

An Analysis of Student Behavior in Two Massive Open Online Courses

James Schaffer*, Brandon Huynh*, John O'Donovan*, Tobias Höllerer*, Yinglong Xia[†] and Sabrina Lin[†]

*Department of Computer Science

University of California Santa Barbara, Santa Barbara, California 93106

Email: {james_schaffer,bhuynh,jod,holl}@cs.ucsb.edu

[†]IBM Thomas J. Watson Research Center

Yorktown Heights, NY 10598

Email: {yxia,sabrinall}@us.ibm.com

Abstract—Massive open online courses (MOOCs) have high potential for improving education worldwide, but understanding of student behavior and situations is difficult to achieve in online settings. Network analytics and visualizations can assist instructors with supporting understanding of student behavior as courses unfold. In this work, we perform a visual comparative analysis of two different MOOC courses to analyze the impacts of course structure differences and demonstrate the benefits of visual network analysis in this context. We present several insights: (1) behavior features that are best for prediction of student attrition varied with course structure, (2) a large proportion (about 35%) of students never received a reply to their original post and this was correlated with an eventual dropout, and (3) students that received a reply to their original post were twice as likely to post again. We contribute several information visualizations of student network data and draw recommendations for MOOC instructors and designers of course systems.

I. INTRODUCTION

MOOCs (massively open online courses) have recently emerged as a community driven alternative for people to access new knowledge¹. An emergent problem in MOOCs is student attrition. Many students (90-95%) stop engaging the course material long before the course has finished and do not obtain a certificate. Thus, MOOCs to fall short of the completion rates for traditional courses (e.g. [1][2][3]).

A common approach to decreasing student attrition is to encourage participation in the online forums that accompany MOOC courses. It is thought that through social interactions, students can overcome obstacles in the course. Research groups have begun to recognize the necessity of understanding forum interactions [4]. Additionally, the MOOC platform NovoED has already taken steps towards giving its instructors a suite of analytics tools to monitor student progress. Despite this, previous experimental attempts at intervention and course modification, e.g. [5] and [6], have not been completely successful, which indicates that MOOC communities are still not understood.

Visualization and network analytics are well positioned to benefit teachers and institutions that are trying to design MOOC

courses and guide student interaction remotely. Teachers could potentially visualize the forum in real time to watch the emergence of superposters [7] in the center of the graph, the development of cliques, and identify unconnected students that could benefit from intervention. Moreover, network analytics [8] are well equipped to provide insight into MOOC forums.

In this work, a visualization and network analytics approach is utilized to highlight the differences between two MOOC courses from Stanford Online Lagunita that vary in structure and curriculum. After an initial exploratory phase, we identified a large percentage (about 35%) of students that had posted in the course forums but never received a reply, indicating a “bad forum experience”. This group of students with “orphaned” posts spurred our investigation. In this paper, we contribute a multidimensional comparative analysis of two MOOC courses, based on the following research questions:

- 1) What features correlate with dropout? Does this change between courses?
- 2) How does the structure of the student interaction network change over time and what are the differences between the two courses?
- 3) To what degree might further forum participation be encouraged if replies are given to a student’s original post?

II. RELATED WORK

The definition of “dropout” in MOOC research is varied between researchers. Recent research has noted that social, motivational, and economic circumstances vary considerably between students. Rivard et al [1] questioned whether looking at “registered vs. completed” students made any sense when trying to determine dropout rate, since MOOCs do not have a substantial impact on college credit, have no prerequisites, and are free to students. Instead, it might make more sense to analyze the different kinds of people that sign up for MOOCs and what their goals are. For instance, Halawa et al [2] noted that MOOC dropouts are very heterogeneous and that internal factors (such as student ability and self regulation) play a large role. Onah et al [3] provided a taxonomy of dropout motivations, many of which could be chalked up to internal factors of students (no real intention to complete, lack of time,

¹<https://www.class-central.com/report/moocs-stats-and-trends-2014/>

IEEE/ACM ASONAM 2016, August 18-21, 2016, San Francisco, CA, USA
U.S. Government work not protected by U.S. copyright

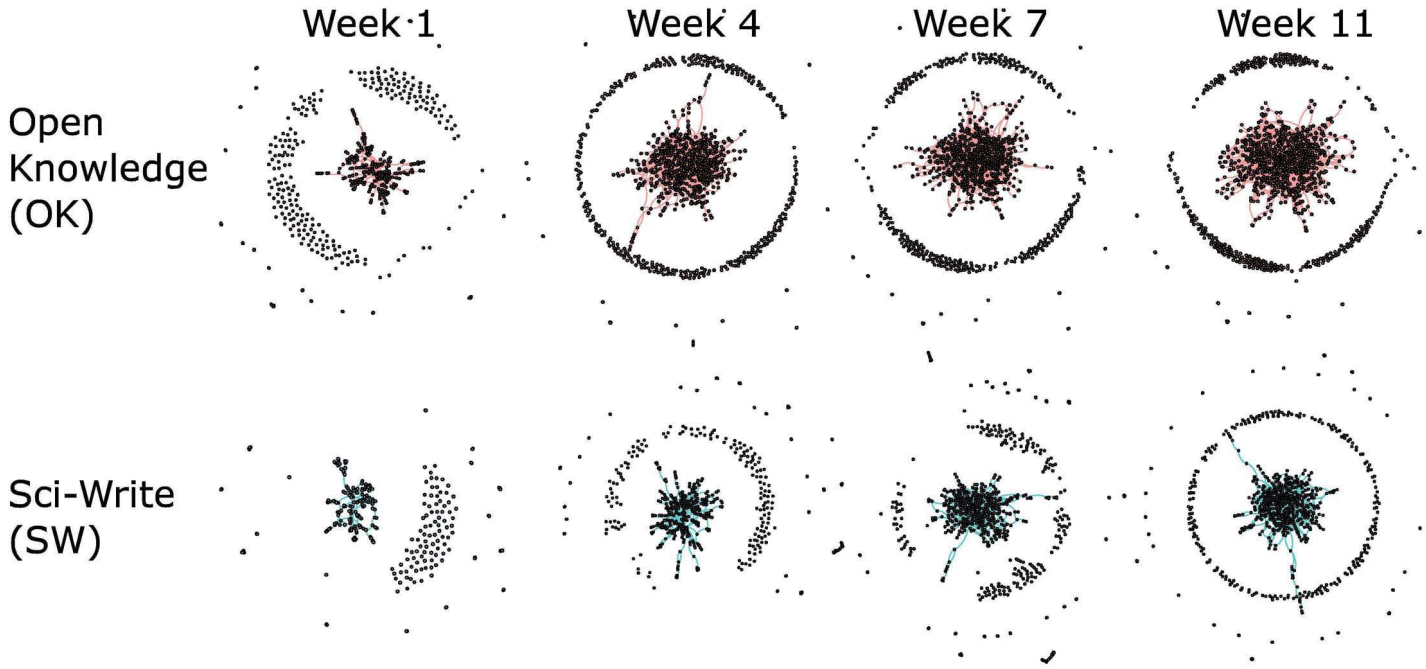


Fig. 1. A Force Atlas layout reveals some key features of the forums throughout the courses. Students with orphaned posts (initial posts by students that received no reply) are shown in the outer circles, and small connected components of students that communicated in the periphery of the forum are shown at the outermost edge. 4 of the 11 weeks of the forum are shown. Keep in mind that the Sci-Write course concluded at week 11 while the Open Knowledge course continued until week 15. The statistics of the final network can be seen in Table II.

lack of digital learning skills, starting late), but others shed some light on how MOOCs might be improved, for example, peer review requirements in a course have been observed to lead to attrition. Finally, Yang et al [9] are optimistic about intervening to help struggling students (those that have stayed in the course longer than one week, but eventually drop out). He writes “supporting the participation of these struggling students may be the first low hanging fruit for increasing the success rate of courses.”

One promising method for helping struggling students is the improvement of online community features. Students have indicated that they are enthusiastic about online forums, chat, and networking, however, the mere presence of such features does not guarantee their usage. Kotturi [10] notes that one of the main problems is limited social translucence online. Coetzee et al [11] indicated that embedded chat features had apparently no effect on a student’s grades, retention, forum participation, or perceived sense of community.

Researchers have also studied how forum interaction and peer bonds help students remain active in the course, for instance, Yang et al [12] has categorized the different ways in which peers form bonds in MOOC forums. He notes that students are much more likely to drop out if they lose close peers. Additionally, students known as “superposters” [7] display high engagement, get better grades, and significantly stimulate forum activity, contributing to its overall health by engaging other students.

We identified two primary gaps in current research on student behavior in MOOCs: 1) differences between course structures are often ignored (this can impact prediction of dropouts)

and 2) intervention efforts are not always successful which means that student use of the forum is not well understood. To address these gaps, we employed network analysis and visualization methods to analyze data from two courses with distinct structures and made several novel findings, in addition to findings which reinforce student behavior patterns found in other work. The tools (Gephi², R³) used in this analysis are openly available to instructors and could be used as new courses unfold to better understand student behavior.

III. MOOC DATASET

Two courses were available to us for analysis: Sci-Write and Open Knowledge, both of which were offered by Stanford Online Lagunita. Both courses were free to the general public, had an online forum where students could discuss the subject matter, occurred over a fixed time period (Fall 2014), and offered students the chance to obtain a certificate of completion if they fulfilled some requirements.

The Open Knowledge course started on August 21, 2014 and concluded on November 28. Students explored broad topics related to MOOCs, including subjects such as digital identity, rapidly changing technology, intellectual property, copyright, and global equity. Students that wanted to get a statement of accomplishment of the course had to complete the first of three different “tracks.” One of the primary requirements of getting the certificate of completion was participation in the forum, and beyond that, social media.

²<https://gephi.org/>

³<https://www.r-project.org/>

The Sci-Write course started on September 2, 2014 and concluded on November 6. The 8 week course covered topics such as principles of effective writing, crafting better sentences, formatting of manuscripts, the scientific peer review process, and other issues in scientific writing such as plagiarism. Unlike the open knowledge course, Sci-Write had concrete assignments to complete, short papers to write, and peer editing exercises.

55318 students registered for the Open Knowledge course. Of these, 125 (0.002%) students received at least some level of the certificate of completion, 5273 (9.5%) watched at least one video, and 1374 (2.4%) ever posted on the forum, for a total of 9917 posts. In the dataset, only 1119 of the students that posted on the forum had a grade recorded. Remember that in the Open Knowledge course, forum participation was required for a statement of completion, which means that students that never posted in the forum were most likely not interested in a certificate of completion. Therefore, the real attrition rate of the course was 91%, with 9% of students that posted in the forum obtaining a certificate of completion. This agrees with completion rates reported in other studies (5-10%).

15105 students registered for the Sci-Write course. Compared with Open Knowledge, a relatively high proportion of subscribers, 1814, received a certificate of completion. Nearly all the students watched at least one video (14709 or 97.3%). There were a total of 2742 posts by the last week, but only 1210 (8%) of students ever posted on the forum. In the dataset we received, only 842 of the 1210 students that posted on the forum had a grade recorded (either pass or fail). If we consider only the students that watched at least one video as intending to complete the course, the completion rate is a relatively high 12.33%. In terms of statistics, Sci-Write seems to have the same problem as other MOOC courses: many learners wanted to complete the course, but only 12% did so, and only 8% of the students ever got involved in the forum.

A. Network Modeling and Student Behavior Features

The relational forum data was modeled as a node-link graph in Neo4j⁴. The Stanford Online Lagunita system employed a forum structure similar to what is seen on sites like Youtube and Facebook: one student makes an original post, and comments are given underneath. For each comment, there may be multiple nested replies. We started the modeling with a bipartite graph of students and posts, only making the assumptions that the content of comments were directed at the original post and that the content of replies were directed at comments. We did not assume that a reply addressed the reply that occurred just before it or to the original post. A graph of posts was connected based on this structure, and students were connected to their posts. From here, this bipartite network needed to be reinterpreted as a homogeneous network. Directed edges were created between students based on the structure of posts between them. This means that the network metrics calculated in Table I are from a student-student network. A visualization of the final result can be seen in Figure 1.

⁴<https://neo4j.com/>

A summary of the observed student behavior features which were recorded during the run of the two courses are shown in Table I. In addition to the “raw” student parameters, we calculated a number of network metrics for each week of each course based on our network model. Network metrics were calculated by exporting the Neo4j database to Gephi. These are shown in Table I as well. We also considered a change in “content” of the posts, using a sentiment score like Wen et al [13]. In our analysis, we calculated sentiment based on a legend of words, each with a sentiment score indicating positive or negative values. We used this legend to attach a positive sentiment and negative sentiment to each post (these features are also shown in Table I).

Interpretations of the network features for the context of MOOC forums are given here. In-degree represented how well a student’s posts attracted replies from other students (not just how many replies, but the total number of other students that replied). Similarly, out-degree counts how many other students to whom a student wrote a reply. Eccentricity represents how much a student became involved in the core discussion of the course, for instance, on the right side of Figure 1 we can see two tails of students coming off of the central component of the graph. These students did not engage the forum “superposters.” Next, although there is some uncertainty about interpretation of centrality in social networks [14], students with a high closeness centrality are the most socially “connected” students - for every other student, it is very likely that they either talked directly with that student or talked to someone who did. Similarly, students with a high betweenness-centrality were influential in the communication structure of the forum and were more likely in positions to act as “go-betweens” for other students. Eigenvector centrality is related but the metric is attempting to capture a student’s “influence” rather than just centrality. For example, a student that makes a post wherein all other superposters reply is likely to have a high eigenvector centrality. Finally, we can interpret authority and hub as students that made engaging posts and students that were likely to comment on engaging posts, respectively.

B. Course Dropouts and Certificate Rates

We constructed a definition of dropout based on the following points. First, students that registered for the course but exerted no effort to watch videos, post on the forum, or complete assignments should not be considered dropouts. It is likely these students never committed to taking the course. Second, students also came into the courses at different times and sometimes participated inconsistently, so defining a dropout as a period of absence from the course would mis-classify many students with low attendance but with a certified completion. Third, MOOC Instructors are also *not* likely to care about exactly *when* the student will drop out given that intervention can be taken as soon as possible.

We define a student’s dropout week as the last week in which activity was recorded. However, by definition, no student that completed the course should be considered a dropout. A dropout is thus a student who spent fewer than one standard deviation

Feature Name	Description
totalSecondsSpentInCourse	The system kept track of the time that students spent logged in.
totalVideosPaused	A paused video can indicate that the student is engaging with the material.
totalVideosPlayed	Not the same as loaded - indicates how many times the play button was pressed.
totalVideosLoaded	The total number of videos that were accessed by the student. This includes re-viewings.
totalVideosSearched	The total number of times the student skipped through a video - this can indicate a higher level of engagement.
totalVideosSpeedChanged	In the web player, students could change the speed of videos. This could indicate higher engagement.
averagePositiveSentiment	Averaged over a student's forum posts, this is a measure of the number of words used that have a positive connotation.
averageNegativeSentiment	Averaged over a student's forum posts, this is a measure of the number of words used that have a negative connotation.
averageEmotion	$\text{averagePositiveSentiment} + \text{averageNegativeSentiment}$
averageSentiment	$\text{averagePositiveSentiment} + \text{averageNegativeSentiment}$
indegree	The total number of other students that replied to this student's posts.
outdegree	The total number of other students that this student replied to.
degree	$\text{indegree} + \text{outdegree}$
eccentricity	In the network, the maximum distance between this student and any other student in the network.
closeness_centrality	A measure of how central a student is in the network - measured by sum of shortest paths to all other nodes.
betweenness_centrality	A measure of how many shortest paths in the network intersect the student.
authority	More authoritative students had a high in-degree and are pointed to by many hubs.
hub	Students with a higher hub score had a high out-degree and pointed to many authoritative students.
clustering_coefficient	A measure of how connected a student's neighbors are (proportion of how many edges of a full clique exist).
eigenvector_centrality	A measure of centrality in the network where connections to a highly-connected node are scored higher.
totalPosts	Total number of original posts created by the student.
totalReplies	Total number of replies created by the student.
totalUpvotes	Total number of upvotes this student received from other students.
connected	Whether or not this student is connected to the main component of the forum graph or if this student is an orphan or part of an orphaned clique
orphan	Whether or not this student is an orphan.

TABLE I
FEATURES OF STUDENT BEHAVIOR THAT WERE ANALYZED.

lower than this group in the total number of weeks in the course. For the OpenKnowledge course, students that completed the course spent 11.7 ± 3.5 , making the cutoff 8 weeks of effort, and a total completion group of size 267 (125 certificate holders). For Sci-Write, passing students spent 9.7 ± 1.6 , again making the cutoff 8 weeks, and a total completion group of size 2111 (1814 certificate holders).

IV. VISUAL AND STATISTICAL ANALYSIS

A. What features correlate with dropout? Does this change between courses?

The chart in Figure 2 shows the correlation coefficient between the student features in Table I and the dropout rates for both courses. Network features correlated more strongly with student attrition in Open Knowledge, while video and effort features correlated more strongly in Sci-Write. For each feature, the closer that the bars are together, the more “portable” the feature is between the two courses, meaning that this feature would be a good predictor of dropout despite the differences in course structure. The most portable features appear to be the total number of posts, more complex video watching behaviors such as speed change, positive sentiment, and total emotion.

B. How does the structure of the student interaction network change over time, and are there differences between the two courses?

Parameters of forum network structure can indicate the “health” of the forum, and monitoring these parameters over time gives insight into what topics, videos, or assignments spur student interest. Looking at differences between the OK and SW can also suggest the effects of requiring student participation

	WK1	WK4	WK7	WK11	WK15
Avg Deg (OK)	0.670	1.274	1.690	2.04	2.240
Avg Deg (SW)	0.779	0.928	1.091	1.210	
Mod (OK)	0.727	0.563	0.495	0.436	0.424
Mod (SW)	0.742	0.732	0.663	0.622	
Avg CC (OK)	0.015	0.021	0.030	0.035	0.036
Avg CC (SW)	0.008	0.014	0.023	0.024	

TABLE II
NETWORK STATISTICS AT VARIOUS POINTS IN EACH COURSE. AVG DEG = AVERAGE DEGREE, MOD = MODULARITY, AVG CC = AVERAGE CLUSTERING COEFFICIENT.

in the forum for course completion. The visual in Figure 1 suggests that the OK forum grew rapidly within the first month, but orphaned students remained present while the well-connected posters continued to write new content. The Sci-Write forum grew relatively more slowly, with a lower overall final density.

To better quantify the changes, we examined degree, modularity, and clustering coefficient over the course weeks. These parameters will best inform if the forums became more active over time (average degree), to what extent were they inclusive (modularity), and to what degree students engaged new peers (clustering coefficient). We provide this data in table format (Table II). When compared with Sci-Write, students in Open Knowledge continued to connect with new peers over time, as evidenced by rising average degree and clustering coefficient, and a gradually decreasing modularity.

The rising degree, decreasing modularity, and increasing clustering coefficients are also reflected in the quantity of orphaned students, shown in Figure 3. While the number of orphans in Open Knowledge starts high, the number quickly

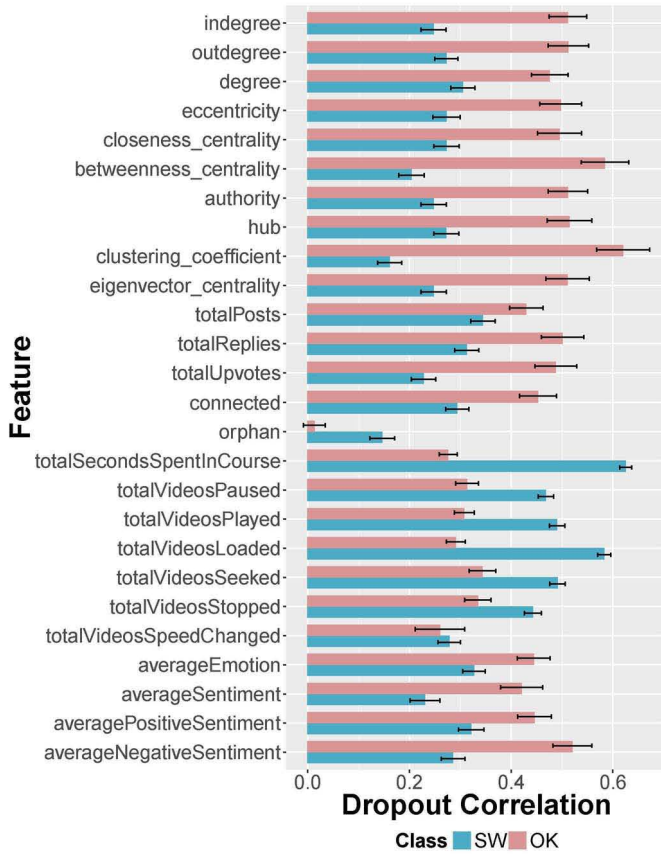


Fig. 2. Spearman correlations with dropout for both Sci-Write (blue) and Open Knowledge (red). Error bars are 95% confidence intervals. Network features correlated with dropout more strongly in Open Knowledge, while video and effort correlated more strongly in Science Writing. Despite Open Knowledge requiring more forum activity for coursework, “orphan” correlated more strongly with dropout in Sci-Write.

decreased when compared with Sci-Write. The hump in the Sci-Write trend may be indicative of many students starting to use the forum for the first time, without a corresponding increase in replies.

C. To what degree might further forum participation be encouraged if replies are given to a student’s original post?

A large number of orphaned students with unanswered posts were identified with visualization. Orphaned students could be considered “forum dropouts” (as opposed to overall course dropout, which would include other activities as well) and an instructor might be interested in knowing if spending the effort to engage such students, either by themselves or with other peers, is worthwhile. To start answering this question, we split the forum datasets based on “orphanhood”, and looked at relative completion rates. 18.6% of students in the central component of OK completed the course, compared to a much smaller 1.25% of the orphans. 65.7% of students in the central component of SW completed the course, compared to the 50.2% in the orphans. Of all orphans in Open Knowledge, only 22.3% ever received an upvote (15.3% for Sci-Write). Those that did

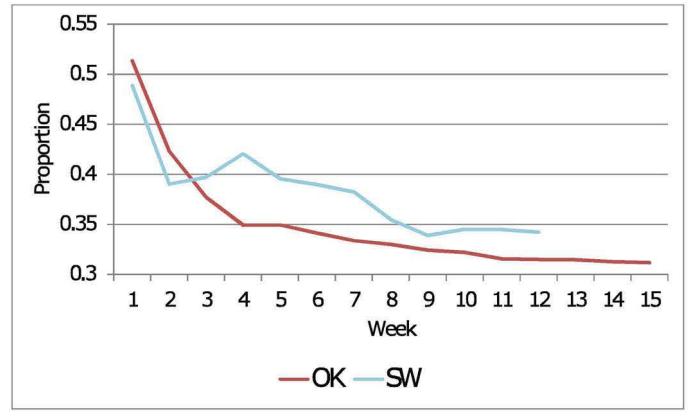


Fig. 3. The relative proportion of orphaned students declined in both courses over time, but remained consistently lower for Open Knowledge once the third week had passed.

receive any upvotes, however, received far more upvotes than the average post in the forums: 1.57 ($s=0.96$) vs 0.01 ($s=0.02$) for Open Knowledge, 1.89 ($s=1.49$) vs 0.01 ($s=0.02$) for Sci-Write. This suggests some orphaned students were creating highly influential posts that simply didn’t elicit a response, but the majority of orphans (about 80%) still received no reply and no upvotes. Even in the case where a response was not elicited by the original poster, community comments/feedback about content shared or upvotes would still encourage future use of the forum. Next, we considered using the semantic content of the posts to explain why some did not receive a reply. We utilized TF-IDF vector space models [15] with cosine similarity to analyze how similar orphan posts were to non-orphaned posts, as they may have been ignored for being duplicates or may have been answered elsewhere. We found that this was likely not the case. Posts with at least one other very similar (> 0.8) non-orphan post made up less than 0.5% of all posts in either course. Those with cosine similarity of > 0.5 were only 1.16% of all posts in Open Knowledge and 3.43% in Sci-Write.

It might be possible that replying to orphaned posts may cause those students to continue using the forum (similar to [5]). We looked at how students reacted in terms of posting behavior when they received replies to their very first post on the forum so we could better understand how intervention might help orphaned students. We also considered the amount of time that elapsed between the original post and the reply. We calculated the proportion of students that used the forum again after their original post had received a reply and compared them against the students who posted a second time without receiving any replies. It was found that the proportion of students that posted after receiving a reply was twice that of the other group, 68.8% vs. 34.24%, and that this proportion did not change much at all as weeks passed.

V. DISCUSSION OF RESULTS

First we will address the first research question from the introduction: from our correlation analysis (Figure 2), it appears

that different features had varying degrees of predictive power in each course. Visual representations of course data helped us understand how successful students used the forum, what types of forum behaviors were correlated with success, and which types of behavior were most highly correlated with success in either course. In the end, the total amount of effort that a student spent on the course correlated very positively with the outcome, but this was unsurprising. To further verify the differences between courses, a best-first feature analysis was done (using Weka⁵) on each course and the subsets of features were compared. As expected, in Open Knowledge, network features were chosen more often and in Science Writing video features were chosen more often. Referring back to Figure 2, the features with the highest “portability” were the total number of posts and whether the student engaged in particular video behaviors (seek, speed change). The majority of the features were not portable between courses. This observation could be explained either by the differences in the courses or the differences in the types of people that take different types of courses, however, we did not have access to data that could help control for these factors. Either way, the implication here is that when researchers design automated methods to flag students that need intervention, course parameters need to be considered as well.

Next we will address research questions (2) and (3). First, Open Knowledge became inclusively connected much faster than Sci-Write. It might be possible that this is due to the course requiring a minimum number of posts to obtain a certificate of completion, which caused more posts to be created, which in turn created more opportunities for replies. However, there are too many hidden variables in the course structures to really pin down a cause/effect relation. As more MOOC courses are conducted, it may be possible to collect enough data to make a recommendation about course structure, especially with regard to forum use. Second, orphaned students remained a problem throughout, but it remains unclear if just replying to their original posts is enough to spur further activity from them. An additional crowd-sourced experiment could conceivably be designed to get approximate measures of the quality of different original forum posts, whether or not it seems appropriate to reply to the content and whether it would be upvoted.

Some recommendations on these last two research questions: instructors might use forum network visualizations to quickly identify highly influential students (e.g. “superposters”) and monitor students that express negative sentiment. Additionally, they could monitor global network parameters such as clustering coefficient, modularity, and total number of orphaned students and compare them with other courses as feedback on their intervention measures. To assist orphaned posters, forums could be designed to highlight posts with no responses, rather than continually showing highly upvoted posts by default. While showcasing upvoted posts might be a valid strategy for content dissemination sites like Reddit and Facebook, they could be creating a divide between students with different

individual motivations that is not necessarily conducive to academic discussions.

VI. CONCLUSION

In this work, data from two MOOC courses was analyzed using multiple methods and multiple visualizations, including local node topography, global network features, video watching behavior, and sentiment. We conclude with a summary of our recommendations: (1) forum networks should be visually monitored over time by instructors to assess inclusiveness and forum success, (2) instructors and course designers should make efforts to identify and help orphaned forum students, and (3) designers of dropout detection systems should take course structure into consideration when making predictions. In this case-study, a visual and network analytics approach was successful in shedding light on student behavior.

REFERENCES

- [1] R. Rivard, “Measuring the mooc dropout rate,” *Inside Higher Ed*, vol. 8, p. 2013, 2013.
- [2] S. Halawa, D. Greene, and J. Mitchell, “Dropout prediction in moocs using learner activity features,” *Experiences and best practices in and around MOOCs*, p. 7, 2014.
- [3] D. F. Onah, J. Sinclair, and R. Boyatt, “Dropout rates of massive open online courses: behavioural patterns,” *EDULEARN14 Proceedings*, pp. 5825–5834, 2014.
- [4] T. Leary, *Turn on, tune in, drop out*. Ronin Publishing, 2009.
- [5] R. F. Kizilcec, E. Schneider, G. L. Cohen, and D. A. McFarland, “Encouraging forum participation in online courses with collectivist, individualist and neutral motivational framings,” *Experiences and best practices in and around MOOCs*, p. 17, 2014.
- [6] D. Coetzee, A. Fox, M. A. Hearst, and B. Hartmann, “Should your mooc forum use a reputation system?” in *Proceedings of the 17th ACM conference on Computer supported cooperative work & social computing*. ACM, 2014, pp. 1176–1187.
- [7] J. Huang, A. Dasgupta, A. Ghosh, J. Manning, and M. Sanders, “Superposter behavior in mooc forums,” in *Proceedings of the first ACM conference on Learning@ scale conference*. ACM, 2014, pp. 117–126.
- [8] J. Scott, *Social network analysis*. Sage, 2012.
- [9] C. P. Rosé, R. Carlson, D. Yang, M. Wen, L. Resnick, P. Goldman, and J. Sherer, “Social factors that contribute to attrition in moocs,” in *Proceedings of the first ACM conference on Learning@ scale conference*. ACM, 2014, pp. 197–198.
- [10] Y. Kotturi, C. Kulkarni, M. S. Bernstein, and S. Klemmer, “Structure and messaging techniques for online peer learning systems that increase stickiness,” in *Proceedings of the Second (2015) ACM Conference on Learning@ Scale*. ACM, 2015, pp. 31–38.
- [11] D. Coetzee, A. Fox, M. A. Hearst, and B. Hartmann, “Chatrooms in moocs: all talk and no action,” in *Proceedings of the first ACM conference on Learning@ scale conference*. ACM, 2014, pp. 127–136.
- [12] D. Yang, M. Wen, and C. Rose, “Peer influence on attrition in massively open online courses,” in *Educational Data Mining 2014*, 2014.
- [13] M. Wen, D. Yang, and C. Rose, “Sentiment analysis in mooc discussion forums: What does it tell us?” in *Educational Data Mining 2014*, 2014.
- [14] L. C. Freeman, “Centrality in social networks conceptual clarification,” *Social networks*, vol. 1, no. 3, pp. 215–239, 1978.
- [15] C. D. Manning, P. Raghavan, H. Schütze *et al.*, *Introduction to information retrieval*. Cambridge university press Cambridge, 2008, vol. 1, no. 1.

This material is based in part upon work supported by the U.S. Army Research Laboratory under Cooperative Agreement No. W911NF-09-2-0053. The views and conclusions contained in this document are those of the authors and should not be interpreted as representing the official policies, either expressed or implied, of the Army Research Laboratory or the U.S. Government. The U.S. Government is authorized to reproduce and distribute reprints for Government purposes notwithstanding any copyright notation here on.

⁵<http://www.cs.waikato.ac.nz/ml/weka/>