



Hatch, Match and Dispatch:

Examining the relationship between student intent, expectations, behaviours and outcomes in six Coursera MOOCs at the University of Toronto



UNIVERSITY OF
TORONTO



The Team

Bodong Chen (@bodongchen)

Stian Håklev (@houshuang)

William Heikoop

Laurie Harrison (@IT4learning)

Hedieh Najafi

Carol Rolheiser

Chris Teplovs (@cteplovs)

Research Objectives

In the context of 6 University of Toronto MOOCs, we examined:

Student intentions, goals and motivation (Hatch) and their relationship (Match) to student behaviour and student performance based on formative and summative assessment (Dispatch)

Research Design

Multiple case study of six MOOCs:

- Aboriginal Worldviews and Education
- Introduction to Psychology
- Learn to Program: The Fundamentals
- Learn to Program: Crafting Quality Code
- The Social Context of Mental Health and Illness
- Statistics: Making Sense of Data

Research Design

Multiple case study of six MOOCs:

- Aboriginal Worldviews and Education
- **Introduction to Psychology**
- Learn to Program: The Fundamentals
- Learn to Program: Crafting Quality Code
- The Social Context of Mental Health and Illness
- Statistics: Making Sense of Data

Data Sources

1. Entry Survey (demographic info, motivation, expected engagement)
2. Student behaviour data (clickstream data)
3. Outcome data (quiz scores, final grades)
4. Exit survey
5. Instructor interviews

Data Sources

1. **Entry Survey** (demographic info, **motivation, expected engagement**)
2. **Student behavior data (clickstream data)**
3. Outcome data (quiz scores, **final grades**)
4. Exit survey
5. Instructor interviews

Data Analysis

Part I: Entry survey (n=70,418 for all MOOCs)

Principal component analysis (PCA) of surveys to reduce complexity of categorization.

Two axes: “betterment” and “enjoyment”

Part II: Clickstream data (n=16,900,926 events for Introductory Psychology only)

Sequential pattern mining

Part III: Comparing frequent patterns from Part II based on classification from Part I

Data Analysis Part Ia: Entry Surveys

This survey is intended for learners using this course as an “archived” MOOC. We consider an “archived” MOOC to be a course in which course materials are available, but for which deadlines have passed. The alternative is a “live” MOOC for which the due dates are upcoming. Your feedback will help us understand how our MOOCs can be better designed to meet the diverse needs and interests of our learners.

1. Which of the following descriptions best characterizes you?

- Student
- Professional
- Researcher
- Academic / Professor
- Computer Engineering
- Lifelong Learner
- None of the above

2. Why did you enroll in this course? For each reason below, please rate on the scale of not at all important to very important.

	Strongly Disagree	Disagree	Neutral	Agree	Strongly Agree
The subject is relevant to my academic field of study					
This class teaches skills that will help my job/career					
I want to earn some sort of credential that I can use to enhance my CV / resume					
Because this course is offered by a prestigious university					
I think taking this course will be fun and enjoyable					

Data Analysis Part Ib: Principal Components Analysis of Survey Data

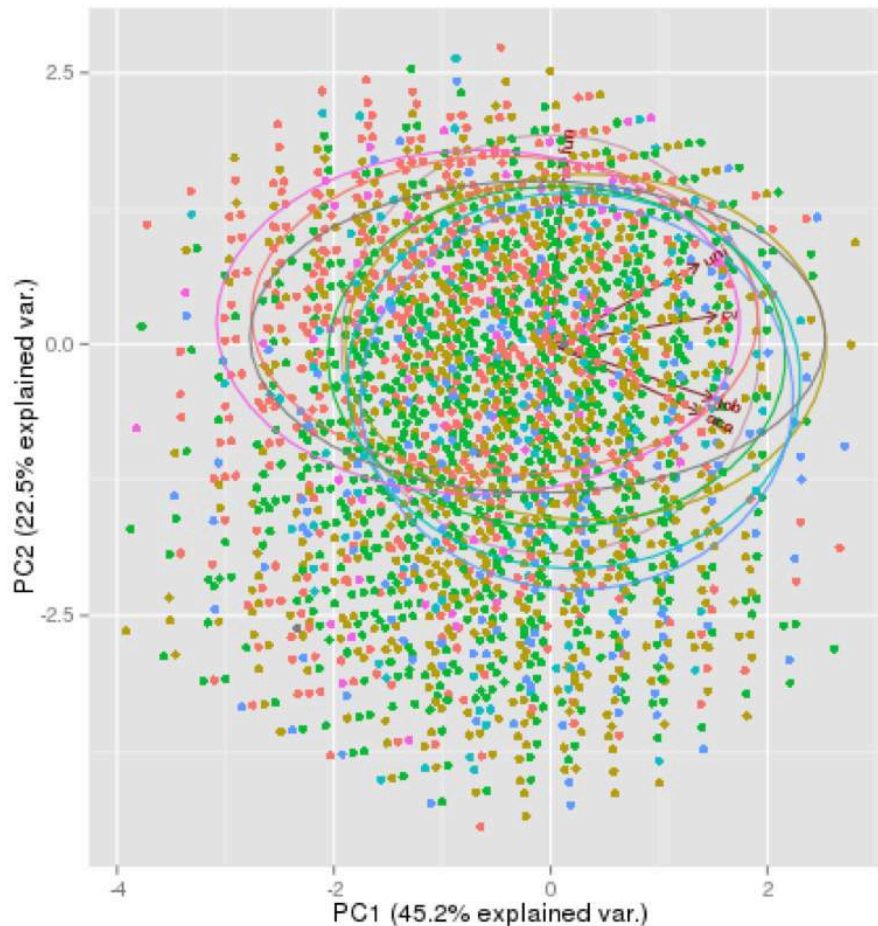
##		PC1	PC2	PC3	PC4	PC5
##	aca	0.4782	-0.3096	0.4917	-0.53174	-0.38852
##	job	0.5187	-0.2458	0.3214	0.48693	0.57459
##	cv	0.5273	0.1293	-0.4000	0.45436	-0.58206
##	uni	0.4714	0.3607	-0.4558	-0.51393	0.41923
##	fun	0.0447	0.8348	0.5359	0.09786	-0.06596

2. Why did you enroll in this course? For each reason below, please rate on the scale of not at all important to very important.

	Strongly Disagree	Disagree	Neutral	Agree	Strongly Agree
The subject is relevant to my academic field of study					
This class teaches skills that will help my job/career					
I want to earn some sort of credential that I can use to enhance my CV / resume					
Because this course is offered by a prestigious university					
I think taking this course will be fun and enjoyable					

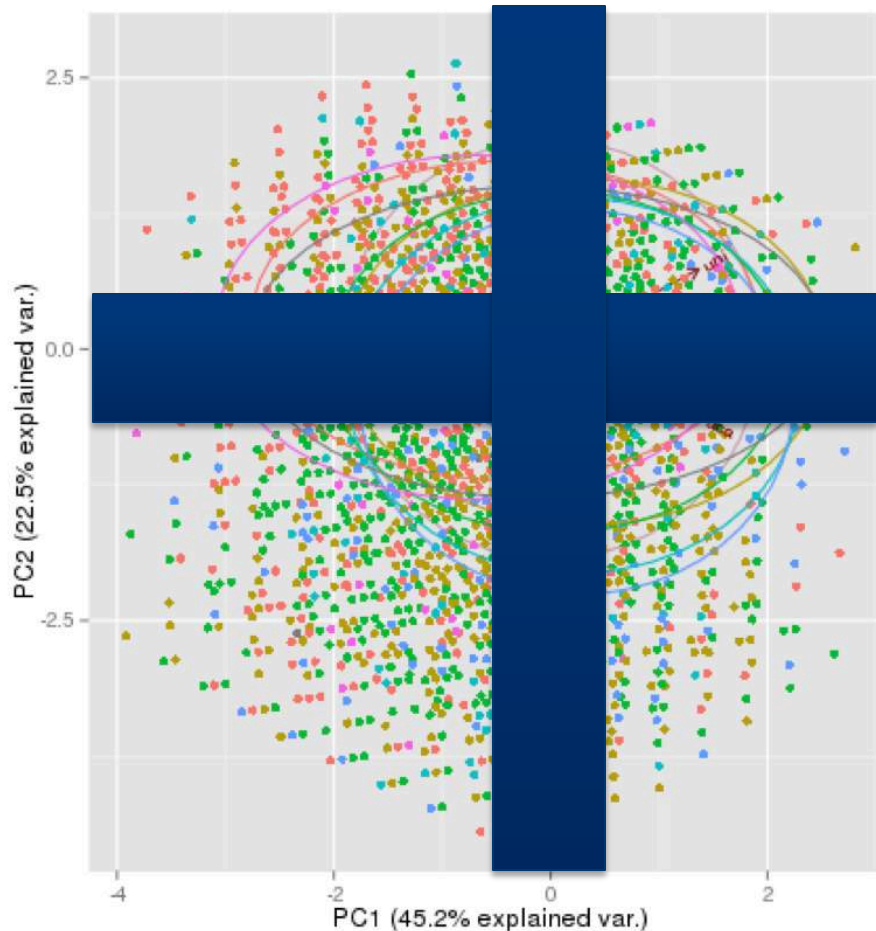
Data Analysis Part I: PCA continued

— Lifelong Learner — Student — Professional — Academic/Professor — Researcher — None of the at

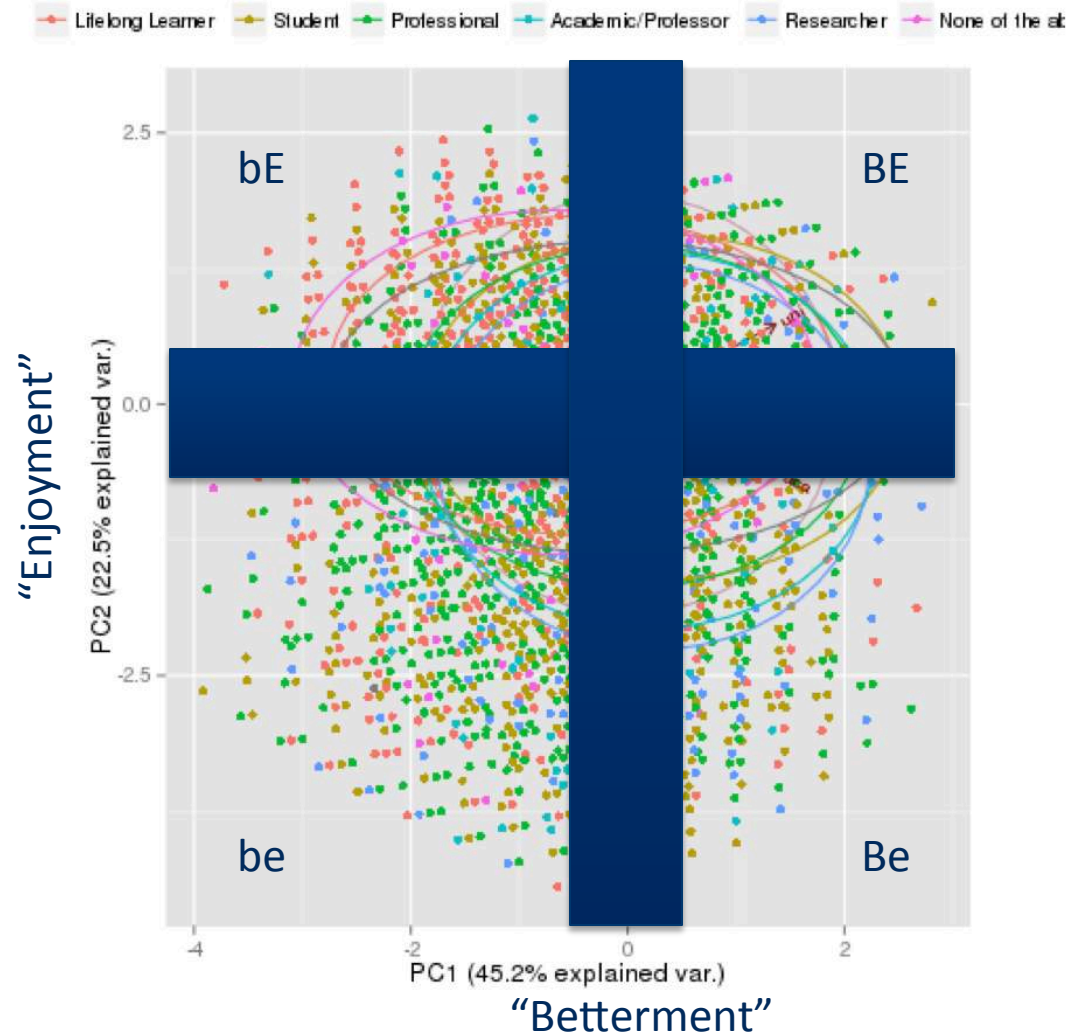


Data Analysis Part 1c: PCA continued

Life-long Learner Student Professional Academic/Professor Researcher None of the at



Data Analysis Part 1c: PCA continued



Data Analysis Part IIa: Clickstream Data

3 Clickstream data export

The clickstream data are given as a single gzipped text file. The text file contains a listing of clickstream events, one per line. Each entry consists of a single JSON-serialized object with the following fields:

- key: a string describing a particular kind of event
- value: a string containing metadata for the event
- username: anonymized user ID (anon_user_id; missing if individual not logged in)
- timestamp: POSIX timestamp indicating when event occurred
- page_url: the webpage associated with an event
- client: a keyword describing the context of an event
- session: browser session cookie
- language: the client browser's language preference
- from: the referer URL
- user_ip: IP address of user
- user_agent: browser user agent string

In general, the export contains two types of events: video-related events and page view events. In both cases, the value field itself is a JSON-serialized object that contains metadata associated with the event. For example, for video events, the value field contains information about the specific type of action taken (e.g., play, pause, seek) and current video settings (e.g., current playback speed).

← Info from Coursera

This is what clickstream data looks like

```
{"key": "pageview", "value": {}, "username": "e2bb8a569a2f7fee440e9747b2d6631db49af802", "timestamp": 1377380881722, "page_url": "https://class.coursera.org/intropsych-001/lecture/view?lecture_id=5", "client": "spark", "session": "5938170350-1374957666948", "language": "en-US,en;q=0.8", "from": "https://class.coursera.org/intropsych-001/lecture/view?lecture_id=5", "user_ip": "5.82.190.197", "user_agent": "Mozilla/5.0 (Windows NT 6.1) AppleWebKit/537.36 (KHTML, like Gecko) Chrome/29.0.1547.57 Safari/537.36", "12": [{"height": 600, "width": 1024}], "13": [0], "14": ["https://www.coursera.org/course/intropsych"]}
{"key": "user.video.lecture.action", "value": {"currentTime": 0, "playbackRate": 1, "paused": false, "error": null, "networkState": 3, "readyState": 0, "endTimeStamp": 137738088225, "initTimestamp": 137738088186, "type": "play", "prevTime": 0}, "username": "e2bb8a569a2f7fee440e9747b2d6631db49af802", "timestamp": 1377380883017, "page_url": "https://class.coursera.org/intropsych-001/lecture/view?lecture_id=5", "client": "spark", "session": "5938170350-1374957666948", "language": "en-US,en;q=0.8", "from": "https://class.coursera.org/intropsych-001/lecture/5", "user_ip": "5.82.190.197", "user_agent": "Mozilla/5.0 (Windows NT 6.1) AppleWebKit/537.36 (KHTML, like Gecko) Chrome/29.0.1547.57 Safari/537.36", "12": [{"height": 600, "width": 1024}], "13": [0], "14": ["https://www.coursera.org/course/intropsych"]}
{"key": "pageview", "value": {}, "username": "87cd7dec5a9ff4b9f4e1904bbc44552c1f5c5b1e", "timestamp": 1377380989756, "page_url": "https://class.coursera.org/intropsych-001/lecture/view?lecture_id=19", "client": "spark", "session": "7746927522-1376902134923", "language": "da-dk", "from": "https://class.coursera.org/intropsych-001/lecture/19", "user_ip": "95.166.8.167", "user_agent": "Mozilla/5.0 (iPad; CPU OS 6_1_3 like Mac OS X) AppleWebKit/536.26 (KHTML, like Gecko) Version/6.0 Mobile/10B329 Safari/8536.25", "12": [{"height": 1024, "width": 768}], "13": [0]}
{"key": "pageview", "value": {}, "username": "87cd7dec5a9ff4b9f4e1904bbc44552c1f5c5b1e", "timestamp": 1377380989764, "page_url": "https://class.coursera.org/intropsych-001/lecture/view?lecture_id=19", "client": "spark", "session": "7746927522-1376902134923", "language": "da-dk", "from": "https://class.coursera.org/intropsych-001/lecture/view?lecture_id=19", "user_ip": "95.166.8.167", "user_agent": "Mozilla/5.0 (iPad; CPU OS 6_1_3 like Mac OS X) AppleWebKit/536.26 (KHTML, like Gecko) Version/6.0 Mobile/10B329 Safari/8536.25", "12": [{"height": 1024, "width": 768}], "13": [0]}
```

Data Analysis Part IIb: Clickstream Data Wrangling Workflow

```
{
  "key": "pageview",
  "value": {},
  "username": "e2bb8a569a2f7fee440e9747b2d6631db49af802",
  "timestamp": 1377380881722,
  "page_url": "https://class.coursera.org/intropsych-001/lecture/view?lecture_id=5",
  "client": "spark",
  "session": "5938170350-1374957666948",
  "language": "en-US,en;q=0.8",
  "from": "https://class.coursera.org/intropsych-001/lecture/view?lecture_id=5",
  "user_ip": "5.82.190.197",
  "user_agent": "Mozilla/5.0 (Windows NT 6.1) AppleWebKit/537.36 (KHTML, like Gecko) Chrome/29.0.1547.57 Safari/537.36",
  "12": [{"height": 600, "width": 1024}],
  "13": [0],
  "14": ["https://www.coursera.org/course/intropsych"]
},
{
  "key": "user.video.lecture.action",
  "value": {"currentTime": 0, "playbackRate": 1, "paused": false, "error": null, "networkState": 3, "readyState": 0, "eventTimestamp": 137738088225, "initTimestamp": 137738088186, "type": "play", "prevTime": 0},
  "username": "e2bb8a569a2f7fee440e9747b2d6631db49af802",
  "timestamp": 1377380883017,
  "page_url": "https://class.coursera.org/intropsych-001/lecture/view?lecture_id=5",
  "client": "spark",
  "session": "5938170350-1374957666948",
  "language": "en-US,en;q=0.8",
  "from": "https://class.coursera.org/intropsych-001/lecture/view?lecture_id=5",
  "user_ip": "5.82.190.197",
  "user_agent": "Mozilla/5.0 (Windows NT 6.1) AppleWebKit/537.36 (KHTML, like Gecko) Chrome/29.0.1547.57 Safari/537.36",
  "12": [{"height": 600, "width": 1024}],
  "13": [0],
  "14": ["https://www.coursera.org/course/intropsych"]
},
{
  "key": "pageview",
  "value": {},
  "username": "56",
  "timestamp": 1377380989764,
  "page_url": "https://class.coursera.org/intropsych-001/lecture/view?lecture_id=19",
  "client": "spark",
  "session": "7746927522-1376902134923",
  "language": "da-dk",
  "from": "https://class.coursera.org/intropsych-001/lecture/view?lecture_id=19",
  "user_ip": "95.166.8.167",
  "user_agent": "Mozilla/5.0 (iPad; CPU OS 6_1_3 like Mac OS X) AppleWebKit/536.26 (KHTML, like Gecko) Version/6.0 Mobile/10B329 Safari/8536.25",
  "12": [{"height": 1024, "width": 768}],
  "13": [0]
},
{
  "key": "pageview",
  "value": {},
  "username": "180750175-1368679630336",
  "timestamp": 1377380989764,
  "page_url": "https://class.coursera.org/intropsych-001/lecture/view?lecture_id=19",
  "client": "spark",
  "session": "7746927522-1376902134923",
  "language": "da-dk",
  "from": "https://class.coursera.org/intropsych-001/lecture/view?lecture_id=19",
  "user_ip": "95.166.8.167",
  "user_agent": "Mozilla/5.0 (iPad; CPU OS 6_1_3 like Mac OS X) AppleWebKit/536.26 (KHTML, like Gecko) Version/6.0 Mobile/10B329 Safari/8536.25",
  "12": [{"height": 1024, "width": 768}],
  "13": [0]
}
```

1. Raw Clickstream



Hadoop & Hive

```
1368679633700 00037b2d03e0700950055ec0268bf62341c93491 180750175-1368679630336 pageview https://class.coursera.org/intropsych-001/class/index {}
1368679636923 00037b2d03e0700950055ec0268bf62341c93491 180750175-1368679630336 pageview https://class.coursera.org/intropsych-001/lecture/index {}
1368679683200 00037b2d03e0700950055ec0268bf62341c93491 180750175-1368679630336 pageview https://class.coursera.org/intropsych-001/lecture/view?lecture_id=45 {}
1368679683905 00037b2d03e0700950055ec0268bf62341c93491 180750175-1368679630336 user.video.lecture.action https://class.coursera.org/intropsych-001/lecture/view?lecture_id=45 {"currentTime": 0, "playbackRate": 1, "paused": false, "error": null, "networkState": 3, "readyState": 0, "eventTimestamp": 137738088225, "initTimestamp": 137738088186, "type": "play", "prevTime": 0}
```

2. Cleaned Clickstream



Python

```
0001b4ecd4595d7f2c9611327e65c33cdabdd01c 1367975197 1 class/index
0001b4ecd4595d7f2c9611327e65c33cdabdd01c 1367975227 1 lecture/index
0001b4ecd4595d7f2c9611327e65c33cdabdd01c 1367975231 1 lecture/view?lecture_id=3
0001b4ecd4595d7f2c9611327e65c33cdabdd01c 1367975958 1 lecture/view?lecture_id=5
0001b4ecd4595d7f2c9611327e65c33cdabdd01c 1367976244 1 lecture/index
0001b4ecd4595d7f2c9611327e65c33cdabdd01c 1368491221 1 class/index
0001b4ecd4595d7f2c9611327e65c33cdabdd01c 1367883156 1 class/index
0001b4ecd4595d7f2c9611327e65c33cdabdd01c 1367883164 1 wiki/view?page=ToDoThisWeek
0001b4ecd4595d7f2c9611327e65c33cdabdd01c 1367883261 1 quiz/index?quiz_type=survey
0001b4ecd4595d7f2c9611327e65c33cdabdd01c 1367883269 1 quiz/attempt?quiz_id=7
0001b4ecd4595d7f2c9611327e65c33cdabdd01c 1367883458 1 quiz/feedback?submission_id=15652
```

3. "Buckets"

4,804,167 transactions

58,691 sequences

Data Analysis Part IIc: Clickstream Analysis (Sequential Pattern Mining)

```
Console ~/
Restarting R session...

Loading required package: arulesSequences
Loading required package: arules
Loading required package: Matrix
Loading required package: lattice

Attaching package: 'arules'

The following objects are masked from 'package:base':

  %in%, write

> load("~/RData")
> summary(ip1p1fs)
set of 60 sequences with

most frequent items:
  quiz/attempt?quiz_id=7 wiki/view?lecture/view?lecture_id=3
  16
  8

most frequent elements:
  {quiz/attempt?quiz_id=7} {wiki/view?lecture/view?lecture_id=3}
  16
  8

element (sequence) size distribution:
sizes
 1 2 3
24 35 1

sequence length distribution:
lengths
 1 2 3
24 35 1

summary of quality measures:
  support
Min. :0.5012
1st Qu.:0.5398
Median :0.5875
Mean :0.6287
3rd Qu.:0.6594
Max. :0.9951

mining info:
 data ntransactions nsequences support
ip1p1 293984 1622 0.
```

Using the 'arulesSequences' package in R, which implements Zaki (2001)...

M. J. Zaki. (2001). SPADE: An Efficient Algorithm for Mining Frequent Sequences. Machine Learning Journal, 42, 31–60.

```
Console ~/
> as(ip1p1fs,"data.frame")

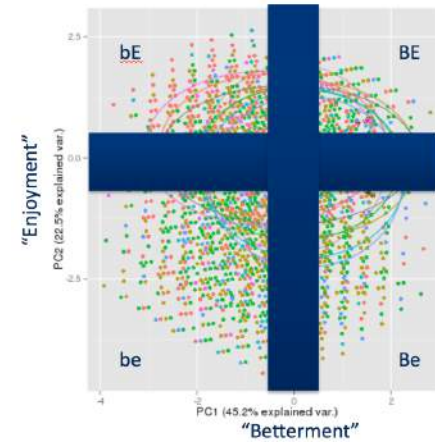
      sequence support
1      <{class/index}> 0.9950678
2      <{forum/index}> 0.6874229
3      <{forum/list?forum_id=2}> 0.7774353
4      <{lecture/39}> 0.5425401
5      <{lecture/index}> 0.9543773
6      <{lecture/view?lecture_id=3}> 0.7953144
7      <{lecture/view?lecture_id=37}> 0.5012330
8      <{lecture/view?lecture_id=39}> 0.6337855
9      <{lecture/view?lecture_id=47}> 0.5400740
10     <{lecture/view?lecture_id=5}> 0.6972873
11     <{lecture/view?lecture_id=7}> 0.5838471
12     <{lecture/view?lecture_id=9}> 0.5246609
13     <{quiz/attempt?quiz_id=49}> 0.5505549
14     <{quiz/attempt?quiz_id=51}> 0.6602959
15     <{quiz/attempt?quiz_id=55}> 0.6732429
16     <{quiz/attempt?quiz_id=7}> 0.9938348
17     <{quiz/index?quiz_type=survey}> 0.8378545
18     <{quiz/start?quiz_id=49}> 0.5856967
19     <{wiki/view?page=AnsweringQuestions}> 0.6171393
20     <{wiki/view?page=courselogistics}> 0.5110974
21     <{wiki/view?page=DigitalLabCoat}> 0.7435265
22     <{wiki/view?page=mTuner}> 0.6288533
23     <{wiki/view?page=pScholar}> 0.5579531
24     <{wiki/view?page=ToDoThisWeek}> 0.9654747
25     <{class/index}, {wiki/view?page=ToDoThisWeek}> 0.5086313
26     <{forum/list?forum_id=2}, {wiki/view?page=ToDoThisWeek}> 0.6183724
27     <{lecture/view?lecture_id=3}, {wiki/view?page=ToDoThisWeek}> 0.5955610
28     <{quiz/attempt?quiz_id=7}, {wiki/view?page=ToDoThisWeek}> 0.8249075
29     <{quiz/index?quiz_type=survey}, {wiki/view?page=ToDoThisWeek}> 0.6048089
30     <{wiki/view?page=ToDoThisWeek}, {wiki/view?page=ToDoThisWeek}> 0.5209618
31     <{quiz/attempt?quiz_id=7}, {wiki/view?page=mTuner}> 0.5425401
32     <{forum/list?forum_id=2}, {wiki/view?page=DigitalLabCoat}> 0.5104809
33     <{lecture/view?lecture_id=3}, {wiki/view?page=DigitalLabCoat}> 0.5030826
34     <{quiz/attempt?quiz_id=7}, {wiki/view?page=DigitalLabCoat}> 0.6578298
35     <{quiz/attempt?quiz_id=7}, {wiki/view?page=AnsweringQuestions}> 0.5314427
```

...we extract frequently occurring sequences

Early Results

We compare different categories, similar to approach taken by Sabourin, Mott & Lester, 2013

	Sequence		Partition		Spread
	BE	Be	be	bE	
<{class/index},<{class/index}>	51%	53%	54%	52%	3%
<{class/index},<{lecture/index}>	54%	55%	57%	57%	3%
<{class/index},<{wiki/view?page=ToDoThisWeek}>	51%		53%	55%	4%
<{forum/list?forum_id=2},<{class/index}>	60%	52%	51%	53%	9%
<{forum/list?forum_id=2},<{lecture/index}>	66%	57%	56%	55%	11%
<{forum/list?forum_id=2},<{quiz/index?quiz_type=survey}>	54%				0%
<{forum/list?forum_id=2},<{wiki/view?page=DigitalLabCoat}>	51%				0%
<{forum/list?forum_id=2},<{wiki/view?page=ToDoThisWeek}>	62%	55%	52%		10%
<{lecture/index},<{class/index}>			51%		0%
<{lecture/index},<{lecture/index}>			52%		0%
<{lecture/view?lecture_id=3},<{class/index}>	59%	56%	59%		4%
<{lecture/view?lecture_id=3},<{lecture/index}>	60%	58%	61%		3%
<{lecture/view?lecture_id=3},<{lecture/view?lecture_id=5}>			54%		0%
<{lecture/view?lecture_id=3},<{lecture/view?lecture_id=5}>	62%	61%	65%	59%	6%
<{lecture/view?lecture_id=3},<{lecture/view?lecture_id=7}>	53%	53%	57%	52%	5%
<{lecture/view?lecture_id=3},<{lecture/view?lecture_id=9}>			51%		0%
<{lecture/view?lecture_id=3},<{wiki/view?page=DigitalLabCoat}>	50%		50%		0%
<{lecture/view?lecture_id=3},<{wiki/view?page=ToDoThisWeek}>	60%	56%	58%	52%	7%
<{lecture/view?lecture_id=5},<{class/index}>			52%		0%
<{lecture/view?lecture_id=5},<{lecture/index}>			53%		0%
<{lecture/view?lecture_id=5},<{lecture/view?lecture_id=7}>	54%	51%	57%	52%	6%
<{lecture/view?lecture_id=5},<{lecture/view?lecture_id=9}>			51%		0%
<{lecture/view?lecture_id=5},<{wiki/view?page=ToDoThisWeek}>	50%		50%		0%
<{lecture/view?lecture_id=7},<{lecture/view?lecture_id=9}>			51%		0%
<{quiz/attempt?quiz_id=51},<{lecture/index}>			50%		0%
<{quiz/attempt?quiz_id=55},<{lecture/index}>			50%		0%
<{quiz/attempt?quiz_id=7},<{class/index}>	79%	78%	80%	75%	5%
<{quiz/attempt?quiz_id=7},<{forum/index}>	83%	84%	86%	83%	4%
<{quiz/attempt?quiz_id=7},<{lecture/index}>	83%	84%	86%	83%	4%
<{quiz/attempt?quiz_id=7},<{lecture/view?lecture_id=3}>	58%	60%	61%	55%	6%
<{quiz/attempt?quiz_id=7},<{lecture/view?lecture_id=39}>	55%	53%	56%		3%
<{quiz/attempt?quiz_id=7},<{lecture/view?lecture_id=5}>	55%	57%	58%	52%	6%
<{quiz/attempt?quiz_id=7},<{lecture/view?lecture_id=7}>			52%		0%
<{quiz/attempt?quiz_id=7},<{quiz/attempt?quiz_id=49}>			52%		0%
<{quiz/attempt?quiz_id=7},<{quiz/attempt?quiz_id=51}>	60%	61%	62%	58%	4%
<{quiz/attempt?quiz_id=7},<{quiz/attempt?quiz_id=55}>	61%	61%	63%	57%	6%
<{quiz/attempt?quiz_id=7},<{quiz/index?quiz_type=survey}>	57%	60%	61%	58%	3%
<{quiz/attempt?quiz_id=7},<{quiz/start?quiz_id=49}>	52%		53%		2%
<{quiz/attempt?quiz_id=7},<{wiki/view?page=AnsweringQuestions}>	53%	52%	55%		3%
<{quiz/attempt?quiz_id=7},<{wiki/view?page=DigitalLabCoat}>	66%	63%	65%	59%	7%
<{quiz/attempt?quiz_id=7},<{wiki/view?page=mTuner}>	54%	54%	55%		2%
<{quiz/attempt?quiz_id=7},<{wiki/view?page=ToDoThisWeek},<{lecture/index}>	51%		54%	50%	4%
<{quiz/attempt?quiz_id=7},<{wiki/view?page=ToDoThisWeek}>	82%	83%	84%	79%	5%
<{quiz/index?quiz_type=survey},<{class/index}>	58%	60%	60%	54%	6%
<{quiz/index?quiz_type=survey},<{lecture/index}>	62%	64%	65%	61%	4%
<{quiz/index?quiz_type=survey},<{quiz/attempt?quiz_id=7}>			52%		0%
<{quiz/index?quiz_type=survey},<{wiki/view?page=ToDoThisWeek}>	60%	59%	61%	57%	4%
<{quiz/start?quiz_id=49},<{quiz/attempt?quiz_id=49}>			52%		0%
<{wiki/view?page=ToDoThisWeek},<{class/index}>	54%	55%	57%	53%	4%
<{wiki/view?page=ToDoThisWeek},<{lecture/index}>	59%	57%	62%	60%	5%
<{wiki/view?page=ToDoThisWeek},<{wiki/view?page=ToDoThisWeek}>	52%	51%	54%	53%	3%



Comparing the prevalence of frequently occurring sequences across the four partitions (BE, Be, bE, and be), we see that “be” shows more “return to outline” behaviours and less engagement with forums.

Number of Cases in Each Partition

Number of Cases in Each Partition

Outcome	Intent		Aboriginal Education	Intro Psych	Intro Stats	Mental Health 1	Mental Health 2	Programming 2
	Betterment	Enjoyment						
Certificate	High	High	283	332	227	290	182	634
No Certificate	High	High	27	1400	1458	54	881	861
Certificate	High	Low	234	234	642	317	206	624
No Certificate	High	Low	25	951	3821	38	1058	1221
Certificate	Low	High	328	973	155	162	120	353
No Certificate	Low	High	36	3702	1008	47	990	441
Certificate	Low	Low	215	533	642	106	121	438
No Certificate	Low	Low	49	2350	4777	44	906	736

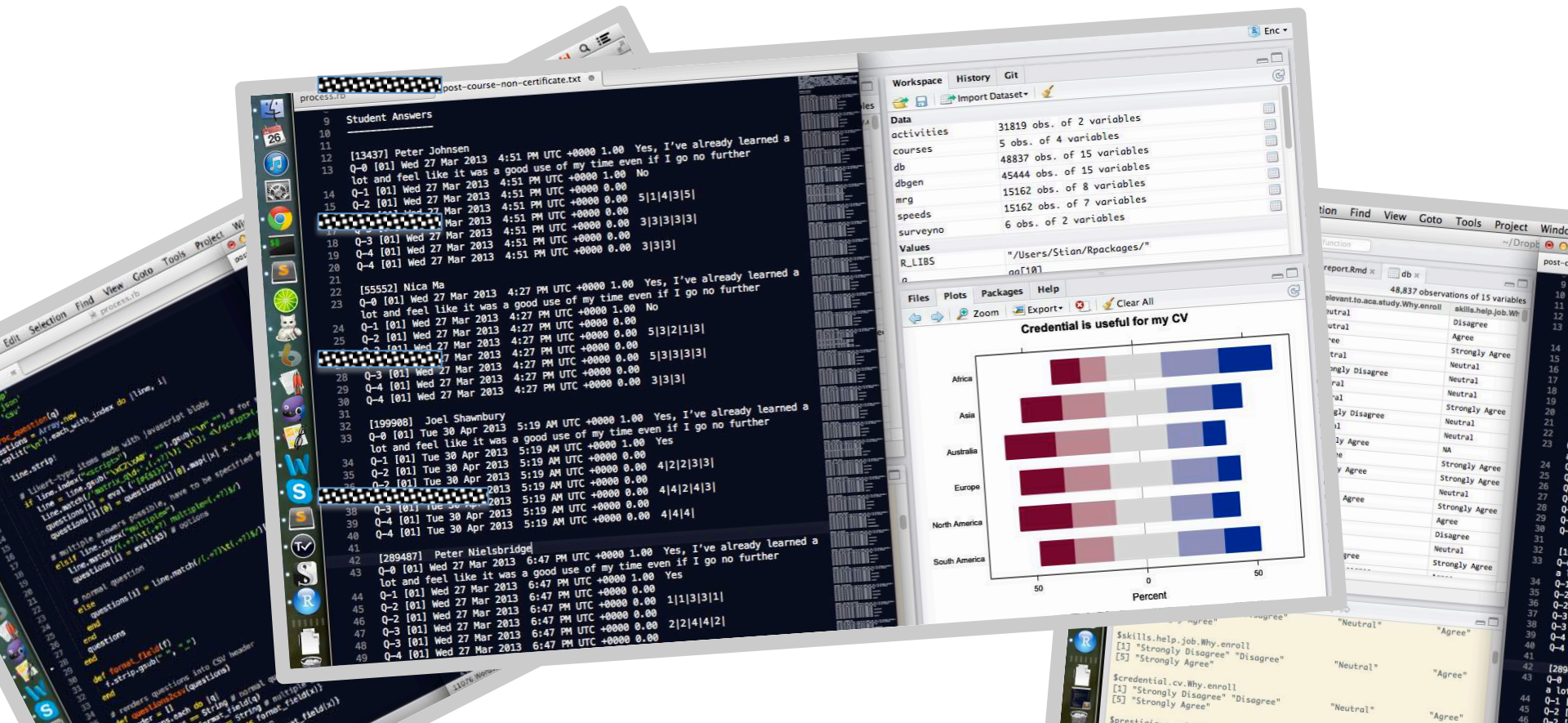
Most Frequent Sequences

- Certificate
 - Common: {wiki/view}{lecture/view}{short-event}{class/index}
 - Exclusive to Low Betterment, Low Enjoyment: {lecture/view, out-of-sequence}{lecture/index}
- No-Certificate
 - Common: {class/index}{lecture/index}{quiz/attempt}{lecture/view}
 - Exclusive to High Betterment, High Enjoyment: {wiki/view}

Data Wrangling and Analysis: A Look Behind the Scenes

Stian Haklev

Open UToronto Institutional Researcher

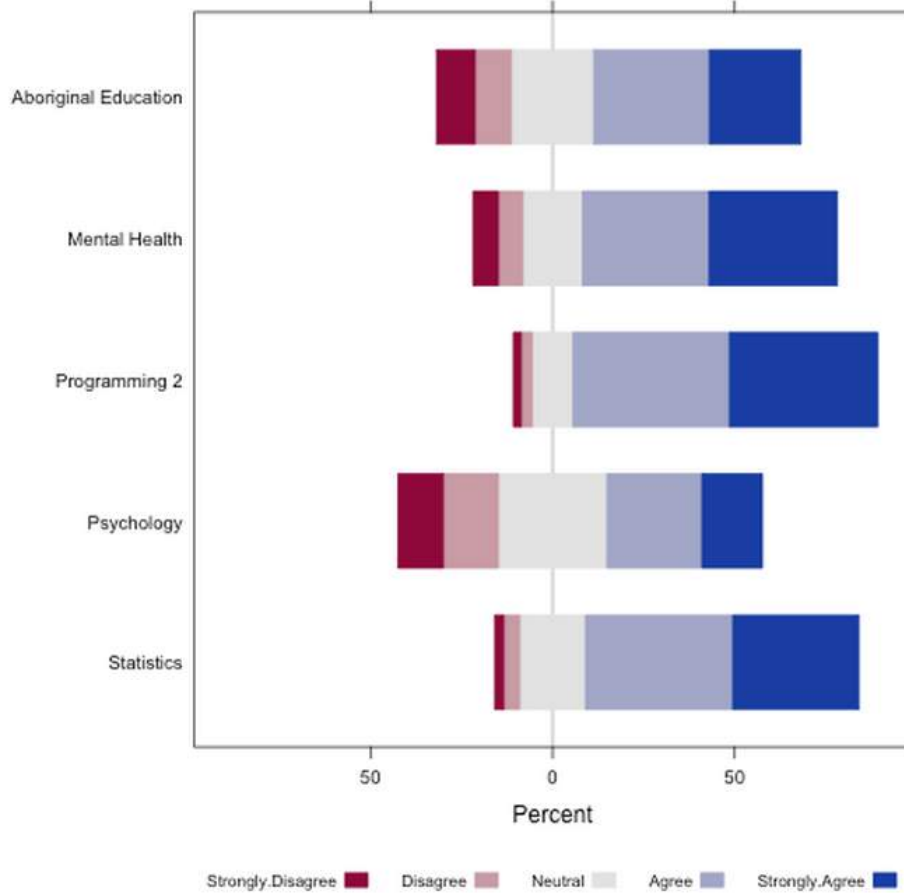


Challenges & What We Are Learning

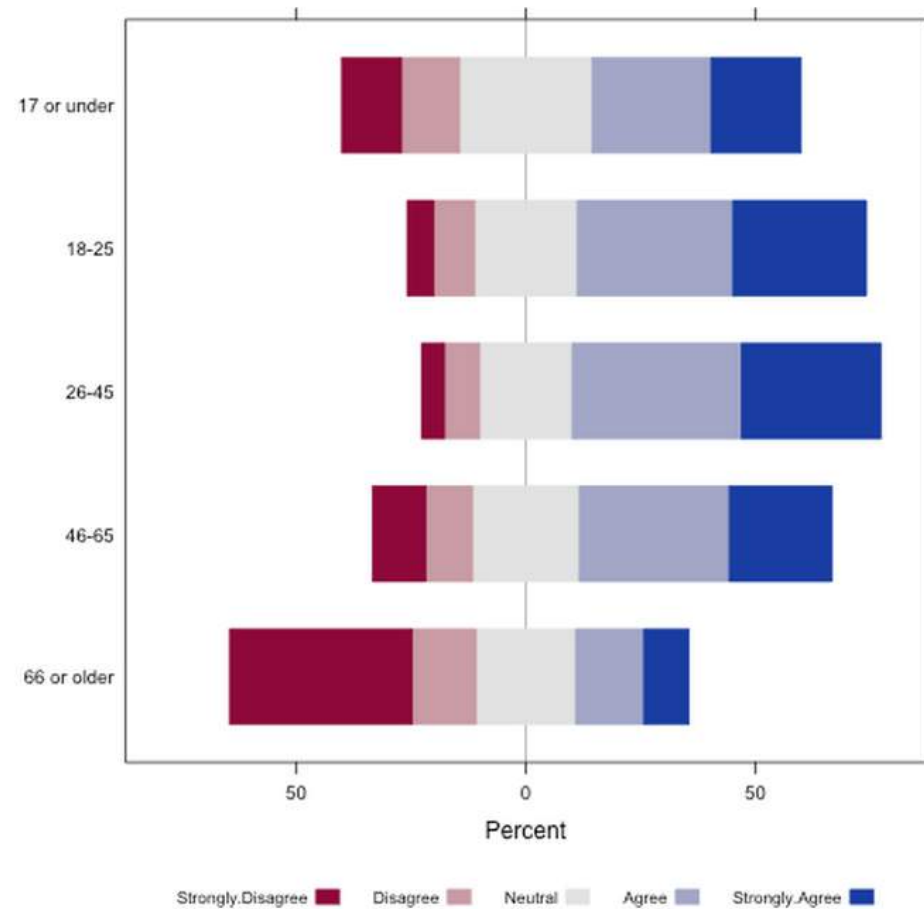
1. Preparing and combining complex data structures for analysis
2. Managing secure and anonymised access to multiple research teams/researchers/ethics agreements
3. Securing enough computing capacity for “big data” analysis
4. Conducting collaborative research in a documented, reproducible way

Simple relationships from survey data

Skills gained will help my job



Skills gained will help my job





Demographic Reports

Demographic Report on University of Toronto Coursera MOOCs

June 16, 2013

Background

In order to support MOOC-related research and evaluation processes, the Office of Online Learning Strategies is involved in administrating access to detailed user data sets from Coursera and cleaning it for analysis. The data is being stored centrally and made available for a range of purposes, including institutional planning processes. It will also be used in a research study among five MOOCs that launched between January and May 2013. The course instructors are all co-PIs on the ethics proposal, and will share the analytics data. For a full description of activities visit: [Open Utoronto MOOC Research and Evaluation](#).

Connecting database with other data

(MySQL 5.5.31-0+wheezy1) mri/sta220/quiz_metadata

sta220 Select Database

Structure Content Relations Triggers Table Info Query

SSH Connected Table History Users Console

Filter

Search: id =

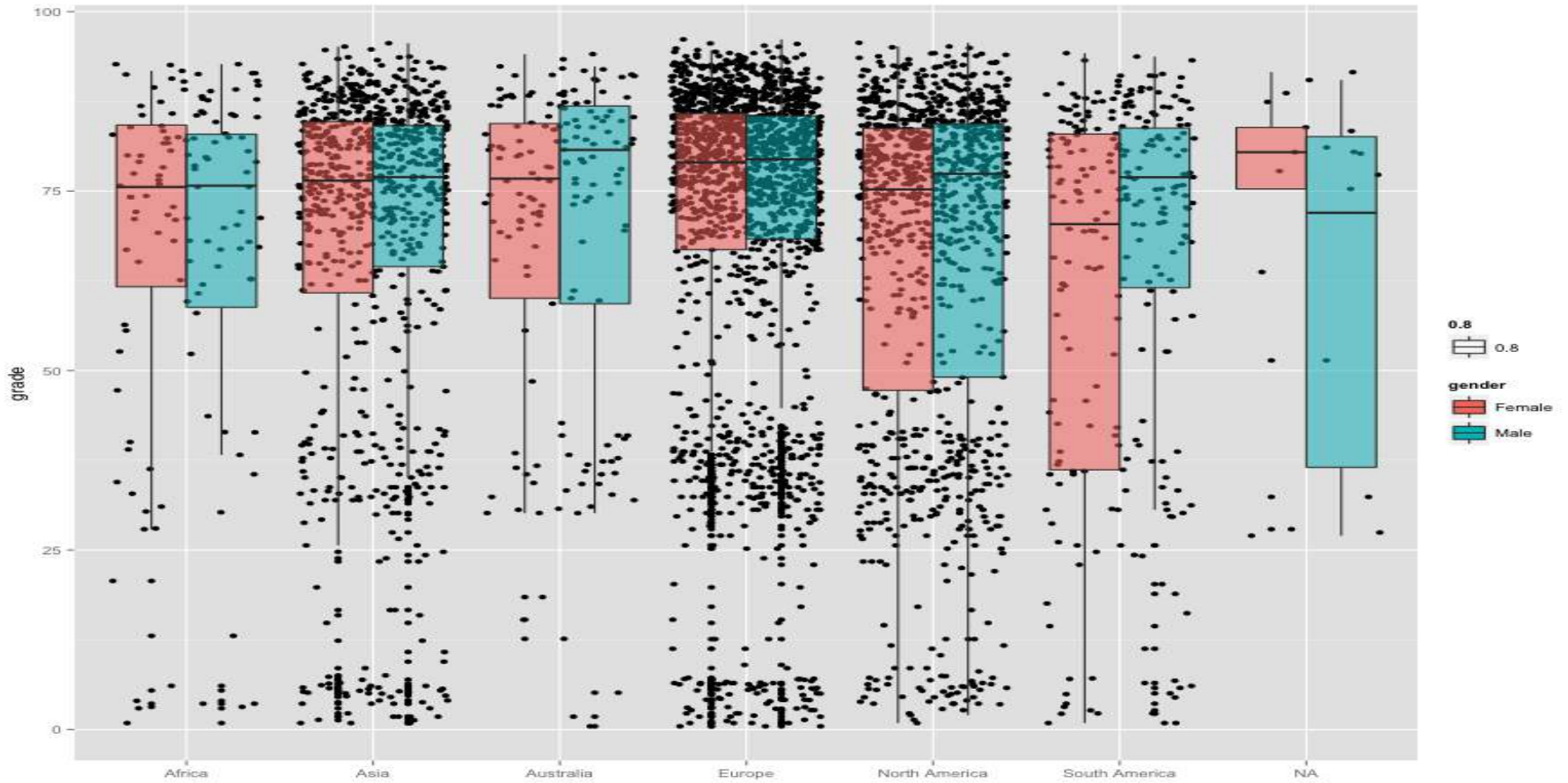
id	parent_id	open_time	soft_close_time	hard_close_time	maximum_submissions	title	duration	quiz_type	proctoring_requirement
1	-1	0	2147483640	2147483640	0	Sample Lecture	0	video	none
2	-1	NULL	2147483647	2147483647	100	Tutorial Quiz	0	quiz	none
3	-1	1341338040	1893484860	1893484860	1	Post Course Survey Template	0	quiz	none
4	-1	1341338040	1893484860	1893484860	1	Original Feedback Survey (Cert Ear...	0	survey	none
5	-1	1345060237	1345665037	1346874637	100	New Quiz	0	quiz	none
6	-1	NULL	1893484860	1893484860	1	Post Course Survey (Non-certificat...	0	survey	none
7	-1	NULL	1372651140	1372651140	1	Pre-course survey	0	survey	none
8	-1	NULL	1893484860	1893484860	1	Post Course Survey (Cert earners)	0	survey	none
9	-1	1341338040	1893484860	1893484860	1	Post Course Survey (Cert Earners)	0	survey	none
10	1	0	2147483640	2147483640	0	Sample Lecture	0	video	none
11	2	NULL	2147483640	2147483640	100	Tutorial Quiz	0	quiz	none
12	6	NULL	1893484860	1893484860	1	Post Course Survey (Non-certificat...	0	survey	none
13	7	0	1372651140	1372651140	1	Pre-course survey	0	survey	none
14	8	NULL	1893484860	1893484860	1	Post Course Survey (Cert earners)	0	survey	none
15	-1	1364788860	1363973340	1365182940	100	1 Introduction (6:01)	0	video	none
16	15	1364788860	1363973340	1365182940	100	1 Introduction (6:01)	0	video	none
17	-1	1364788861	1364240880	1365450480	100	2 Five Number Summary (8:44)	0	video	none
18	17	1364788861	1364240880	1365450480	100	2 Five Number Summary (8:44)	0	video	none
19	-1	1364788861	1364270445	1365480045	100	R tutorial for 2 Five Number Sum...	0	video	none
20	19	1364788861	1364270445	1365480045	100	R tutorial for 2 Five Number Sum...	0	video	none
21	-1	1364788862	1364272431	1365482031	100	R tutorial for 3 The Centre of The...	0	video	none
22	21	1364788862	1364272431	1365482031	100	R tutorial for 3 The Centre of The...	0	video	none
23	-1	1364788861	1364273360	1365482960	100	R tutorial for 4 The Spread of The...	0	video	none
24	23	1364788861	1364273360	1365482960	100	R tutorial for 4 The Spread of The...	0	video	none
25	-1	1364788861	1364274432	1365484032	100	R tutorial for 5 The Shape of The...	0	video	none
26	25	1364788861	1364274432	1365484032	100	R tutorial for 5 The Shape of The...	0	video	none
27	-1	1364788861	1364275815	1365485415	100	R tutorial for 6 Categorical Variabl...	0	video	none
28	27	1364788861	1364275815	1365485415	100	R tutorial for 6 Categorical Variabl...	0	video	none
29	-1	1364788861	1364320883	1365530483	100	Installing R - Mac OSX (2:10)	0	video	none
30	29	1364788861	1364320883	1365530483	100	Installing R - Mac OSX (2:10)	0	video	none
31	-1	1364788861	1364320995	1365530595	100	Installing R - PC (1:55)	0	video	none
32	31	1364788861	1364320995	1365530595	100	Installing R - PC (1:55)	0	video	none
33	-1	NULL	1364408460	1365618060	100	Week 1 Quiz	0	quiz	none
34	33	NULL	1364408460	1365618060	100	Week 1 Quiz	0	quiz	none
35	-1	1363806565	1364411365	1365620965	100	NULL	0	video	none

TABLE INFORMATION

- created: 18/02/14
- engine: InnoDB
- rows: 220
- size: 48.0 KiB
- encoding: utf8mb4
- auto_increment: 221

Rows 1 - 100 of 220 from table

Integrated surveys allow us to link intention and demographics with outcomes and behaviour



Analyzing the clicklog

Loading the files

```
In [4]: import pandas as pd
store = pd.HDFStore("/Users/Stian/src/clickpy/mentalhealth_002.h5")

In [5]: store2 = pd.HDFStore("/tmp/mentalhealth_002.h5", complib="zlib", complevel=9)

In [6]: db=store['db']
store2['db'] = db
store2.close()

In [65]: lec5 = db[db.lecture_id == 5]

In [66]: lec5['timestamp'] = pd.to_datetime(lec5.timestamp, un

In [67]: lec5.set_index(pd.DatetimeIndex(lec5.timestamp), inplace=True)
```

Categorizing video views

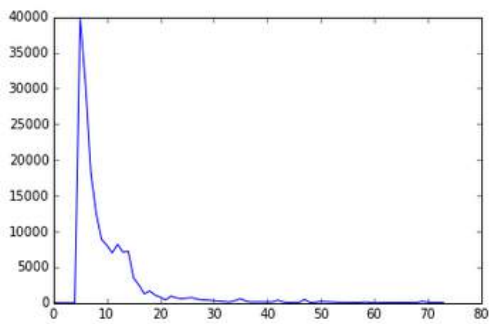
```
In [68]: lecnorm = lec5[lec5.playbackRate == 1]
lecfast = lec5[lec5.playbackRate > 1]
norm=lecnorm.action.resample("D", how="count")
fast=lecfast.action.resample("D", how="count")
tot=lec5.action.resample("D", how="count")
```

Video views

This is how we find everyone watching the video by day

```
In [69]: plot(lec5.action.resample("1D", how="count"))

Out[69]: [<matplotlib.lines.Line2D at 0x11f366490>]
```



Tools to collaborate and document



Pages / Hatch, Match and Dispatch

Coherence of Research Question and Analysis

Added by Hedieh Najafi, last edited by Hedieh Najafi on Jan 09, 2014 (view change)

Hedieh's note: Here, I am first creating a map that shows the linkage between research objectives, res that we would use.

Red font color means questions or things that I will have to double check with team members or my no

- Research goals from the proposal:
- Research Questions:
- Data Sources for Each Question and Data Analysis
 - Research Question 1:
 - Research Question 2:
 - Research Question 3:
 - Research Question 4:

The screenshot shows the RStudio environment with several windows:

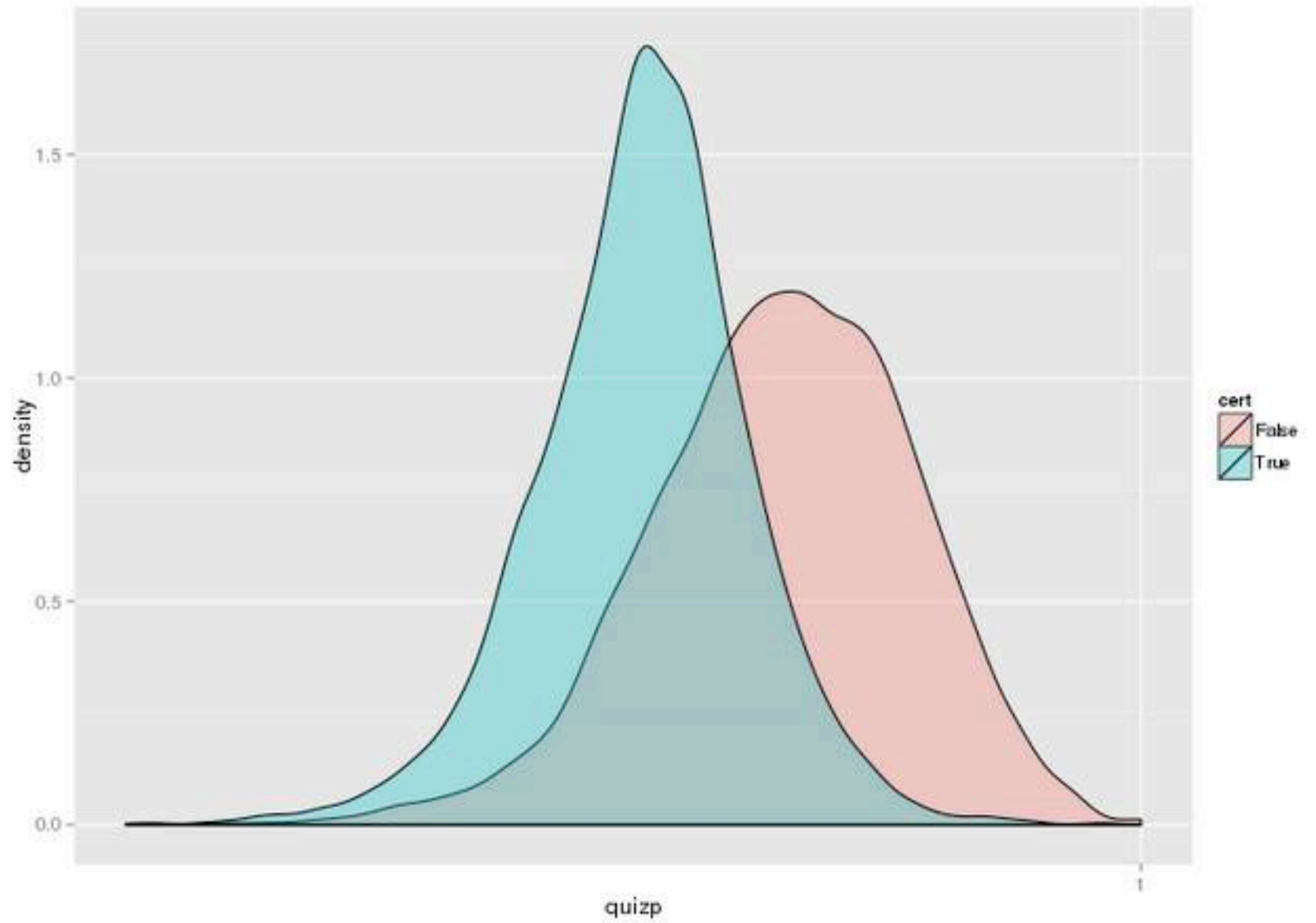
- Environment:** Shows the current workspace with variables like 'ave (stats)' and 'R Documentation'.
- Console:** Displays R code execution output, including a warning message: "Warning messages: 1: In if (columns == "ALL") { : the condition has length > 1 and only the first element will be used".
- Plot Window:** Displays a horizontal bar chart titled "Course was at a prestigious university". The y-axis lists regions: Africa, Asia, Australia, Europe, North America, and South America. The x-axis is labeled "Percent" and ranges from 50 to 50 (representing 0% to 100%). The bars are segmented by color: red (Strongly Disagree), pink (Disagree), grey (Neutral), light blue (Agree), and dark blue (Strongly Agree).

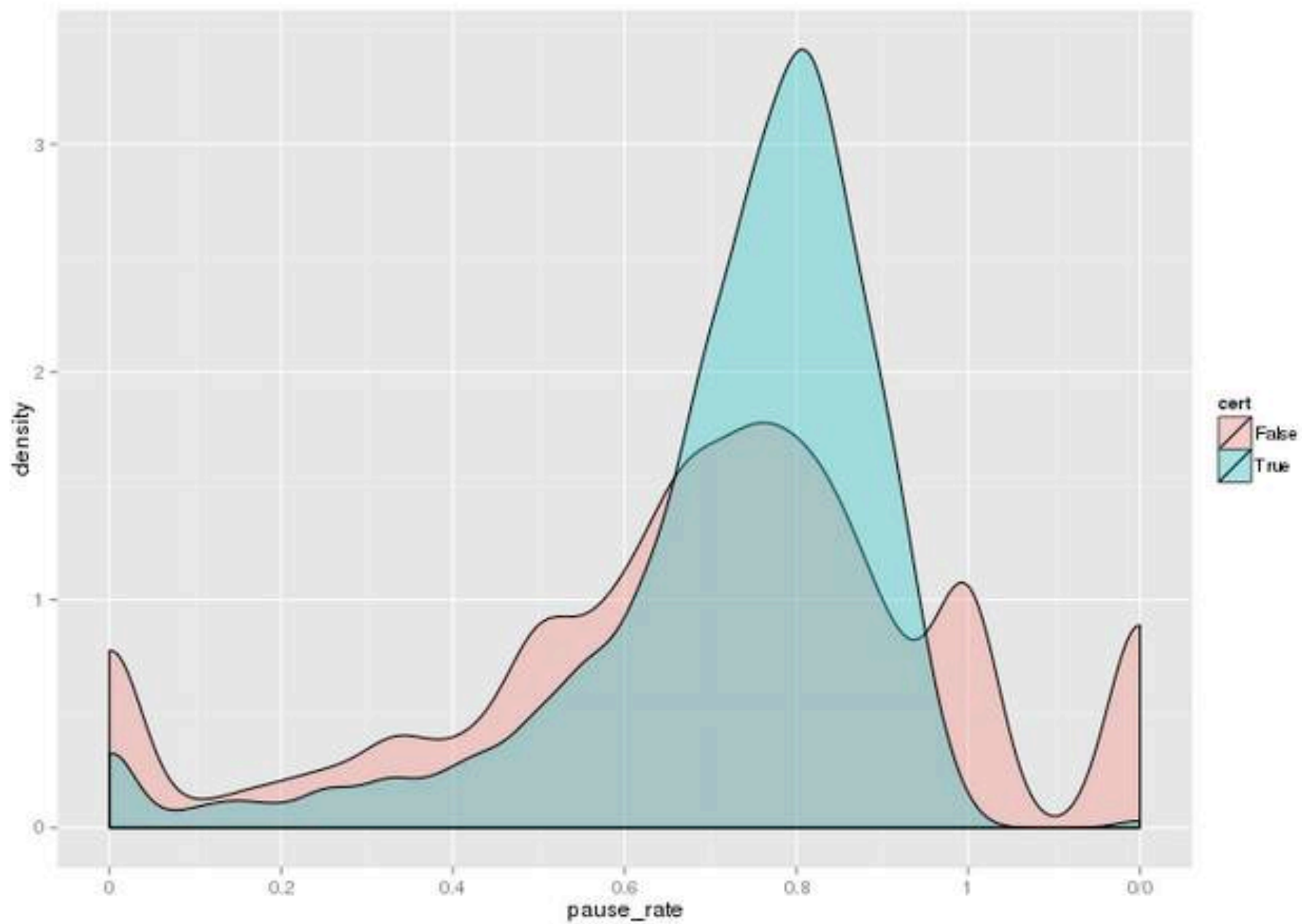
Clicklog: big data, making it queryable, increasing levels of abstraction

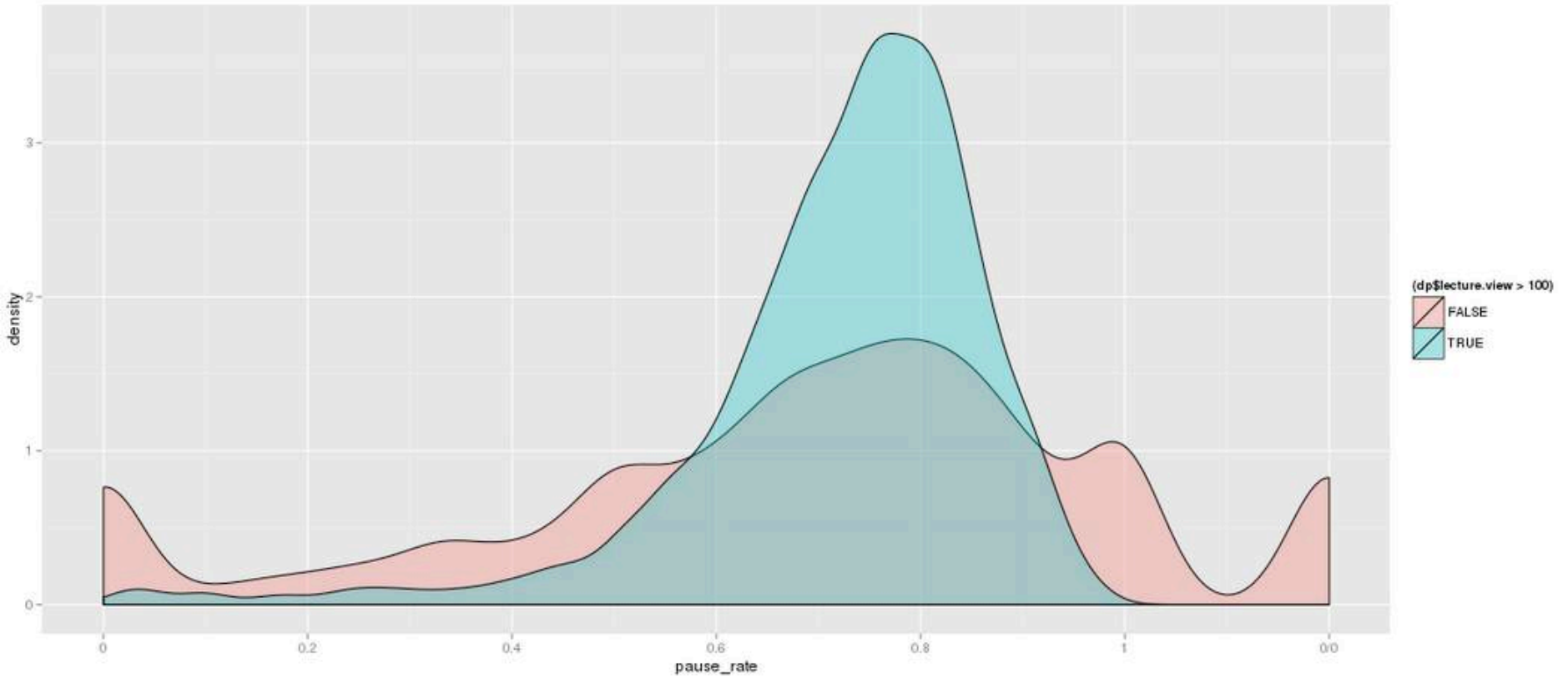
action	duration	action_val	lecture_id_val	type_val
0	16118	survey/attempt	NaN	NaN
1	71.000	lecture/view	out-of-sequence, sought, pause	out-of-sequence
1	31.000	lecture/view	out-of-sequence	out-of-sequence
1	33.000	lecture/view	out-of-sequence, pause, short-event	out-of-sequence
1	31.000	lecture/view	seen-before, sought, short-event	seen-before, sought, short-event
1	33.000	lecture/view	seen-before, short-event	seen-before, short-event
1	31.000	lecture/view	seen-before, sought, pause	seen-before, sought, pause
1	33.000	lecture/view	seen-before, pause, short-event	seen-before, pause, short-event
1	35.000	lecture/view	out-of-sequence, pause, short-event	out-of-sequence, pause, short-event
1	33.000	lecture/view	seen-before, sought, short-event	seen-before, sought, short-event
1	35.000	lecture/view	seen-before, short-event	seen-before, short-event
1	33.000	lecture/view	seen-before, sought, pause	seen-before, sought, pause

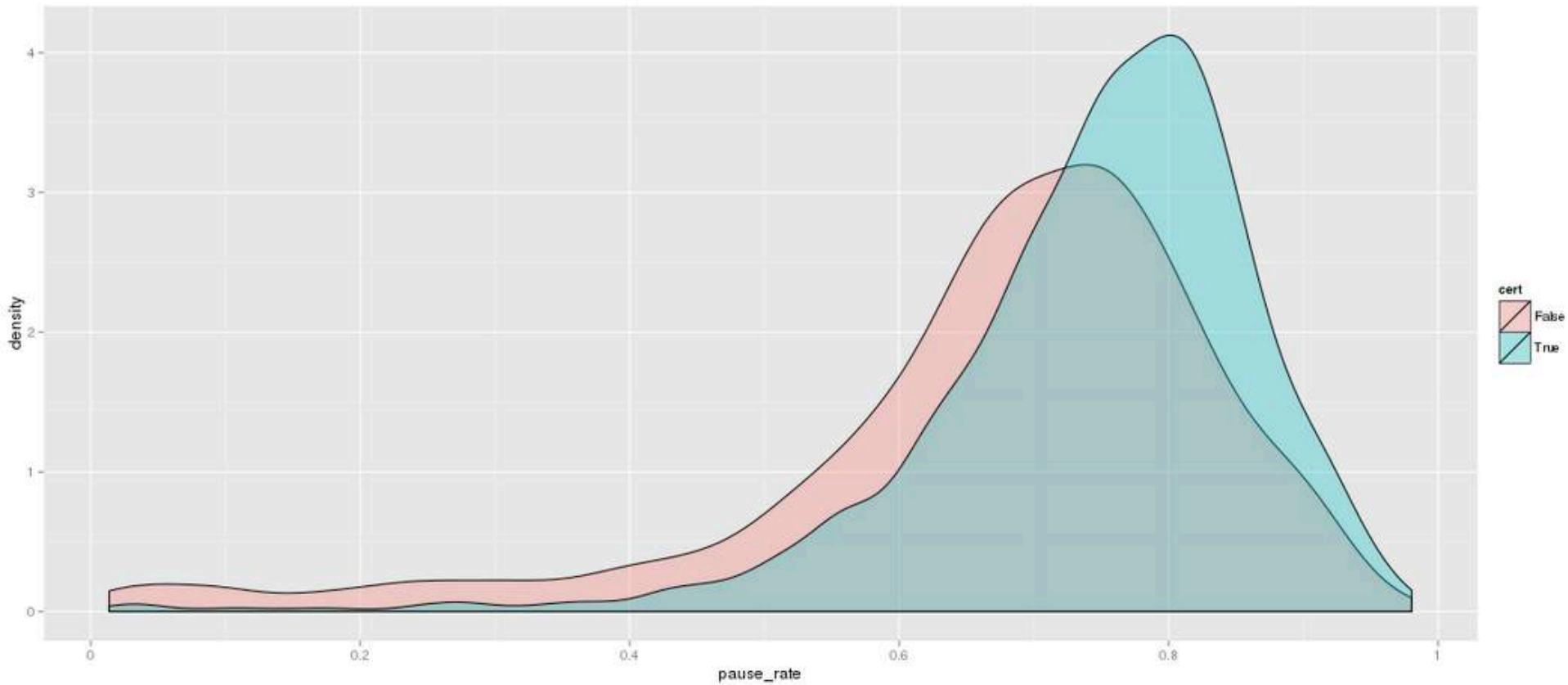
[1 rows x 4 columns]

duration	action_val	lecture_id_val	type_val
169202	human_grading/	NaN	NaN

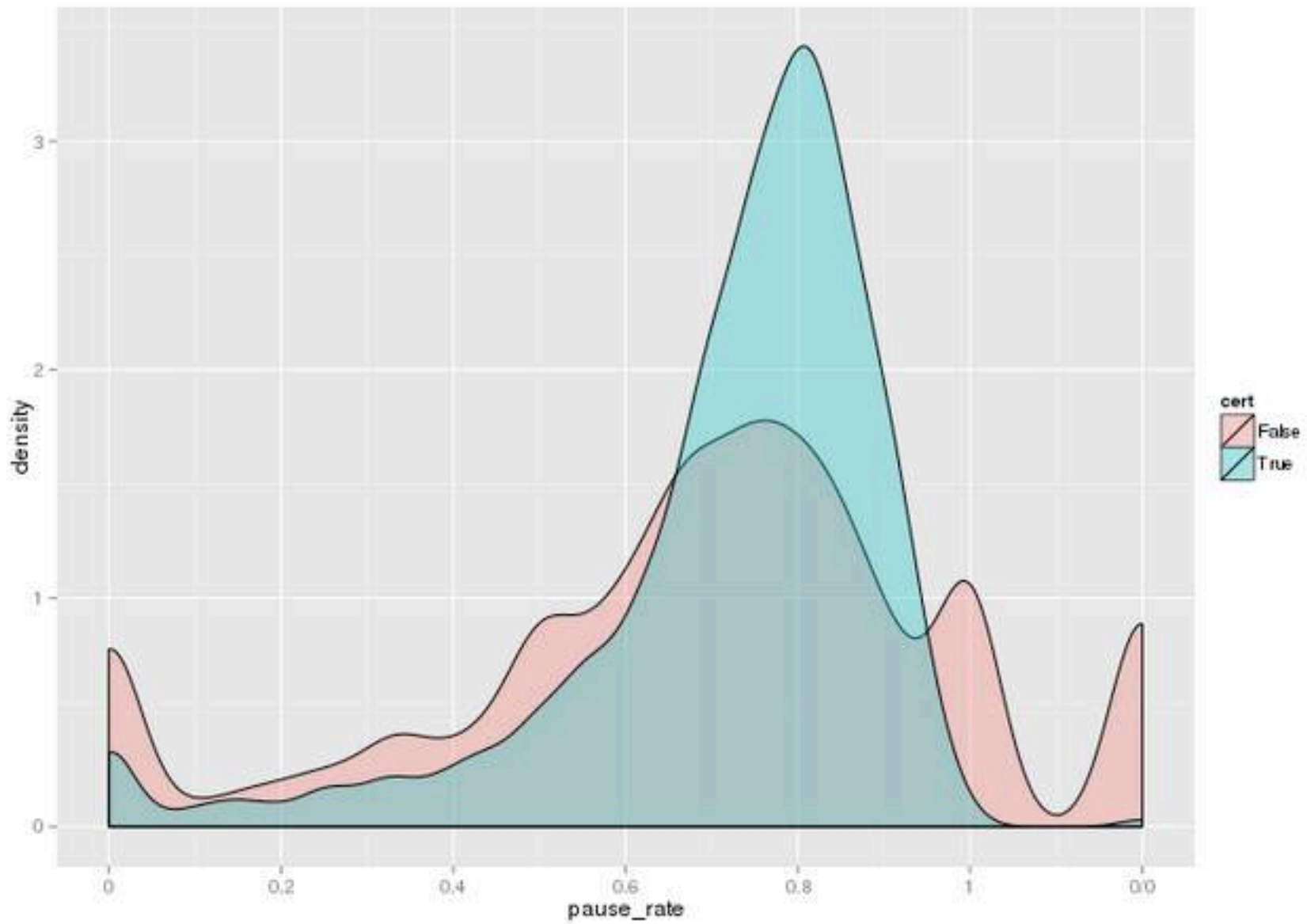


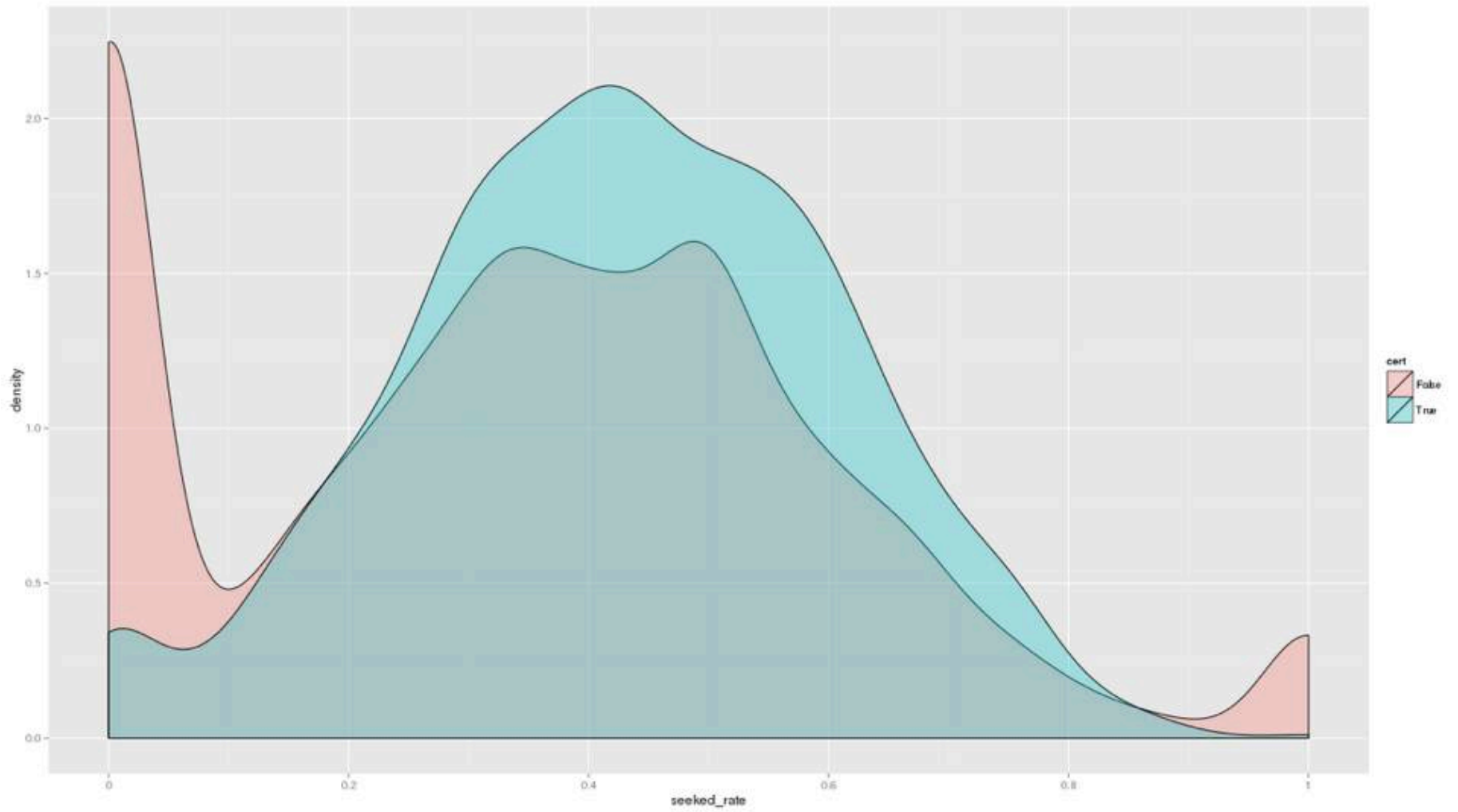


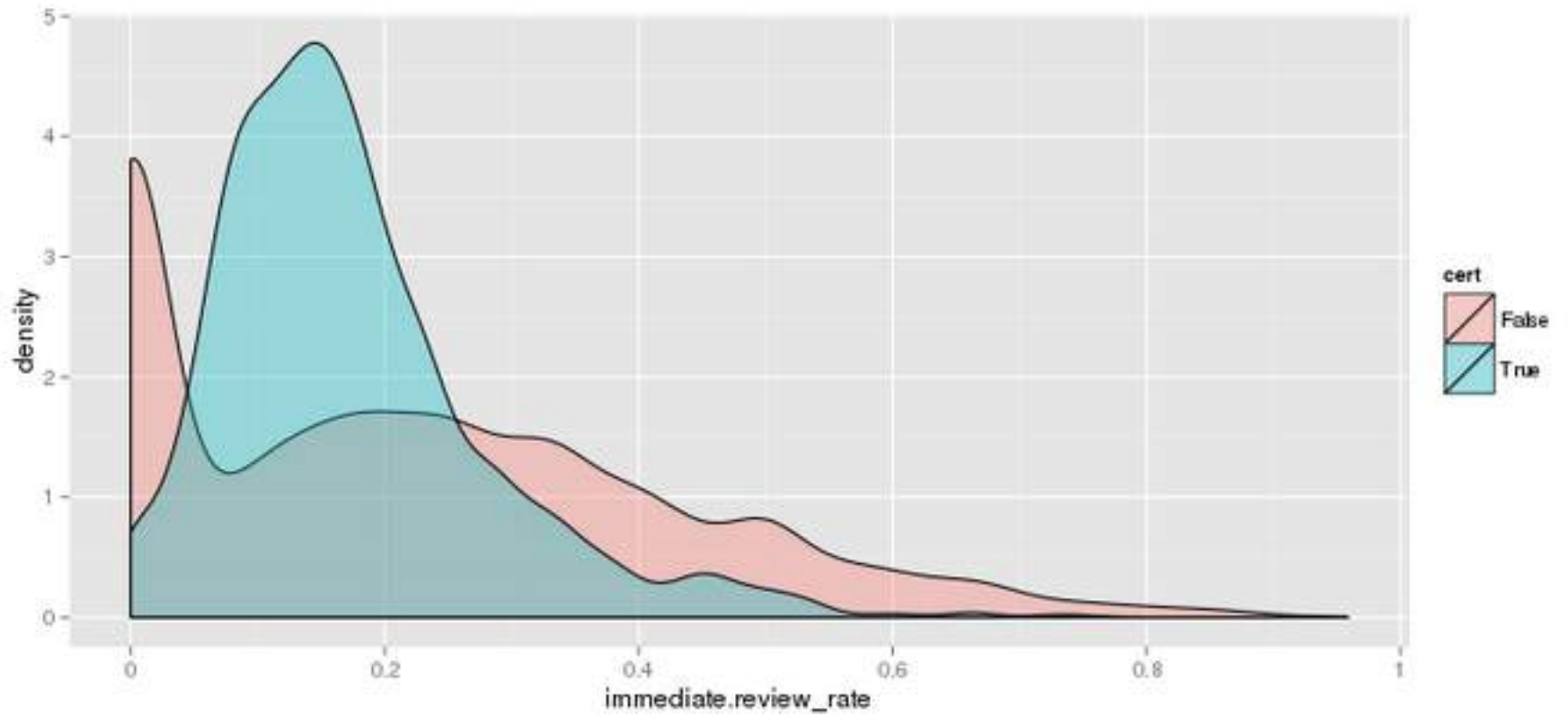




Students with $100 > \text{lecture.view} > 200$







sequence	hi-hi-cert-sorted	hi-hi-nocert-sorted	hi-hi-sorted	hi-lo-cert-sorted	hi-lo-nocert-sorted	hi-lo-sorted	lo-hi-cert-sorted	lo-hi-nocert-sorted	lo-hi-sorted	lo-lo-cert-sorted
<{class/index},{lecture/index}>	0.590655885	0.583867422	0.5859498797	0.5816372682	0.5838588957	0.5714839961	0.5619515097	0.8184888037	0.5972074443	0.5653710247
<{class/index},{wiki/view}>	0.5471698113	0.5565626634	0.5534949234	0.5024875622	0.5394021739	0.5235693501	0.4948683623	0.5313136151	0.5175934741	0.5041224971
<{class/index},{lecture/view}>	0.5399820305	0.5217012376	0.5276718117	0.5142469471	0.4575407609	0.4818622696	0.4983886162	0.5223326548	0.5133001064	0.5197290931
<{lecture/index},{lecture/view}>	0.4941599281	0.4734181628	0.4801924996	0.4825870647	0.4385190217	0.4574199806	0.4700282612	0.4850317327	0.4793825017	0.4787985866
<{class/index},{lecture/view,pause}>	0.4858939802	0.4295799198	0.4479722988	0.4595205789		0.4019398642	0.4434032426	0.4430906478	0.4432062123	0.4670200236
<{lecture/view},{lecture/view}>	0.4733153639	0.4191214921	0.4368214097	0.4522840344		0.4065955383	0.4361148297	0.4161477667	0.423662056	0.472614841
<{class/index},{lecture/index},{lecture/view}>	0.4567834681	0.4332403695	0.440929632	0.4364540932		0.409699321	0.4193564381	0.44713208	0.4366728267	0.4331566549
<{lecture/index},{lecture/view,pause}>	0.4486972147	0.4008192435	0.4164563648	0.4332881049			0.4262977837	0.422045264	0.423662056	0.4361012956
<{class/index},{lecture/view,out-of-sequence}>	0.4379155436	0.4102318285	0.4192734315	0.4265038444			0.4156874411	0.4229732966	0.420208695	0.445229682
<{lecture/view},{lecture/view,pause}>	0.4283917341			0.4147444595						0.4284452297
<{lecture/view,pause},{lecture/view}>	0.4273135669			0.4093170511						0.4260895171
<{class/index},{lecture/index},{lecture/view,pause}>	0.4145552561									
<{class/index},{lecture/index},{pause}>	0.4145552561									
<{lecture/index},{lecture/view,out-of-sequence}>	0.4041329739			0.4066033469						0.4175500589
<{class/index},{lecture/view,seeked}>	0.4023360288									
<{class/index},{lecture/view},{lecture/view}>	0.4016172507									
<{lecture/view,out-of-sequence},{lecture/view}>										0.4093050648
<{lecture/view},{lecture/view,out-of-sequence}>										0.410188457

sequence	lo.lo.cert.sorted	lo.lo.nocert.sorted
<{class/index},{forum/index}>	0.2618	0.18585
<{forum/index},{forum/thread}>	0.2506	0.15308
<{class/index},{forum/thread}>	0.2494	0.18696
<{class/index},{forum/index},{forum/thread}>	0.2288	0.13930
<{forum/thread},{forum/thread}>	0.2264	0.14630
<{forum/index},{forum/thread},{forum/thread}>	0.2011	0.10864
<{class/index},{forum/thread},{forum/thread}>	0.1932	0.12875
<{forum/thread},{forum/thread},{forum/thread}>	0.1829	0.10531
<{class/index},{forum/index},{forum/thread},{forum/thread}>	0.1820	0.09776
<{forum/index},{forum/list}>	0.1699	0.10775
<{forum/index},{forum/index}>	0.1673	0.10076
<{wiki/view},{forum/index},{forum/thread}>	0.1670	0.11264
<{forum/index},{forum/thread},{forum/thread},{forum/thread}>	0.1649	0.08009
<{forum/index},{forum/list},{forum/thread}>	0.1572	0.09453
<{class/index},{forum/thread},{forum/thread},{forum/thread}>	0.1572	0.09309
<{forum/thread},{forum/index}>	0.1555	0.09842
<{class/index},{wiki/view},{forum/index},{forum/thread}>	0.1549	0.10409
<{class/index},{forum/index},{forum/list}>	0.1537	0.10009
<{class/index},{forum/index},{forum/index}>	0.1519	0.09042
<{forum/thread},{forum/thread},{forum/thread}>	0.1493	0.07876
<{forum/index},{forum/index},{forum/thread}>	0.1487	0.08054
<{forum/index},{forum/thread},{forum/index}>	0.1472	0.08354
<{class/index},{forum/index},{forum/list},{forum/thread}>	0.1422	0.08831
<{forum/thread},{forum/index},{forum/thread}>	0.1419	0.08087
<{lecture/view},{forum/index},{forum/thread}>	0.1402	0.08787
<{class/index},{forum/thread},{forum/index}>	0.1390	0.08731
<{forum/index},{forum/thread},{wiki/view}>	0.1381	0.08509
<{lecture/view,pause},{forum/thread}>	0.1378	0.08776

Next Steps

- Examine outcome measures
- Repeat with other MOOCs
- Refine feature identification