

Graph Mining to Characterize Competition for Employment

Andrew Toulis and Lukasz Golab
University of Waterloo
Waterloo, Ontario, Canada
{aptoulis,lgolab}@uwaterloo.ca

ABSTRACT

In this paper, we discuss a novel application of graph analytics to characterize competition in the workforce. We propose a methodology that relies on finding communities in a graph representing prospective employees (with edges connecting people who interviewed for the same job) and communities in a graph representing available jobs (with edges connecting jobs that interviewed the same person). We then apply the proposed methodology to a real dataset corresponding to co-operative internships offered to undergraduate students at a North American post-secondary institution, illustrating the benefits of our approach.

1. INTRODUCTION

Many applications involve relationships that can naturally be expressed as graphs: friend/follower relationships in social media, hyperlink relationships in the World Wide Web, chemical structure and protein interactions in bioinformatics, etc. Graph mining techniques such as clustering and community detection can identify interesting structure in such relationships.

We discuss a novel application of graph mining in the context of the workforce to understand competition for employment. This is an increasingly important application domain: the job market has become very competitive due to forces such as globalization and technological change [1]. As a result, employers are becoming more data-driven in their hiring processes [3]. Current data analysis efforts focus on measuring time to hire, hiring manager satisfaction, and performance of a hire [6], as well as measuring the effectiveness of employee benefits programs [3]. However, there remain critical gaps in the job market that employers alone cannot tackle. For example, employers may not have a good understanding of the available talent pool and may not be allocating their recruiting resources effectively. Likewise, employees may not be aware of the extent of competition for various types of jobs and therefore they may not know which jobs are realistically within their reach.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

NDA'17 May 19, 2017, Chicago, IL, USA

© 2017 ACM. ISBN 978-1-4503-4990-1/17/05... 15.00

DOI: <http://dx.doi.org/10.1145/3068943.3068946>

LinkedIn is perhaps the most common example of the potential of graph analytics in the context of the job market: its users establish connections to other users, thereby creating professional networks whose properties have been investigated [11, 12]. LinkedIn data have been used to develop graph-based skill extraction algorithms that are useful for hiring [9] and have inspired link recommendation algorithms [13]. Furthermore, a natural example of network effects is employee referrals, which have been shown to generate higher quality hires with better retention and which LinkedIn offers as a service [6].

In this paper, we go a step further and perform an end-to-end analysis of an entire job market. Our analysis is enabled by a unique dataset from a North American post-secondary institution, corresponding to 4100 undergraduate students competing for 2000 co-operative internships. The dataset includes information about every job interview and hiring decision that took place during the summer 2016 season. The questions we want to answer include:

1. Are there natural clusters of employers and employees, and if so, what are the defining characteristics of each cluster?
2. Can we rank employees and jobs into tiers based on the corresponding networks?
3. Which employers attract good prospective employees and which prospective employees obtain interviews for sought-after positions?

To answer these questions, we propose a graph-oriented methodology to characterize workforce competition. Our methodology is based on finding communities in graphs induced by the job market. We construct an employee graph by connecting two employees if they interview for at least one job in common. Similarly, we construct a job graph by connecting two jobs that interview at least one employee in common. We then perform community detection on both graphs, describe each community using its unique characteristics, and identify relationships among the community structures in the two graphs.

To summarize, we make the following three contributions:

1. We apply graph mining to a novel application domain of employment/workforce.
2. We propose a methodology to characterize competition for jobs and for prospective employees using graph mining (using community detection in the corresponding networks).

Table 1: Job and student data

Table	Attributes
Students	<u>StudentID</u> , Academic Year
Jobs	<u>JobID</u> , Title, Employer, Location, Industry
Interviews	<u>JobID</u> , <u>StudentID</u> , Hired

3. We apply our methodology to a unique dataset consisting of undergraduate students as prospective employees and employers offering co-operative internships. We show that our techniques lead to actionable insight for the benefit of employees and employers.

The remainder of this paper is organized as follows. We discuss related work in Section 2; we describe our dataset in Section 3; we present our methodology in Section 4 followed by our results in Section 5; and we conclude in Section 6.

2. RELATED WORK

This paper is related to two bodies of work: graph mining and workforce/employment analytics. In the context of graph mining, there are several standard analysis techniques such as clustering and community detection, which have been widely successful in understanding complex networks [4]. We also use standard analysis techniques but we apply them to a novel application domain.

Social media services like LinkedIn are becoming increasingly used by employers for recruiting and by employees for job-hunting, and the importance of network effects such as employee referrals are becoming increasingly important to employers [6]. Employers recognize the importance of understanding competition, and there has been prior work on analyzing the LinkedIn graph [9, 11, 12, 13]. However, we are not aware of any previous work on characterizing competition among employees and employers in a large job market.

Finally, there has been prior work on mining co-operative employment data [7, 8], but competition in the co-operative workplace has not been studied at the employee and employer level.

3. DATA

The dataset used in this analysis consists of 4100 undergraduate students from a North American post-secondary institution. These students competed for nearly 2000 jobs from 700 distinct employers over a two-month period in summer 2016. On average, each student had 3.5 interviews, and each job interviewed 7 students. Some jobs interviewed many more than seven students per job and hired multiple students for the same position. Approximately 50 percent of students were hired by some employer and 75 percent of jobs were filled with at least one student. This is indicative of a highly competitive environment.

The information we have for students and jobs is summarized in Table 1, with primary keys underlined. We have the year of study (first year through fourth year) of each student. For each job opening, we have the job title, employer name, location, and an industry label which indicates whether the job is related to Information Technology (IT). The industry label was manually assigned by the data providers and contained incorrect and missing values, which we will treat through community detection.

Table 2: Student academic year sizes and hire rates

Academic Year	Students	Students Hired	Hired (%)
1	375	180	48
2	1450	695	48
3	975	495	51
4	675	370	55

ALGORITHM 1: Graph mining to analyze the job and student graphs

Data: Student Graph S and Job Graph G

Result: A labeled set of student and job communities

1. Perform community detection on S and G separately;
 2. Using domain knowledge, rank the job communities in G into tiers: Top, Mid, Low, and other;
 3. Label students in S who were interviewed/hired with each job tier;
 4. Describe the unlabeled student communities based on the state of their nodes as determined by (3).
-

For each interview, we have information about the job, the participating student, and whether the student was hired for the job.

Hiring statistics for each academic year are shown in Table 2. Senior students (year 3 and especially 4) were more likely to be hired than junior students (years 1 and 2). Some student records were missing the academic year, which explains why the total number of students in Table 2 does not sum up to 4100.

4. METHODOLOGY

Algorithm 1 summarizes our graph mining methodology. We begin by constructing job and employee graphs. Edges are undirected and unweighted in both graphs. The job graph is formed by connecting jobs that interviewed at least one student in common. The student graph is formed by connecting students who have at least one interview in common.

Table 3 is an example set of interviews across 9 students (labeled 1-9) and 8 jobs (labeled A-H). The corresponding job and student networks are illustrated in Figures 1 and 2 respectively. Ignore the coloured communities for now.

With a job and student graph in hand, we run community detection on both graphs. We use the Louvain Method, implemented in the Networkx Python package [5], which aims to optimize connections within communities while minimizing connections between communities [2, 10]. We then inspect the characteristics of students and jobs in each community to identify representative features. This step helps handle missing data and noisy signals since students and jobs that are similar tend to compete with each other. For example, we observed that a community of jobs labeled as Information Technology jobs were actually oriented around user experience (details in Section 5).

The community features of the job network were easily interpretable given our domain knowledge of co-operative education. To scope down the manual labeling process, we focused on the largest industry in the dataset: Information Technology (IT). Given the institution’s specialty in

Table 3: Example table of interviews

Job ID	Student IDs
A	1, 2
B	1, 2
C	2
D	1, 3, 4, 5
E	5, 6
F	6, 7, 8, 9
G	7, 8, 9
H	7

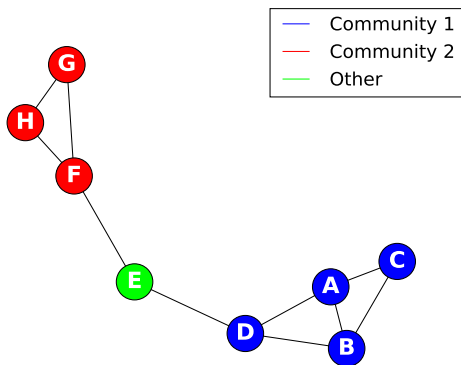


Figure 1: A job network based on the data from Table 3, coloured by job community

computing, internships offered by IT companies in California and Seattle were deemed the “best”. Such companies include, for example, Amazon, Apple, Facebook, Google, LinkedIn, Microsoft, Snapchat, Uber and Yelp. As it turns out, these employers competed in a single community of top jobs, as discussed in Section 5.

This domain knowledge-driven labeling step made it possible to rank the communities of jobs by their quality. In return, this new information was used to categorize communities in the student network. This was done by inspecting the distribution of job quality that students in a given community were competing (interviewing) for. Labeling can also be done based on where students were hired, although analyzing either distribution yielded similar results.

Recall Figures 1 and 2. The example job network is coloured based on communities detected. The communities in the student network are then colored based on the job communities in which the students had the most interviews. For example, student community 1, containing students 1–5, is colored blue because these students interviewed for jobs in job community 1 which is colored blue.

To zoom in on important nodes in each community, we extend our approach by performing centrality analysis in the job network. Centrality can be used to identify the most important nodes in a graph. In our case, central jobs are typically multi-disciplinary roles which interview a diverse set of students. While communities allow us to characterize competing groups of jobs and students, looking at central nodes gives us a more concrete view into a community. We find that central nodes complement community detection

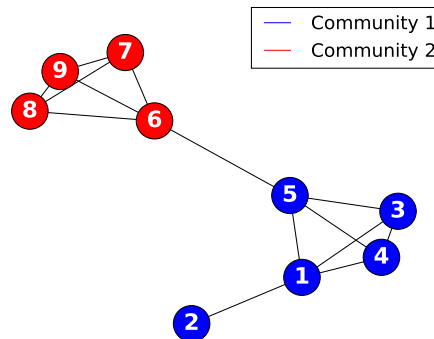


Figure 2: A student network based on the data from Table 3, coloured by the job communities from Figure 1

by acting as representatives for their communities. This provides us a chance to validate the results of the Louvain Method which is influenced highly by central nodes.

5. RESULTS AND DISCUSSION

5.1 Job Communities

In the job network, the Louvain Method found eight communities, three of which we labeled “IT communities” after inspecting the job titles. Table 4 summarizes each community in terms of its size, density and percentage of jobs labeled as IT.¹ IT was the only discipline with several communities, which naturally invites us to rank them. Furthermore, the IT industry label had missing and incorrect values, which is why the IT communities have less than 60% of their jobs previously known to be in IT.

5.1.1 Job Network Properties

With the exception of Communities H and E, it is not obvious how the communities differ from each other based on their network properties shown in Table 4. Community H is small and is mostly composed of a single job type (user experience). Community E is dense, has the lowest diameter for its size, and has a high average degree per node. This establishes Community E as the most competitive community.

5.1.2 Job Community Labeling

Community E is the Top Tier IT Community: it contains the best jobs at top companies. A large majority of students interviewing for these jobs were in their fourth (final) year of study, with the remaining students mostly in their third year of study. We ranked Community B second in IT (Mid Tier). It contained small IT companies who interviewed more junior students, mostly in their second year. The final IT community, Community G, had mostly quality assurance and software testing roles, which is perceived by students as less desirable work (Low Tier). Most students competing for

¹Community H had a large percentage of jobs labeled as IT, but was completely composed of user experience roles so we adjusted the percentage to zero.

Table 4: Summary of the job communities

Community	Community Type	All Jobs	Edges	Density (%)	Diameter	Avg. Degree	IT Jobs (%)
All		1970	25700	1.3	6	26	29
A	Physical Engineering	390	3325	4.4	6	17	12
B	Mid Tier IT	365	4150	6.3	5	23	60
C	Finance	280	2550	6.4	7	18	9
D	Research & Teaching	280	2225	5.6	6	16	9
E	Top Tier IT	240	4200	14.5	4	35	60
F	Civil Engineering	200	1550	7.6	6	15	0
G	Low Tier IT	185	1425	8.4	5	15	42
H	User Experience	30	150	36.5	4	10	0

these jobs were in their first year and had little prior work experience.

Community H exclusively contained user experience roles. Community A had primarily mechanical engineering, manufacturing and hardware jobs, which we label as “Physical Engineering”. Community C was primarily composed of financial firms. Although Communities A and C were non-IT communities, many of their most interviewed jobs were IT-like. For example, the most competitive roles in the Finance community were in data science.

5.1.3 Job Centrality Analysis

To zoom into each community, we extracted the top ten nodes with the highest closeness centrality (that is, the jobs with the smallest average shortest path length to other jobs). In the Finance community, central nodes were data analyst jobs at a consulting firm and bank trading floor. In the Physical Engineering community, the central node was a manufacturing job from a large automotive company. In the Top Tier IT Community, a data scientist position was the most central role.

These central jobs are interdisciplinary, competing with many types of jobs and interviewing a diversity of students such as younger students. An advantage of this in co-operative education is that it enables students to change industries, say from software engineering to data science. Furthermore, finding central jobs allows us to explain interesting phenomena, such as lower-year students competing for top jobs.

5.1.4 Job Cross-Community Competition

The community network, where nodes are the communities detected, is shown in Figure 3. Node sizes correspond to the number of jobs in the corresponding community. The thickness of an edge visualizes the overlap of students interviewed between communities (in a directed way).² Defining $S()$ as the set of students who interviewed for at least one job in a given job community, we calculate the directed edge weight between communities X and Y as the conditional probability of a student interviewing for a job in Y given that he or she interviewed for a job in X .

$$Weight(X, Y) = |S(X) \cap S(Y)| / |S(X)|.$$

As seen in Figure 3, the job community network is a clique, which was not expected due to the large differences

²The edges leaving Community H were made thin since many students who interviewed for these jobs also interviewed elsewhere due to its small size.

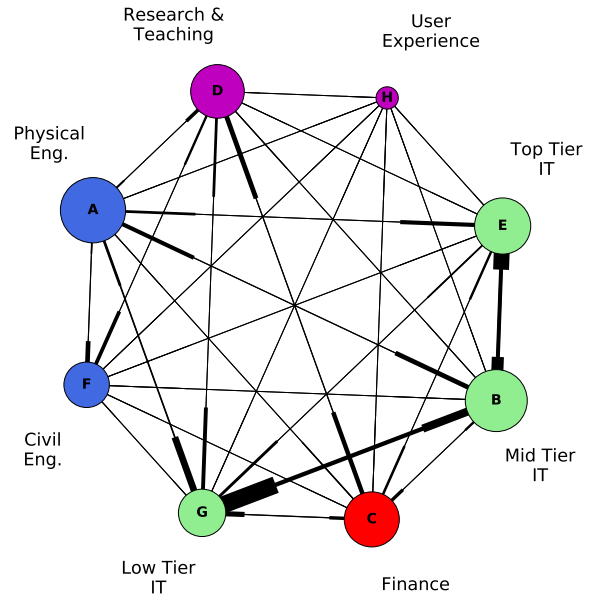


Figure 3: The labeled job community network

between communities and the exclusivity of interviewing for Top Tier IT jobs. Furthermore, we find that the Low Tier IT Community interviewed many students who are also seeking jobs from other communities, especially from the Mid Tier IT Community. Likewise, there are connections between the Top and Mid Tier IT Communities, where top students could be applying to lower tier jobs as a safety net. Validating our ranking, there are few connections between the Top and Low Tier IT Communities.

5.1.5 Job Outlier Detection

We now zoom into the Top Tier IT Community and inspect top jobs that did not hire any student.

Manual inspection of these unfilled “top” jobs showed that they were outliers: they were interesting enough to attract students, but their actual quality was lower. The students these jobs interviewed were hired by other top jobs. Manual extraction of job descriptions showed that these jobs commonly exaggerated the role a student would be hired for (e.g., “Big Data Hacker”).

It was also common to find IT start-ups competing against large employers in California and Seattle. These large em-

Table 5: Hiring statistics of the different job communities

Community	Interviews/Job	Hired/Job
Network	7	1.10
A (Physical Engineering)	7	1.12
B (Mid Tier IT)	7	1.08
C (Finance)	6	1.09
D (Research & Teaching)	6	0.90
E (Top Tier IT)	11	1.43
F (Civil Engineering)	6	1.18
G (Low Tier IT)	6	0.95
H (User Experience)	5	0.79

employers are able to provide the best offers to students, including high salaries and free living accommodations. Furthermore, these large employers commit more resources to interviewing and hiring students. Table 5 shows the number of interviews per job and number of students hired per job in each community. While some start-ups appeared in the Top Tier IT Community, most of them interviewed a below-average number of students. Thus they are not fully aware of the extent of competition, creating missed opportunities for hiring.

5.2 Student Communities

There are nine student communities in the network derived from the Louvain Method. We proceed to use the information from the job communities to label the student communities. Table 6 shows the job community makeup of each student community in the “Job Interview Distribution” column, as per step four of our algorithm. Immediately we find that several student communities (2, 3, 4, 5, 6) interviewed mostly with a single job community. This allows us to instantly understand these student communities based on their corresponding job communities. With this information, there are four clear IT student communities.

Before we discuss this further, we look at information we would have about student communities if we did not perform community detection on the job graph. The number of jobs students competed for, the number of interviews they had, and the percentage of students hired for any job are tabulated on the far right of Table 6. It is not obvious how the student communities are different if we only considered their success rates. For example, Community 1 has a higher success rate than Community 8, although we find that the tiers of IT jobs interviewed for are lower in Community 1.

5.2.1 Student Network Properties

From Table 6 we see immediately that Community 6 is extremely dense and has a low diameter. In particular, the average degree between students in the community is 132. Certainly this student community is exceptionally competitive. Community 8 is also dense and has a low diameter which is partially due to its small size.

5.2.2 Student Community Labeling

It is not obvious from the success rates that student Community 6 had nearly all (90%) students interviewed by top employers and thus is the best student community. Despite not having access to student grades in our dataset, manual extraction of resumes showed these students had exceptional

academic performance and extracurricular involvement compared to students in other IT communities.

Community 8 contained several students who interviewed with IT start-ups. This community consisted mostly of second year students, which shows that IT start-ups do not compete for the most experienced students. There were only three large Top Tier IT employers who interviewed and hired students in Community 8.

Students in Community 1 mainly interviewed for Mid and Low Tier IT jobs at smaller companies, often involving software testing. These students were in their first and second year of study. This demonstrates that smaller IT companies are unable to compete for the same batch of students as Top Tier IT companies.

Community 7 had engineering students who competed for both Physical Engineering and IT jobs. These students, as well as students in Community 2 (Finance), showed an interest in IT jobs based on their distribution of interviews. This reveals that hardware and finance talent is trending towards IT. Likewise, we have noticed an increase of employers in these industries hiring students for IT jobs in data science. On the other hand, certain types of engineering communities, especially Civil Engineering, remain less interested in IT.

5.2.3 Student Cross-Community Competition

The student community network, where nodes are the communities detected, is shown in Figure 4. Node sizes indicate the number of students in the corresponding community. The thickness of an edge visualizes the overlap of jobs competed for between communities.³ Defining $J()$ as the set of jobs that interviewed at least one student in a given student community, the directed edge weight between student communities X and Y is the conditional probability of a job interviewing a student from community Y given that it interviewed a student from community X :

$$Weight(X, Y) = |J(X) \cap J(Y)| / |J(X)|.$$

The student community network is not a full clique, since Community 5 (Civil Engineering) did not compete for any jobs that students in Community 8 (Top/Mid Tier IT) competed for. This strengthens the notion that Civil Engineering students are not trending into IT as much as other students. On the other hand, Community 2 (Finance) has a strong connection to Community 1 (Mid/Low Tier IT), which indicates a shift in the industry. The strongest discovery is that many physical engineering students in Community 7 are trending heavily towards IT, due to their connections with all IT student communities. We see that this community is very different than the other physical engineering one (Community 4) due to their low connection.

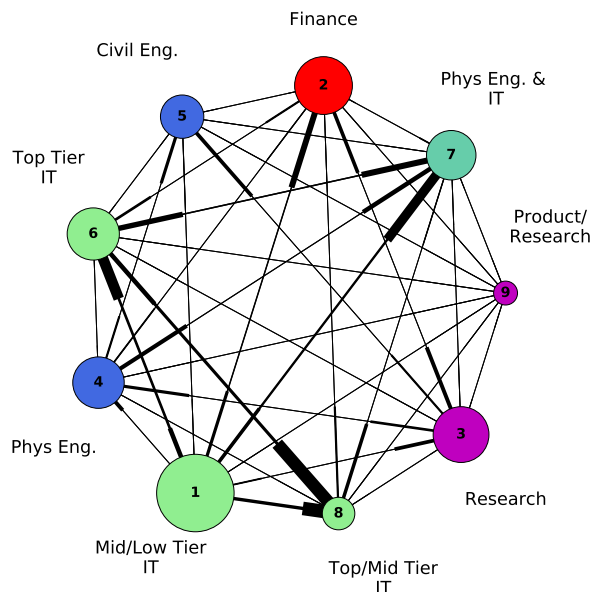
We also see that IT students in Community 8 interviewed for more lower tier IT jobs than students in Community 6. In fact, 75% of the jobs students interviewed for in Community 8 were also interviewed for by students in Community 1.

Finally, we obtained results that strengthen our findings from cross-community competition in jobs (Section 5.1.4). Students from Community 1 (Mid/Low Tier IT) are generally unable to compete for the same jobs as students in Com-

³The edges leaving Community 9 were made thin since its students tended to compete with other student communities due to its small size.

Table 6: Student communities in the network

Community	Job Interview Distribution	Students	Edges	Density (%)	Diameter	Avg Degree	Jobs	Interviews	Hired %
All		4120	80075	0.9	7	39	1970	14275	52
1	Mid/Low Tier IT (55/26%)	1005	11125	2.2	6	22	860	3725	56
2	Finance (83%)	640	7225	3.5	6	23	470	2075	48
3	Research (81%)	625	5475	2.8	8	17	440	1850	47
4	Phys Eng. (93%)	555	8250	5.4	6	30	370	1625	47
5	Civil Eng. (93%)	415	2975	3.5	7	14	260	1200	55
6	Top Tier IT (90%)	410	27025	32.4	4	132	380	2125	60
7	Phys Eng. (65%) & IT (Mixed)	325	2825	5.4	5	17	340	1175	52
8	Top/Mid Tier IT (64/30%)	85	1450	42.3	4	35	145	275	53
9	Product/Research (52/23%)	60	275	14.8	8	9	75	225	53

**Figure 4: The student community network with communities labeled using the job communities**

munity 6 (Top Tier IT). However, top tier students could interview for lower tier jobs as a back-up plan.

5.2.4 Student-Job Asymmetry

We conclude by analyzing what happened to the best students who were not hired. There is an asymmetry in the causes of unfilled students compared to unfilled jobs. Manually extracting student records after the fact, we found many top students were hired outside of the institution’s internship system. These students were so exceptional that they could find employment on their own. This was not the case for other classes of students, who indeed were unable to find a job. The unmatched top jobs (recall Section 5.1.5), on the other hand, were outliers in the other direction due to having fewer resources for interviewing and hiring, or due to being lower quality jobs.

6. CONCLUSIONS

In this paper, we proposed a graph mining methodology for workforce data and demonstrated its effectiveness on a large co-operative internship dataset. We uncovered the nature of competition in the network, including the top employee and employer communities. We successfully found the types of jobs that compete with each other, showing that start-ups, for example, often cannot compete with the largest Information Technology companies in the world. We also showed that natural clusters in the job network can be utilized to rank prospective employees who are otherwise equal. We overcame data limitations to understand each community of employees better and form better-justified tiers of employees and jobs.

Our dataset contained many students who were not hired and jobs that did not hire any student. We discovered several types of jobs that were unaware of the level of competition for students, including small companies and start-ups. Thus we identify several missed opportunities in the hiring process, which can be naturally uncovered through network analysis. In particular, future work can be done to form recommendations via link prediction.

7. REFERENCES

- [1] D. Acemoglu, D. Dorn, G. H. Hanson, B. Price, et al. Import competition and the great us employment sag of the 2000s. Technical report, National Bureau of Economic Research, 2014.
- [2] V. D. Blondel, J.-L. Guillaume, R. Lambiotte, and E. Lefebvre. Fast unfolding of communities in large networks. *Journal of Statistical Mechanics: Theory and Experiment*, 2008(10):P10008, 2008.
- [3] T. H. Davenport, J. Harris, and J. Shapiro. Competing on talent analytics. *Harvard Business Review*, 88(10):52–58, 2010.
- [4] S. Fortunato. Community detection in graphs. *Physics reports*, 486(3):75–174, 2010.
- [5] A. A. Hagberg, D. A. Schult, and P. J. Swart. Exploring network structure, dynamics, and function using NetworkX. In *Proceedings of the 7th Python in Science Conference (SciPy)*, pages 11–15, 2008.
- [6] S. Howell. Global recruiting trends 2016, September 2015. https://business.linkedin.com/content/dam/business/talent-solutions/global/en_us/c/pdfs/GRT16_GlobalRecruiting_100815.pdfbusiness.linkedin.com

[Online; posted 22-September 2015].

- [7] Y. Jiang and L. Golab. On competition for undergraduate co-op placements: A graph mining approach. In *Proc. Int. Conf. on Educational Data Mining (EDM)*, pages 294–299, 2016.
- [8] Y. Jiang, S. Lee, and L. Golab. Analyzing student and employer satisfaction with cooperative education through multiple data sources. *Asia-Pacific Journal of Cooperative Education*, 16(4):225–240, 2015.
- [9] I. Kivimäki, A. Panchenko, A. Dessy, D. Verdegem, P. Francq, C. Fairon, H. Bersini, and M. Saerens. A graph-based approach to skill extraction from text. *Graph-Based Methods for Natural Language Processing*, page 79, 2013.
- [10] R. Lambiotte, J.-C. Delvenne, and M. Barahona. Laplacian dynamics and multiscale modular structure in networks. *arXiv preprint arXiv:0812.1770*, 2008.
- [11] P. Lops, M. de Gemmis, G. Semeraro, F. Narducci, and C. Musto. Leveraging the linkedin social network data for extracting content-based user profiles. In *Proc. of the fifth ACM conference on Recommender systems*, pages 293–296, 2011.
- [12] M. A. Russell. *Mining the Social Web: Data Mining Facebook, Twitter, LinkedIn, Google+, GitHub, and More.* ” O’Reilly Media, Inc.”, 2013.
- [13] Z. Yin, M. Gupta, T. Weninger, and J. Han. Linkrec: a unified framework for link recommendation with user attributes and graph structure. In *Proc. of the 19th international conference on World Wide Web*, pages 1211–1212, 2010.