

Concept Graph Learning from Educational Data

Yiming Yang, Hanxiao Liu, Jaime Carbonell and Wanli Ma

School of Computer Science
Carnegie Mellon University

February 3, 2015

Outline of the Talk

- Motivation
- Concept Representation Schemes
- Concept Graph Learning
- Experiments & Empirical Results
- Future Work

Scenario: Massive course materials are online available from different course providers

- Universities, Coursera, Edx, MIT OpenCourseWare ...

Challenge: How to integrate the scattered information?

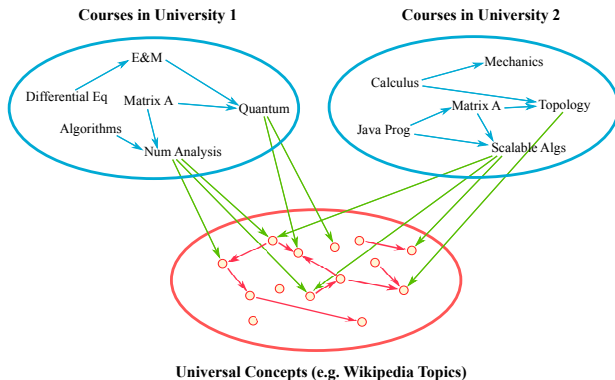
A CMU graduate: *“After completing courses A, B on Coursera, what course shall I take next at CMU?”*

- Lack a method to measure the course overlapping and the course prerequisite relations across institutions.

We address this by putting cross-institutional courses under a canonical language—concept.

Introduction

Concept Graph Learning



Goal: Learning a **universal graph of concepts** based on

- 1 Course-level prerequisite relations
- 2 Concept representation of courses

Outline of the Talk

- Motivation
- **Concept Representation Schemes**
- Learning the Concept Graph
- Experiments & Empirical Results
- Future Work

Representation Schemes

Word-based Representation



Home » Courses » Electrical Engineering and Computer Science » Great Ideas in Theoretical Computer Science

Great Ideas in Theoretical Computer Science

COURSE HOME



Instructor(s)
Prof. Scott Aaronson

↓ crawling and parsing

```
<course>
<id>6.080</id>
<name>Great Ideas in Theoretical Computer Science</name>
<tag>Electrical Engineering and Computer Science</tag>
<description>This course provides a challenging introduction to some of the central ideas of
theoretical computer science. It attempts ... </description>
<keywords>computer science,theoretical computer science,logic,turing machines,computability,
finite automata,godel,complexity,polynomial time,efficient algorithms ... </keywords>
<calendar>Introduction # Logic # Circuits and finite automata # Turing machines # Reducibility
and Godel # Minds and machines # Complexity # Polynomial time # P and NP # NP-
completeness # NP-completeness in practice # Space complexity and more ... </calendar>
</course>
```

Concepts $\stackrel{\text{def}}{=} \text{Words}$

- Vocabularies not controlled
 - CMU 10-715: shattering coefficient; MIT 15.097: growth function
- Words are in multiple granularities \implies interpretability ↓

Representation Schemes

Latent Space Representation

Schemes based on dimensionality reduction

- Sparse Coding of Words
 - Trained on the given courses—purely unsupervised
- Distributed Word Embedding
 - Trained on Wikipedia articles—leverages exterior info

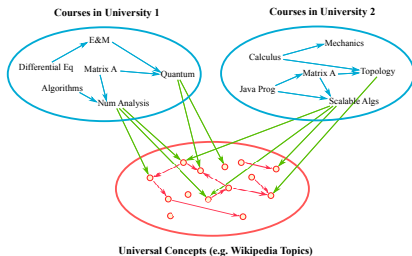
Concepts $\stackrel{\text{def}}{=} \text{Dimensionality-reduced vectors}$

- Controlled “vocabulary”
 - Words are mapped onto a unified latent space
- Concept granularity can be controlled by latent dimensionality
- Hard to interpret

Outline of the Talk

- Motivation
- Concept Representation Schemes
- Learning the Concept Graph
- Experiments & Empirical Results
- Future Work

Problem Formulation



Observed course-level relations \mathcal{O}
Concept representation of courses \mathbf{X} n by p
Concept graph \mathbf{A} p by p

How to evaluate \mathbf{A} ?

- 1 Map \mathbf{A} to an estimated course graph $\hat{\mathcal{O}}$ (n by n) via \mathbf{X} .
- 2 Then, evaluate the quality of $\hat{\mathcal{O}}$ with \mathcal{O} .

How to map the concept graph \mathbf{A} to an estimated course graph Θ ?

—through a bilinear form:

$$\Theta \stackrel{\text{def}}{=} \mathbf{XAX}^\top$$

Explanation:

$$\theta_{ij} \stackrel{\text{def}}{=} \sum_{u,v} a_{uv} x_{iu} x_{jv}$$

- θ_{ij} : strength from course j to course i
- a_{uv} : strength from concept v to concept u

Each course-level prerequisite θ_{ij} is defined as the cumulative effect of multiple concept-level prerequisites $\sum_{u,v} a_{uv} x_{iu} x_{jv}$

How to evaluate the estimated course graph Θ with our observed course-level relations \mathcal{O} ?

i.e. how to define the loss function over Θ w.r.t. \mathcal{O} ?

Problem

- Only positive examples are available
- Treating unobserved course relations as negative examples leads to highly skewed label set

Solution: ranking

- We hope $\theta_{ij} > \theta_{ik}$ if $j \in \text{prereq}(i)$ and $k \notin \text{prereq}(i)$

CGL objective with p^2 variables

$$\begin{aligned} \min_{\mathbf{A} \in \mathcal{R}^{p \times p}} \quad & C \sum_{(i,j,k)} \ell(\theta_{ij} - \theta_{ik}) + \frac{1}{2} \|\mathbf{A}\|_F^2 \\ \text{s.t.} \quad & \Theta = \mathbf{XAX}^\top \end{aligned}$$

Problem

- \mathbf{A} can be a huge dense matrix (e.g. is 15,396 by 15,396 for words-based concept representation)
- Dual space? #dual variables = $O(n^3)$

Solution: derive an equivalent optimization problem with only n^2 ($n^2 \ll p^2$) variables.

Introduce slack variable $\mathbf{S} \in \mathcal{R}^{n \times n}$ for constraint $\mathbf{\Theta} = \mathbf{XAX}^\top$.

The Lagrangian is

$$\mathcal{L} = C \sum_{(i,j,k)} \ell(\theta_{ij} - \theta_{ik}) + \frac{1}{2} \|\mathbf{A}\|_F^2 + \langle \mathbf{S}, \mathbf{\Theta} - \mathbf{XAX}^\top \rangle$$

$\frac{\partial \mathcal{L}}{\partial \mathbf{A}}$ should vanish at the stationary point

- $\implies \mathbf{A}^* = \mathbf{X}^\top \mathbf{S}^{*\top} \mathbf{X}$
- $\mathbf{A}^* \in \mathcal{R}^{p \times p}$ only has n^2 ($n^2 \ll p^2$) degrees of freedom!

Equivalent CGL objective with n^2 variables

$$\begin{aligned} \min_{\mathbf{S} \in \mathcal{R}^{n \times n}} \quad & C \sum_{(i,j,k)} \ell(\theta_{ij} - \theta_{ik}) + \frac{1}{2} \text{tr}(\mathbf{\Theta} \mathbf{S}^\top) \\ \text{s.t.} \quad & \mathbf{\Theta} = \mathbf{KSK} \end{aligned}$$

We choose the squared hinge loss $\ell(x) = (\max(1 - x, 0))^2$

- large-margin property: strong generalization ability
- smoothness: allows Nesterov's accelerated GD
 - GD: 37.3min & 1490 iterations on MIT
 - accelerated GD: 3.08 min & 103 iterations on MIT

An alternative of GD: Inexact Newton Method

- To avoid the huge Hessian—use a matrix-free Conjugate Gradient to compute the Newton direction

Outline of the Talk

- Motivation
- Concept Representation Schemes
- Learning the Concept Graph
- Experiments & Empirical Results
- Future Work

Table : Datasets Statistics¹

University	Department	# Courses	# Prerequisites	# Words
MIT ²	*	2322	1173	15396
Caltech	*	1048	761	5617
CMU	CS, STATS	83	150	1955
Princeton	MATH	56	90	454

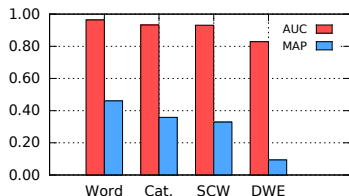
Metrics for evaluation: MAP and AUC

¹available at <http://nyc.lti.cs.cmu.edu/teacher/dataset/>

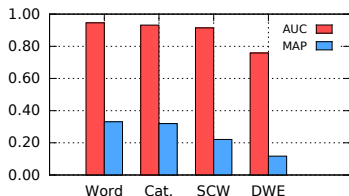
²MIT OpenCourseWare <http://ocw.mit.edu/index.htm>

Experiments

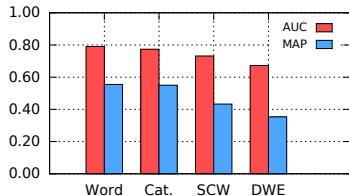
Comparison among Concept Representation Schemes



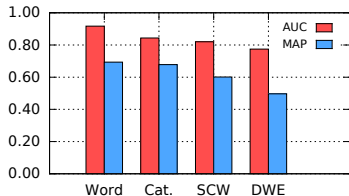
(a) CGL.Rank on MIT Data



(b) CGL.Rank on Caltech Data



(c) CGL.Rank on CMU Data



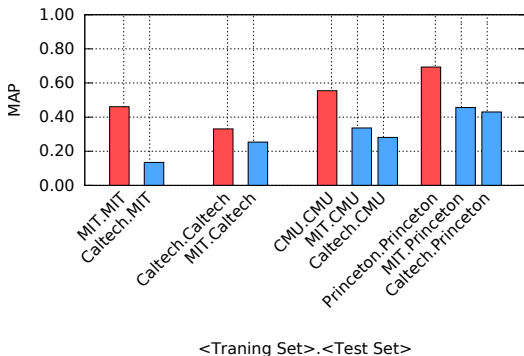
(d) CGL.Rank on Princeton Data

Words \succeq Categories \succ Sparse Coding \succ Distributed Word Embedding

Experiments

Cross-institutional Prerequisite Prediction

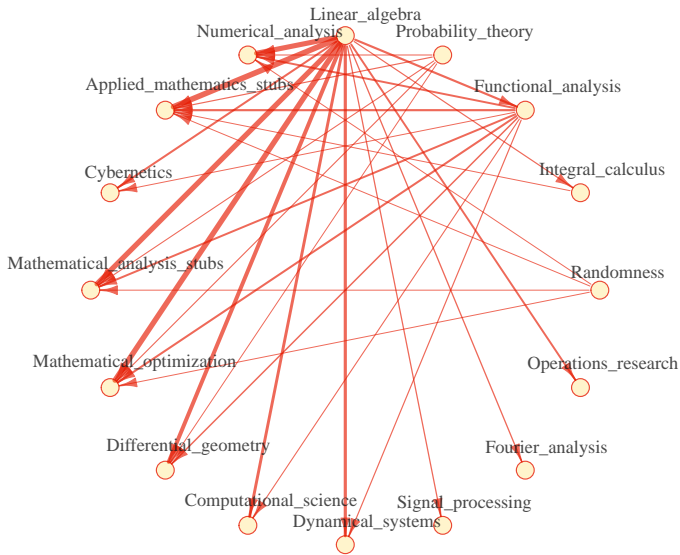
A good concept graph should be universal, thus should be transferable across different institutions



- There is always a performance loss if we go across institutions.
- We do get good transfer.

Empirical Results

Concept Graph for MIT



Outline of the Talk

- Motivation
- Concept Representation Schemes
- Learning the Concept Graph
- Experiments & Empirical Results
- **Future Work**

- Deploying the induced concept graph for personalized curriculum planning (on-going work)
 - Student's academic background/goal $\stackrel{\text{def}}{=} \text{bag-of-concepts}$
 - Find an optimal sequence of courses?
- Cross-language transfer learning by using Wikipedia categories (concepts) as the interlingua.

Thanks!

hanxiaol@cs.cmu.edu