

WHAT ARE WE TEACHING?

Automated Evaluation of CS Curricula Content Using Topic Modeling

Jean Michel Rouly

April 30, 2015

George Mason University

OVERVIEW

- From 2007, number of new CS undergrads has increased 28.9%
- Fifth straight year of increase

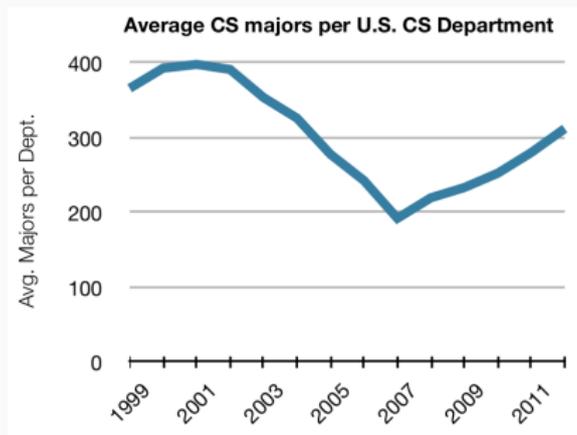


Figure 1: From Zweben, 2011.

Evaluating a department on its conceptual coverage involves...

Evaluating a department on its conceptual coverage involves...

Relative Standing comparing against other, similar departments

Evaluating a department on its conceptual coverage involves...

Relative Standing comparing against other, similar departments

Absolute Performance benchmark against standardized expectations (published by the ACM)

Evaluating a department on its conceptual coverage involves...

Relative Standing comparing against other, similar departments

Absolute Performance benchmark against standardized expectations (published by the ACM)

Both of these require data about the topics covered in a course.

What we've got:

What we've got:

1. Widely available university course description data.

What we've got:

1. Widely available university course description data.
2. Descriptions detail what concepts are taught in a course.

What we've got:

1. Widely available university course description data.
2. Descriptions detail what concepts are taught in a course.
3. **Human-readable** descriptions require manual inspection.

What we've got:

1. Widely available university course description data.
2. Descriptions detail what concepts are taught in a course.
3. **Human-readable** descriptions require manual inspection.

So what to do?

Through the application of probabilistic machine learning methods, specifically LDA topic modeling, a corpus of unstructured course descriptions can be digested and mined for topics. In this scenario, each topic represents a core concept covered by the courses.

A research framework will be constructed to read data from the Internet, digest into a common internal format, pipeline into an LDA topic model, and ultimately visualize in an interactive manner.

Ultimately the automatically discovered topics can be used in end-user university evaluation processes.

In other words

Build a tool that automatically...

In other words

Build a tool that automatically...

- ingests large collection of course descriptions

In other words

Build a tool that automatically...

- ingests large collection of course descriptions
- **infers topics from course descriptions**

In other words

Build a tool that automatically...

- ingests large collection of course descriptions
- infers topics from course descriptions
- **computes comparisons between departments**

In other words

Build a tool that automatically...

- ingests large collection of course descriptions
- infers topics from course descriptions
- computes comparisons between departments
- **evaluates departments on their concepts**

BACKGROUND

Topic Modeling

Attempts to discover the abstract **topics** of a dataset.

Topics

A **topic** is a frequency distribution over terms, roughly representing a concept taught in a course.

Overview

Latent Dirichlet Allocation (LDA) is a form of *Latent Variable Modeling* that can infer topics from within a document.

LDA takes a generative approach to latent variable modeling, assuming the topics occur in some proportion within each document.

DATA

University	Course Count	Web
American University	32	american.edu
George Mason University	145	gmu.edu
Kansas State University	83	ksu.edu
Louisiana State University	59	lsu.edu
Portland State University	190	pdx.edu
Rensselaer Polytechnic Institute	61	rpi.edu
University of South Carolina	64	sc.edu
Stanford University	69	stanford.edu
University of Utah	142	utah.edu
University of Tennessee, Knoxville	29	utk.edu
ACM Exemplar Courses	68	—

Table 1: University course descriptions

Raw Course Description

Capstone course focusing on design and successful implementation of major software project, encompassing broad spectrum of knowledge and skills, developed by team of students. Requires final exhibition to faculty-industry panel.

Cleaned course description

capston focus design success implement major softwar project
encompass broad spectrum knowledg skill develop team student
requir final exhibit faculti industri panel

TRAJECTORY

Trajectory is a tool that automatically ingests course description data from the Internet and presents an accessible interface for departmental evaluation.

Lines of Python	3193
Lines of Java	631
Lines of HTML/CSS/JS	1828
Lines of JSON	3219
Lines of Bash	165
Size on disk	6.7M

Table 2: Code statistics

Four primary modules:

Scrape web-scrape online university catalogs

Import/Export pass structured data between Learn and Scrape

Learn estimate LDA topic model on data

Web visualization tool

Underneath the entire system is a structured relational database layer.

Browse collected data by university or department

Understand courses through inferred topics

Analyze conceptual overlap in prerequisite chains

Compare departments based on conceptual composition

Evaluate departments against ACM benchmarks

(NOT SO) LIVE DEMO

Images have been modified to fit in this presentation.

Visit trajectory.rouly.net on your smart device to follow along.

BROWSE

#	Department	Course Count	Web
1	Applied Information Technology	39	-
2	Computer Science	106	-

Figure 2: Browse a university's departments.

Course List

#	Number	Title
1	100	Principles of Computing
2	101	Preview of Computer Science
3	105	Computer Ethics and Society
4	112	Introduction to Computer Programming
5	123	Computing: From the Abacus to the Web
6	211	Object-Oriented Programming
7	222	Computer Programming for Engineers

Figure 3: Browse a department's courses.

UNDERSTAND

GMU CS 105

Computer Ethics and Society

/ UNIVERSITIES / GMU / CS / 105

Inferred Topics

▶ (0.759646) ethic, comput, issu, profession, social, technolog, privaci, legal, relat, digit

29

Figure 4: View a course's inferred topics.

Inferred Topics

▶ (0.734343) comput, ethic, issu, social, profession, technolog, impact, privati, properti, legal 34

(0.912) **PDX CS 305** *Social, Ethical, and Legal Implications of Computing*

(0.833) **SC CSCE 390** *Professional Issues in Computer Science and Engineering*

(0.734) **GMU CS 105** *Computer Ethics and Society*

(0.680) **KSU CIS 415** *Ethics and Computing Technology*

(0.670) **LSU CSC 1200** *Ethics in Computing*

(0.629) **Stanford CS 181** *Computers, Ethics, and Public Policy*

(0.625) **SC CSCE 190** *Computing in the Modern World*

Figure 5: View related courses.

Inferred Topics

▶ (0.748265) data, mine, algorithm, larg, web, effici, techniqu, search, time, cluster

16

(0.748) RPI CSCI 6390 Database Mining

(0.714) Utah CS 6140 Data Mining

(0.695) RPI CSCI 4390 Database Mining

(0.644) Utah CS 5140 Data Mining

(0.619) Stanford CS 276 Information Retrieval and Web Search (LINGUIST 286)

(0.614) LSU CSC 7442 Data Mining and Knowledge Discovery

(0.588) GMU CS 674 Data Mining on Multimedia Data

Figure 6: View related courses.

Inferred Topics

▶ (0.402983) game, video, player, film, solut, graphic, emphasi, creat, anim, varieti

19

(0.703) **PDX CS 442** *Advanced Artificial Intelligence: Combinatorial Games*

(0.660) **PDX CS 542** *Advanced Artificial Intelligence: Combinatorial Games*

(0.509) **Utah CS 6665** *Character Animation*

(0.433) **Utah CS 6660** *Physics-based Animation*

(0.403) **GMU CS 225** *Culture and Theory of Games*

(0.365) **KSU CIS 580** *Fundamentals of Game Programming*

(0.351) **RPI CSCI 4540** *Game Development II*

Figure 7: View related courses.

ANALYZE

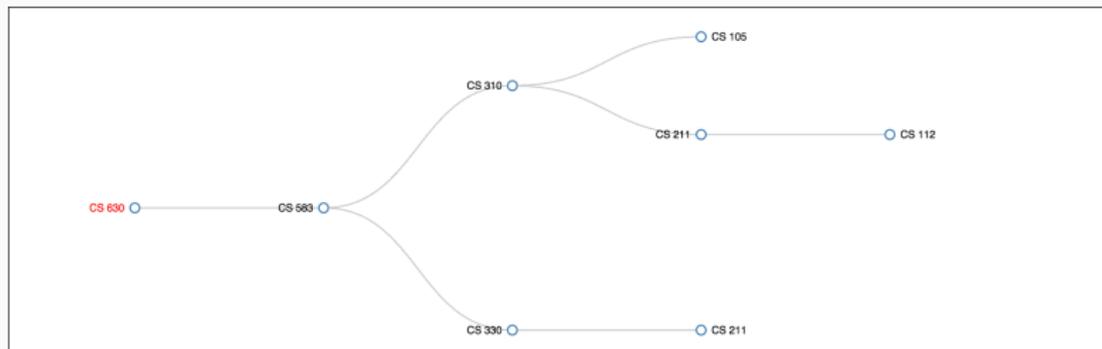


Figure 8: Course prerequisite tree.

COMPARE

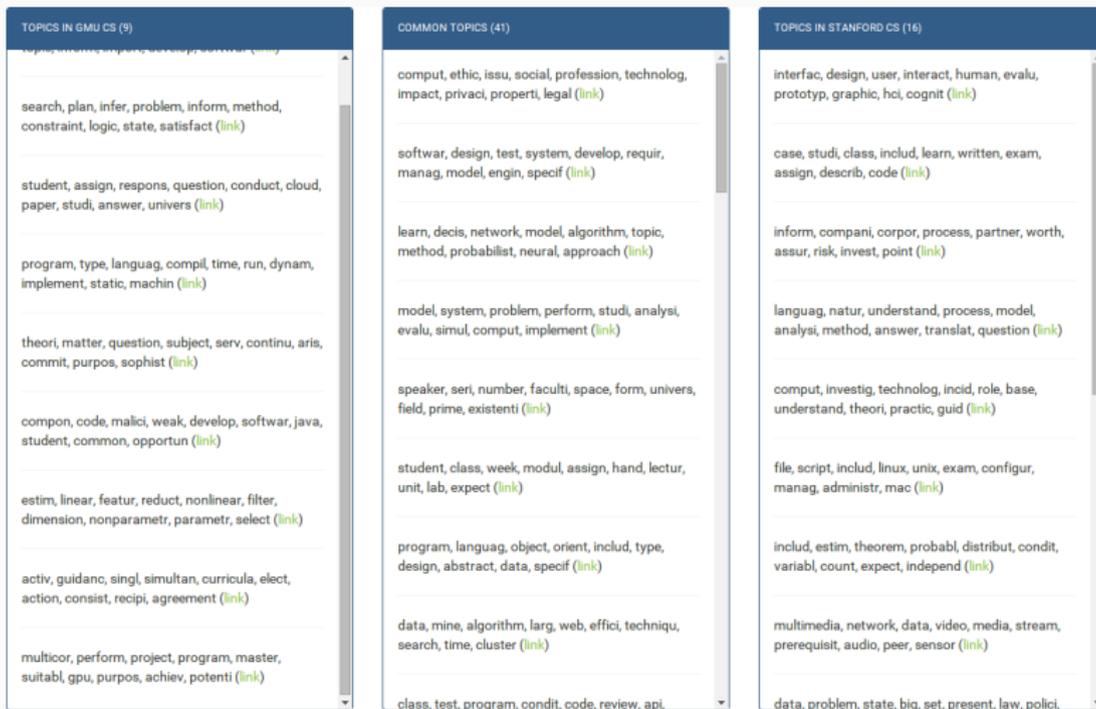


Figure 9: Compare departmental concept coverage.

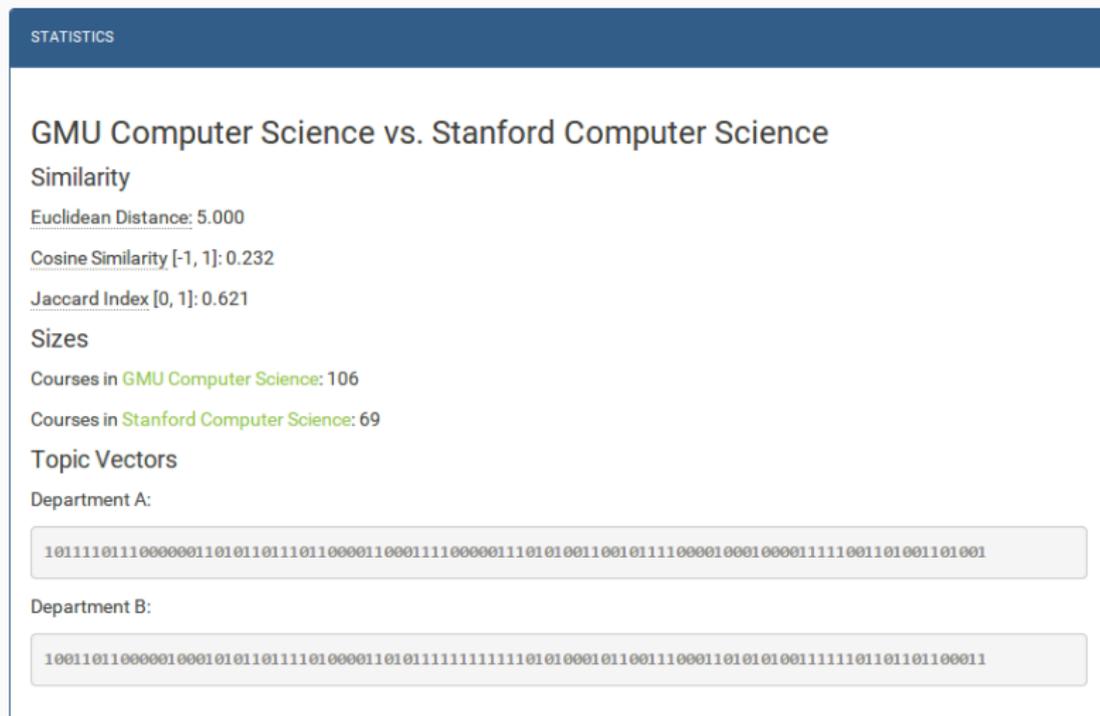


Figure 10: Quick departmental comparison.

EVALUATE

#	Number	Title	Knowledge Areas	Ground Truth
1	100	Principles of Computing	AL, CN, SP, IM	None.
2	101	Preview of Computer Science	CN	None.
3	105	Computer Ethics and Society	SE, SP	IAS, NC, SP
4	112	Introduction to Computer Programming	AL, CN, SP, IM	AL, AR, CN, PL, SDF
5	123	Computing: From the Abacus to the Web	SP	None.
6	211	Object-Oriented Programming	AL, CN, IM	AL, AR, CN, OS, PL
7	222	Computer Programming for Engineers	AL, CN, IM	PL
8	225	Culture and Theory of Games	SP	None.

Figure 11: Evaluate department against Knowledge Areas.

WRAP UP

- Learning Outcomes meta-analysis
- Alternative methods
- Topic summarization
- Extension to workforce preparedness
- Smart web scraping

Try out **Trajectory** online at

`trajectory.rouly.net`

Get the source of this presentation and the **Trajectory** project at

`github.com/jrouly/trajectory`

Trajectory is licensed under the Apache version 2.0 license.

Co-authors:

- Huzefa Rangwala
- Aditya Johri

Presentation theme:

- github.com/matze/mtheme

- D. M. Blei, A. Y. Ng, and M. I. Jordan. Latent Dirichlet Allocation. *J. Mach. Learn. Res.*, 3:993–1022, Mar. 2003.
- ACM/IEEE-CS Joint Task Force on Computing Curricula. Computer science curricula 2013. Technical report, ACM Press and IEEE Computer Society Press, December 2013.
- S. Zweben. Computing degree and enrollment trends. Technical report, Computing Research Association, 2011.

QUESTIONS?