledge efficiently in a reasonable order. As a result, recent work has tried to learn prerequisite relations among concepts, or which concept is needed to learn another concept within a *concept graph* (Liang et al., 2017; Gordon et al., 2016; AlSaad et al., 2018). Figure 1 shows an illustration of prerequisite chains as a directed graph. In such a graph, each node is a concept, and the direction of each edge indicates the prerequisite relation. Consider two concepts $p$ and $q$, we define $p \rightarrow q$ as $p$ is a prerequisite concept of $q$. For example, the concept *Variational Autoencoders* is a prerequisite concept of the concept *Variational Graph Autoencoders*. If someone wants to learn about the concept *Variational Graph Autoencoders*, the prerequisite concept *Variation Autoencoder* should appear in the prerequisite concept graph in order to create a proper study plan.

Recent work has attempted to extract such prerequisite relationships from various types of materials including Wikipedia articles, university course dependencies or MOOCs (Massive Open Online Courses) (Pan et al., 2017; Gordon et al., 2016; Liang et al., 2017). However, these materials either need additional steps for pre-processing and cleaning, or contain too many noisy free-texts, bringing more challenges to prerequisite relation learning or extracting. Recently, Li et al. (2019) presented a collection of university lecture slide files mainly in NLP lectures with related prerequisite concept annotations. We expanded this dataset as we believe these lecture slides offer a concise yet comprehensive description of advanced topics.

Deep models such as word embeddings (Mikolov et al., 2013) and more recently contextualized word embeddings (Devlin et al., 2018) have achieved great success in the NLP tasks as demonstrate a stronger ability to represent the semantics of the words than other traditional models. However, recent prerequisite learning approaches fail to make use of distributional semantics and advances in deep learning representations (Labutov et al., 2017; Pan et al., 2017). In this paper, we investigate deep node embeddings within a graph structure to better capture the semantics of concepts and resources, in order to learn accurate the prerequisite relations.
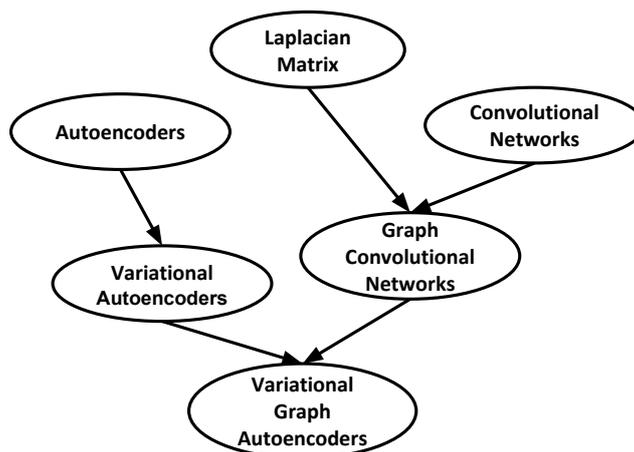
Figure 1: An illustration of prerequisite chains: we show six concepts and the relations. For example, the concept *Variational Autoencoders* is a prerequisite concept of the concept *Variational Graph Autoencoders*.

In addition to learning node representations, there has been growing research in geometric deep learning (Bronstein et al., 2017) and graph neural networks (Gori et al., 2005), which apply the representational power of neural networks to graph-structured data. Notably, Kipf and Welling (2017) proposed Graph Convolutional Networks (GCNs) to perform deep learning on graphs, yielding competitive results in semi-supervised learning settings. TextGCN was proposed by (Yao et al., 2018) to model a corpus as a heterogeneous graph in order to jointly learn word and document embeddings for text classification. We build upon these ideas for constructing a resource-concept graph[1]. Additionally, most of the mentioned methods require a subset of labels for training, a setting which is often infeasible in the real world. Limited research has been investigated learning prerequisite relations without using human annotated relations during training (AlSaad et al., 2018). In practice, it is very challenging to obtain annotated concept-concept relations, as the complexity for annotating is $O(n^2)$ given $n$ concepts. To tackle this issue, we propose a method to learn prerequisite chains without any annotated concept-concept relations, which is more applicable in the real word.

Our contributions are two-fold: 1) we expand upon previous annotations to increase coverage for prerequisite chain learning in five categories, including AI (artificial intelligence), ML (machine learning), NLP, DL (deep learning) and IR (information retrieval). We also expand a previous corpus of lecture files to include an additional 5000 more lecture slides, totaling 1,717 files. More importantly, we add additional concepts, totaling 322 concepts, as well as the corresponding annotations of each concept pair, which totals 103,362 relations. 2) we present a novel graph neural model for learning prerequisite relations in an unsupervised way using deep representations as input. We model all concepts and resources in the corpus as nodes in a single heterogeneous graph and define a propagation rule to consider multiple edge types by eliminating concept-concept relations during training, making it possible to perform unsupervised learning. Our model leads to improved performance over a number of baseline models. Notably, it is the first graph-based model that attempts to make use of deep learning representations for the task of unsupervised prerequisite learning. Resources, annotations and code are publicly available online[2].

---

[1]We use the term *resource* instead of *document* for generalization.
[2]https://github.com/Yale-LILY/LectureBank/tree/master/LectureBank2

| Domain | #courses | #files | #tokens | #pages | #tokens/page |
|---|---|---|---|---|---|
| NLP | 45 | 953 | 1,521,505 | 37,213 | 40.886 |
| ML | 15 | 312 | 722,438 | 12,556 | 57.537 |
| DL | 7 | 259 | 450,879 | 7,420 | 60.765 |
| AI | 5 | 98 | 139,778 | 3,732 | 37.454 |
| IR | 5 | 95 | 205,359 | 4,107 | 50.002 |
| **Overall** | 77 | 1,717 | 3,039,959 | 65,028 | 46.748 |

Table 1: Dataset Statistics. In each category, we have a given number of courses (#courses); each course consists of lecture files (#files); each lecture file has a number of individual slides (#pages). We also show the number of total tokens (#tokens) and average token number per slide (#tokens/page).

## 2 Related Work

### 2.1 Deep Models for Graph-structured Data

There has been much research focused on graph-structured data such as social networks and citation networks (Sen et al., 2008; Akoglu et al., 2015; Defferrard et al., 2016), and many deep models have achieved satisfying results. Deepwalk (Perozzi et al., 2014) was a breakthrough model which learns node representations using random walks. Node2vec (Grover and Leskovec, 2016) was an improved scalable framework, achieving promising results on multi-label classification and link prediction. Besides, there has been some work like graph convolution neural networks (GCNs), which target on deep-based propagation rules within graphs. A recent work applied GCN for text classification (Yao et al., 2018) by constructing a single text graph for a corpus based on word co-occurrence and document word relations. The experimental results showed that the proposed model achieved state-of-the-art methods on many benchmark datasets. We are inspired by this work in that we also attempt to construct a single graph for a corpus, however, we have different types of nodes and edges.

### 2.2 Prerequisite Chain Learning

Learning prerequisite relations between concepts has attracted much recent work in machine learning and NLP field. Existing research focuses on machine learning methods (i.e., classifiers) to measure the prerequisite relations among concepts (Liang et al., 2018; Liu et al., 2016; Liang et al., 2017). Some research integrates feature engineering to represent a concept, inputting these features to a classic classifier to predict relationship of a given concept pair (Liang et al., 2017; Liang et al., 2018). The resources to learn those concept features include university course descriptions and materials as well as online educational data (Liu et al., 2016; Liang et al., 2017). Recently, Li et al. (2019) introduced a dataset containing 1,352 English lecture files collected from university-level lectures as well as 208 manually-labeled prerequisite relation topics, initially introduced in (Fabbri et al., 2018). To avoid feature engineering, they applied graph-based methods including GAE and VGAE (Kipf and Welling, 2017) which treat each concept as a node thus building a concept graph. They pretrained a Doc2vec model (Le and Mikolov, 2014) to infer each concept as a dense vector, and then trained the concept graph in a semi-supervised way. Finally, the model was able to recover unseen edges of a concept graph. Different from their work, we wish to do the prerequisite chain learning in an unsupervised manner, while in training, no concept relations will be provided to the model.

## 3 Dataset

### 3.1 Resources

We manually collected English lecture slides mainly on NLP-related courses in recent years from known universities. We treated them as individual slide file in PDF or PowerPoint Presentations format. Our new collection has 529 additional files from 17 courses, which we combined with the data provided by (Li et al., 2019). We ended up with a total number of 77 courses with 1,717 English lecture slide files,
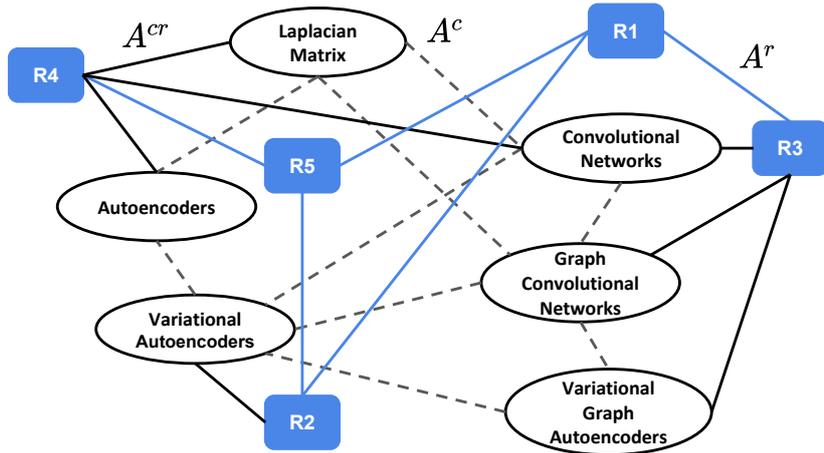
Figure 2: Concept-resource graph for prerequisite chain learning: oval nodes indicate concept nodes, the blue rectangular nodes indicate resource nodes. We show three types of edges: the blue edge between two resource nodes $A^r$, the black solid edge between a concept node and a resource node $A^{cr}$ and the black dashed edge between two concept nodes $A^c$. In the graph, the resource nodes R1 to R5 are example resources used to illustrate the idea. In practice there may be more edges, we show a part of the them for simplicity.

covering five domains. We show the final statistics in Table 1. For our experiments, we converted those files into TXT format which allowed us to load the free texts directly.

### 3.2 Concepts

We manually expanded the size of concept list proposed by (Li et al., 2019) from 208 to 322. We included concepts which were not found in their version like *restricted boltzmann machine* and *neural parsing*. Also, we re-visited their topic list and corrected a small number of the topics. For example, we combined certain topics (e.g. *BLUE* and *ROUGE*) into a single topic (*machine translation evaluation*). We asked two NLP PhD students to re-evaluate existing annotations from the old corpus and to provide labels for each added concept pair in the new corpus. A Cohen kappa score (Cohen, 1960) of 0.6283 achieved between our annotators which can be considered as a substantial agreement. We then took the union of the annotations, where if at least one judge stated that a given concept pair $(A, B)$ had $A$ as a prerequisite of $B$, then we define it a positive relation. We believe that the union of annotations makes more sense for our downstream application, where we want users to be able to mark which concepts they already know and displaying all potential concepts is essential. We have 1,551 positive relations on the 322 concept nodes.

## 4 Method

### 4.1 Problem Definition

In our corpus, every concept $c$ is a single word or a phrase; every resource $r$ is free text extracted from the lecture files. We then wish to determine for a given concept pair $(c_i, c_j)$, whether $c_i$ is a prerequisite concept of $c_j$. We define the concept-resource graph as $G = (X, A)$, where $X$ denotes node features or representations and $A$ denotes the adjacency matrix. In our case, the adjacency matrix is the set of relations between each node pair, or the edges between the nodes. In Figure 2, we build a single, large graph consisting of concepts (oval nodes) and resources (rectangular nodes) as nodes, and the corresponding relations as edges. So there are three types of edges in $A$: the edge between two resource nodes $A^r$ (blue line), the edge between a concept node and a resource node $A^{cr}$ (black solid line), and the edge between two concept nodes $A^c$ (black dashed line). Our goal is to learn the relations between concepts only ($A^c$), so prerequisite chain learning can be formulated as a link prediction problem. Our

unsupervised setting is to exclude any direct concept relations ($A^c$) during training, and we wish to predict these edges through message passing via the resource nodes indirectly.

## 4.2 Preliminaries

**Graph Convolutional Networks** (GCN) (Kipf and Welling, 2017) is a semi-supervised learning approach for node classification on graphs. It aims to learn the node representation $H = \{h_1, h_2, ..h_n\}$ in the hidden layers, given the initial node representation $X$ and the adjacency matrix $A$. The model incorporates local graph neighborhoods to represent a current node. In a simple GCN model, a layer-wise propagation rule can be defined as the following:

$$H^{(l+1)} = \sigma(AH^{(l)}W^{(l)}) \tag{1}$$

where $l$ is the current layer number, $\sigma(\cdot)$ is a non-linear activation function, and $W$ is a parameter matrix that can be learned during training. We eliminate the $\sigma(\cdot)$ for the last layer output. For the task of node classification, the loss function is cross-entropy loss. Typically, a two-layer GCN (by plugging Equation 1 in) is defined as:

$$GCN(X, A) = H^2 = \tilde{A}H^1W^1 = \tilde{A}\sigma(AXW^0)W^1 \tag{2}$$

where $\tilde{A}$ is the new adjacency matrix at the second graph layer.

**Relational Graph Convolutional Networks** (R-GCNs) (Schlichtkrull et al., 2018) expands the types of graph nodes and edges based on the GCN model, allowing operations on large-scale relational data. In this model, an edge between a node pair $i$ and $j$ is denoted as $(v_i, rel, v_j)$, where $rel \in Rel$ is considered a relation type, while in GCN, there is only one type. Similarly, to obtain the hidden representation of the node $i$, we consider the local neighbors and itself; when multiple types of edges exist, different sets of weight will be considered. So the layer-wise propagation rule is defined as:

$$h_i^{(l+1)} = \sigma\left(\frac{1}{M}\sum_{rel\in Rel}\sum_{j\in N_i^{rel}}(W_r^{(l)}h_j^{(l)} + W_0^{(l)}h_i^{(l)})\right) \tag{3}$$

where $Rel$ is the set of relations or edge types in the graph, $N_i^{rel}$ denotes the neighbors of node $i$ with relation $rel$, $W_r^{(l)}$ is the weight matrix at layer $l$ for nodes in $N_i^{rel}$, $W_0^{(l)}$ is the shared weight matrix at layer $l$, $M$ is the number of weight matrices in each layer.

**Variational Graph Auto-Encoders** (V-GAE) (Kipf and Welling, 2016) is a framework for unsupervised learning on graph-structured data based on variational auto-encoders (Kingma and Welling, 2013). It takes the adjacency matrix and node features as input and tries to recover the graph adjacency matrix $A$ through the hidden layer embeddings $Z$. Specifically, the non-probabilistic graph auto-encoder (GAE) model calculates embeddings via a two-layer GCN encoder: $Z = GCN(X, A)$, which is given by Equation 2.

Then, in the variational graph auto-encoder, the goal is to sample the latent parameters $z_i \in Z$ from a normal distribution:

$$q(z_i|X, A) = \mathcal{N}(z_i|\boldsymbol{\mu}_i, \text{diag}(\sigma_i^2)) \tag{4}$$

where $\mu = \text{GCN}_\mu(X, A)$ is the matrix of mean vectors, and $\log_\sigma = \text{GCN}_\sigma(X, A)$. The training loss then is given as the KL-divergence between the normal distribution and the sampled parameters $Z$:

$$\mathcal{L}_{\text{latent}} = \sum_{i\in\mathcal{N}}\text{KL}\left(\mathcal{N}\left(\boldsymbol{\mu}_i, \text{diag}(\boldsymbol{\sigma}_i)^2\right)\|\mathcal{N}(\mathbf{0}, \mathbf{I})\right) \tag{5}$$

In the inference stage, the reconstructed adjacency matrix $\hat{A}$ is the inner product of the latent parameters $Z$: $\hat{A} = \sigma(ZZ^T)$.

### 4.3 Proposed Model

To take multiple relations into consideration and make it possible to do unsupervised learning for concept relations, we propose our R-VGAE model. Our model builds upon R-GCN and VGAE by taking the advantages of both: R-GCN is a supervised model that deals with multiple relations; VGAE is an unsupervised graph neural network. We then make it possible to directly to train on a heterogeneous graph in an unsupervised way for link prediction, in order to learn the prerequisite relations for the concept pairs.

Our model first applies the R-GCN in Equation 3 as the encoder to obtain the latent parameters $Z$, given the initial node features $X$ and adjacency matrix $A$: $Z = \text{R-GCN}(X, A)$. In terms of the variational verison, as opposed to the standard VGAEs, we parameterize $\mu$ by the RGCN model: $\mu = \text{R-GCN}_\mu(X, A)$, and $\log_\sigma = \text{R-GCN}_\sigma(X, A)$.

To predict the link between a concept pair $(c_i, c_j)$, we followed the DistMult (Yang et al., 2014) method: we take the last layer output node features $\hat{X}$, and define the following score function to recover the adjacency matrix $\hat{A}$ by learning a trainable weight matrix $R$:

$$\hat{A} = \hat{X}^\intercal R \hat{X} \tag{6}$$

The loss consists of the cross-entropy reconstruction loss of adjacency matrix ($\mathcal{L}_{cross}$) and the loss from the latent parameters defined in Equation 5:

$$\mathcal{L} = \mathcal{L}_{cross}(A, \hat{A}) + \mathcal{L}_{latent} \tag{7}$$

We compare two variations of our R-GAE model. **Unsupervised**: only the concept-resource edges $A^{cr}$ and resource-resource edges $A^r$ are provided during training. This is an *unsupervised* model because no concept-concept edges are used. **Semi-supervised**: the model has access to concept-resource edges $A^{cr}$ and resource-resource edges $A^r$, as well as a percentage of the available concept-concept edges $A^c$, described later.

### 4.4 Node Features $X$

**Sparse Embeddings** We used TFIDF (term frequencyinverse document frequency) to get sparse embeddings for all nodes. We restricted the global vocabulary to be the 322 concept terms only, which means that the dimension of the node features is 322, as we aim to model keywords.

**Dense Embeddings** As the concepts in our corpus often consist of phrases such as *dynamic programming*, we made use of Phrase2vec (Artetxe et al., 2018). Phrase2vec (P2V) is a generalization of skip-gram models (Mikolov et al., 2013) which learns n-gram embeddings during training, and here we aim to infer the embeddings of the concepts in our corpus. We trained the P2V model using only our corpus by treating each slide file as a short document as a sequence of tokens. For each resource node, we take an element-wise average of the P2V embeddings of each single token and phrases that resource covered. Similarly, for each concept node, we took element-wise average of the embeddings of each individual token and the concept phrase. In addition, we then utilized the BERT model (Devlin et al., 2018) as another type of dense embedding. We fine-tuned the masked language modeling of BERT using our corpus.

### 4.5 Adjacency Matrix $A$

To construct the adjacency matrix $A$, for each node pair $(v_i, v_j)$, we applied cosine similarity based on enriched TFIDF features[3] as the value $A_{ij}$. Previous work has applied cosine similarity for vector space models (García-Pablos et al., 2018; Zuin et al., 2018; Bhatia et al., 2016), so we believe it is a suitable method in our case. This way we were able to generate concept-resource edge values ($A^{cr}$) and resource-resource edge values ($A^r$). Note that for concept-concept edge values $A^c$: 1 if $c_i$ is a prerequisite of $c_j$, 0 otherwise. These values are not computed in the unsupervised setting.

---

[3]This means that the TFIDF features are calculated on an extended vocabulary that includes all possible tokens appeared in the corpus.

| Method | Acc | F1 | MAP | AUC |
|---|---|---|---|---|
| *Concept embedding + classifier* | | | | |
| P2V (lb1) | 0.5927 | 0.5650 | 0.5623 | 0.5929 |
| P2V (lb2) | <u>0.6369</u> | <u>0.5961</u> | <u>0.6282</u> | <u>0.6370</u> |
| BERT (lb1) | 0.6540 | 0.6099 | 0.6475 | 0.6540 |
| BERT (lb2) | <u>0.6558</u> | <u>0.6032</u> | <u>0.6553</u> | <u>0.6558</u> |
| BERT (original) | **0.7088** | **0.6963** | **0.6779** | **0.7090** |
| *Graph-based methods* | | | | |
| DeepWalk (Perozzi et al., 2014) | 0.6292 | 0.5860 | 0.6270 | 0.6281 |
| Node2vec (Grover and Leskovec, 2016) | 0.6209 | 0.6181 | 0.5757 | 0.6259 |
| VGAE (Li et al., 2019) | 0.6055 | 0.6030 | 0.5628 | 0.6055 |
| GAE (Li et al., 2019) | **0.6307** | **0.6275** | **0.5797** | **0.6307** |
| R-GCN (Schlichtkrull et al., 2018) | 0.5387 | 0.4784 | 0.5203 | 0.5387 |
| *R-VGAE (Our proposed model)* | | | | |
| US+BERT (fine-tuned) | 0.5704 | 0.5704 | 0.5579 | 0.5955 |
| US+BERT (original) | 0.5669 | 0.5668 | 0.5658 | 0.6164 |
| US+TFIDF | 0.6495 | 0.6458 | 0.7069 | 0.5507 |
| US+P2V | 0.7694* | 0.7638* | 0.8919* | 0.9126* |
| SS+BERT (fine-tuned) | 0.6942 | 0.6942 | 0.6613 | 0.7412 |
| SS+BERT (original) | 0.6839 | 0.6839 | 0.6556 | 0.7372 |
| SS+TFIDF | 0.7252 | 0.7082 | 0.8181 | 0.7625 |
| SS+P2V | **0.8065** | **0.8010** | **0.9380** | **0.9454** |

Table 2: Accuracy (*Acc*), macro F1, MAP and AUC scores on balanced test set including 10% of prerequisite edges. Bold values are the best results within its experiment group. Underscored values indicate a better performance compared with the trained corpus. Values with an asterisk mean the best performance in the unsupervised setting.

## 5 Evaluation

We compare our proposed model with two groups of baseline models. We report accuracy, F1 scores, the macro averaged Mean Average Precision (MAP) and Area under the ROC Curve (AUC) scores in Table 2, as done by previous research (Chaplot et al., 2016; Pan et al., 2017; Li et al., 2019). We split the $1,551$ positive relations into 9:1 (train/test), and randomly select negative relations as negative training samples, and then we run over five random seeds and report the average scores, following the same setting with Kipf and Welling (2016) and Li et al. (2019).

**Concept embedding + classifier** The first group is the concept embedding with traditional classifiers including Support Vector Machines, Logistic Regression, Naïve Bayes and Random Forest. For a given concept pair, we concatenate the dense embeddings for both concepts as input to train the classifiers, and then we report the best result. We compare Phrase2Vec (P2V) and BERT embeddings. We have two corpora: one is the old version (*lb1*) provided by (Li et al., 2019), another one is our version (*lb2*). For the BERT model, we applied both the original version from Google (*original*) [4], and the fine-tuned language models version on our corpora (*lb1, lb2*) from Xiao (2018), and perform inference on the concepts. The P2V embeddings have 150 dimension, and the BERT embeddings have 768 dimensions. We show improvements on the BERT and P2V baselines by using our additional data via the underscored values. This indicates that the concept relations can be more accurately predicted when enriching the training corpus to train better embeddings. In our following experiments, if not specified, we applied *lb2* as the training corpus.

**Graph-based methods** We apply the classic graph-based embedding methods DeepWalk (Perozzi et al., 2014) and Node2vec (Grover and Leskovec, 2016), by considering the concept nodes only. Then the

---

[4] https://github.com/google-research/bert

positive concept relations in training set are the known sequences, allowing to train both models to infer node features. Similarly, in the testing phrase, we concatenate the node embeddings given a concept pair, and utilize the mentioned classifiers to predict the relation and report the performance of the best one. We then include VGAE and GAE methods for prerequisite chain learning following Li et al. (2019). Both methods construct the concept graph in a semi-supervised way. We apply P2V embeddings to replicate their methods, though it is possible to try additional embeddings, this is not our main focus. Finally, we compare with the original R-GCN model for link prediction proposed by Schlichtkrull et al. (2018) and apply the same embeddings with the VGAE and GAE methods. Other semi-supervised graph methods such as GCNs require node labels and thus are not applicable to our setting. We can see that the GAE method achieves the best results among the baselines. Compare with the first group, BERT (*original*) still has a better performance due to its ability to represent phrases.

**R-VGAE** Our model can be trained in both unsupervised (*US+\**) and semi-supervised (*SS+\**) way. We also utilize various types of embeddings include P2V, TFIDF, BERT (fine-tuned) and BERT (original). The best performed model in the unsupervised setting is with P2V embeddings, marked with asterisks, and it is better than all the baseline supervised methods with a large margin. In addition, our semi-supervised setting models boost the overall performance. We show that the *SS+P2V* model performs the best among all the mentioned methods, with a significant improvement of 9.77% in accuracy and 10.47% in F1 score compared with the best baseline model *BERT (original)*. This indicates that R-VGAE model does better on link prediction by bringing extra resource nodes into the graph, while the concept relation can be improved and enhanced indirectly via the connected resource nodes. We also observe that with BERT embeddings, the performance lags behind the other embedding methods for our approach. A reason might be that the dimensionality of the BERT embeddings is relatively large compared to P2V and may cause overfitting, especially when the edges are sparse; and it might not suitable to represent resources as they are a list of keywords when fine-tuning the language modeling. The P2V embeddings outperform TFIDF for both unsupervised and semi-supervised models. This shows that compared with sparse embeddings, dense embeddings can better preserve the semantic features when integrated within the R-GAE model, thus boosting the performance. Besides, as a variation of R-GCN and GAE, our model surpasses them by taking the advantages of both, comparing with R-GCN and GAE results reported in the second group.

## 6 Analysis

We then take the recovered concept relations from our best performed model *R-VGAE (SS+P2V)* in Table 2), and compare them with the gold annotated relations. Note that here we only look at concept nodes. The average degree for gold graph concept nodes is 9.79, while our recovered one has an average degree of 6.10, and this means our model predicts fewer edges. We also check the most popular concepts that have the most degrees. We select *dependency parsing* and *tree adjoining grammar* as examples. In Table 3, we show a comparison of the prerequisites from the annotations and our model's output. The upper group illustrates results for *dependency parsing*, where one can notice that the predicted concepts all appear in the gold results, missing only a single concept. This shows that even though our model predicts less number of relations, it still predicts correct relations. The lower group shows the comparison for the concept *tree adjoining grammar*, our model gives precise prerequisite concepts among all eight concepts from the gold set. When a concept has a certain amount number of prerequisite concepts, our model is able to provide a comprehensive concept set with a good quality. In the real word, especially in a learner's scenario, he or she wants to learn the new concept with enough prerequisite knowledge, which our model tends to provide.

## 7 Conclusion and Future Work

In this paper we introduced an expanded dataset for prerequisite chain learning with additional an 5,000 lecture slides, totaling 1,717 files. We also provided prerequisite relation annotations for each concept pair among 322 concepts. Additionally, we proposed an unsupervised learning method which makes use of advances in graph-based deep learning algorithms. Our method avoids any feature engineering

| Concept | Gold Prerequisite Concepts | Model Output Concept |
|---|---|---|
| *dependency parsing* | syntax<br>classic parsing methods<br>linguistics basics<br>parsing<br>nlp introduction<br>chomsky hierarchy<br>linear algebra<br>conditional probability | syntax<br>classic parsing methods<br>linguistics basics<br>parsing<br>nlp introduction<br>chomsky hierarchy<br>linear algebra |
| *tree adjoining grammar* | classic parsing methods<br>linguistics basics<br>parsing<br>nlp introduction<br>context free grammar<br>probabilistic context free grammars<br>chomsky hierarchy<br>context sensitive grammar | classic parsing methods<br>linguistics basics<br>parsing<br>nlp introduction<br>context free grammar<br>probabilistic context free grammars<br>chomsky hierarchy<br>context sensitive grammar |

Table 3: A comparison of prerequisite concepts of *dependency parsing* (upper group) and *tree adjoining grammar* (lower group) from our annotated gold labels and labels recovered by our best unsupervised model.

to learn concept representations. Experimental results demonstrate that our model performs well in an unsupervised setting and is able to further benefit when labeled data is available. In future work, we would like to perform a more comprehensive model comparison and evaluation by bringing other possible variations of graph-based models to learn a concept graph. Another interesting direction is to apply multi-task learning to the proposed model by adding a node classification task if there are node labels available. A part of the future work would also include developing educational applications for learners to find out their study path for certain concepts.

# References

Leman Akoglu, Hanghang Tong, and Danai Koutra. 2015. Graph based anomaly detection and description: a survey. *Data mining and knowledge discovery*, 29(3):626–688.

Fareedah AlSaad, Assma Boughoula, Chase Geigle, Hari Sundaram, and ChengXiang Zhai. 2018. Mining MOOC Lecture Transcripts to Construct Concept Dependency Graphs. *International Educational Data Mining Society*.

Mikel Artetxe, Gorka Labaka, and Eneko Agirre. 2018. Unsupervised Statistical Machine Translation. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing, Brussels, Belgium, October 31 - November 4, 2018*, pages 3632–3642.

Shraey Bhatia, Jey Han Lau, and Timothy Baldwin. 2016. Automatic labelling of topics with neural embeddings. *arXiv preprint arXiv:1612.05340*.

Michael M. Bronstein, Joan Bruna, Yann LeCun, Arthur Szlam, and Pierre Vandergheynst. 2017. Geometric Deep Learning: Going Beyond Euclidean Data. *IEEE Signal Process. Mag.*, 34(4):18–42.

Devendra Singh Chaplot, Yiming Yang, Jaime Carbonell, and Kenneth R Koedinger. 2016. Data-driven Automated Induction of Prerequisite Structure Graphs. *International Educational Data Mining Society*.

Jacob Cohen. 1960. A Coefficient of Agreement for Nominal Scales. *Educational and psychological measurement*, 20(1):37–46.

Michaël Defferrard, Xavier Bresson, and Pierre Vandergheynst. 2016. Convolutional neural networks on graphs with fast localized spectral filtering. In *Advances in neural information processing systems*, pages 3844–3852.

Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2018. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*.

Alexander Fabbri, Irene Li, Prawat Trairatvorakul, Yijiao He, Weitai Ting, Robert Tung, Caitlin Westerfield, and Dragomir Radev. 2018. Tutorialbank: A Manually-Collected Corpus for Prerequisite Chains, Survey Extraction and Resource Recommendation. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 611–620.

Aitor García-Pablos, Montse Cuadros, and German Rigau. 2018. W2vlda: almost unsupervised system for aspect based sentiment analysis. *Expert Systems with Applications*, 91:127–137.

Jonathan Gordon, Linhong Zhu, Aram Galstyan, Prem Natarajan, and Gully Burns. 2016. Modeling Concept Dependencies in a Scientific Corpus. In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, volume 1, pages 866–875.

Marco Gori, Gabriele Monfardini, and Franco Scarselli. 2005. A new model for learning in graph domains. In *Proceedings. 2005 IEEE International Joint Conference on Neural Networks, 2005.*, volume 2, pages 729–734. IEEE.

Aditya Grover and Jure Leskovec. 2016. node2vec: Scalable feature learning for networks. In *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*.

Diederik P Kingma and Max Welling. 2013. Auto-encoding variational bayes. *arXiv preprint arXiv:1312.6114*.

Thomas N Kipf and Max Welling. 2016. Variational Graph Auto-Encoders. *Bayesian Deep Learning Workshop (NIPS 2016)*.

Thomas N. Kipf and Max Welling. 2017. Semi-Supervised Classification with Graph Convolutional Networks. In *5th International Conference on Learning Representations, ICLR 2017, Toulon, France, April 24-26, 2017, Conference Track Proceedings*.

Igor Labutov, Yun Huang, Peter Brusilovsky, and Daqing He. 2017. Semi-Supervised Techniques for Mining Learning Outcomes and Prerequisites. In *Proceedings of the 23rd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pages 907–915. ACM.

Quoc Le and Tomas Mikolov. 2014. Distributed representations of sentences and documents. In *International conference on machine learning*, pages 1188–1196.

Irene Li, Alexander R Fabbri, Robert R Tung, and Dragomir R Radev. 2019. What Should I Learn First: Introducing Lecturebank for NLP Education and Prerequisite Chain Learning. *In 33rd AAAI Conference on Artificial Intelligence (AAAI-19)*.

Chen Liang, Jianbo Ye, Zhaohui Wu, Bart Pursel, and C Lee Giles. 2017. Recovering Concept Prerequisite Relations from University Course Dependencies. In *Thirty-First AAAI Conference on Artificial Intelligence*.

Chen Liang, Jianbo Ye, Shuting Wang, Bart Pursel, and C Lee Giles. 2018. Investigating active learning for concept prerequisite learning. In *Thirty-Second AAAI Conference on Artificial Intelligence*.

Hanxiao Liu, Wanli Ma, Yiming Yang, and Jaime Carbonell. 2016. Learning concept graphs from online educational data. *Journal of Artificial Intelligence Research*, 55:1059–1090.

Tomas Mikolov, Ilya Sutskever, Kai Chen, Greg S Corrado, and Jeff Dean. 2013. Distributed Representations of Words and Phrases and Their Compositionality. In *Advances in neural information processing systems*, pages 3111–3119.

Liangming Pan, Chengjiang Li, Juanzi Li, and Jie Tang. 2017. Prerequisite Relation Learning for Concepts in MOOCs. In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1447–1456.

Bryan Perozzi, Rami Al-Rfou, and Steven Skiena. 2014. Deepwalk: Online learning of social representations. In *Proceedings of the 20th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, KDD '14, pages 701–710, New York, NY, USA. ACM.

Michael Schlichtkrull, Thomas N Kipf, Peter Bloem, Rianne Van Den Berg, Ivan Titov, and Max Welling. 2018. Modeling Relational Data with Graph Convolutional Networks. In *European Semantic Web Conference*, pages 593–607. Springer.

Prithviraj Sen, Galileo Namata, Mustafa Bilgic, Lise Getoor, Brian Galligher, and Tina Eliassi-Rad. 2008. Collective classification in network data. *AI magazine*, 29(3):93–93.

Han Xiao. 2018. bert-as-service. `https://github.com/hanxiao/bert-as-service`.

Bishan Yang, Wen-tau Yih, Xiaodong He, Jianfeng Gao, and Li Deng. 2014. Embedding entities and relations for learning and inference in knowledge bases. *arXiv preprint arXiv:1412.6575*.

Liang Yao, Chengsheng Mao, and Yuan Luo. 2018. Graph Convolutional Networks for Text Classification. *In 33rd AAAI Conference on Artificial Intelligence (AAAI-19)*.

Gianlucca Zuin, Luiz Chaimowicz, and Adriano Veloso. 2018. Learning transferable features for open-domain question answering. In *2018 International Joint Conference on Neural Networks (IJCNN)*, pages 1–8. IEEE.