

CSE 6242: Team #131 (GeoSpecial): Project Proposal

What are you trying to do? Articulate your objectives using absolutely no jargon.

Using Newman's Ground Truth Algorithm (NGTA) on a patent network extracted from the PatentsView database, we will filter prolific inventors from network noise. Because the patent network is noisy, we are skeptical that it accurately reflects reality. NGTA will allow us to quantify our confidence that each edge is "real" and not noise. NGTA will produce statistics that quantify the reliability of the network.

How is it done today; what are the limits of current practice?

Noisy graphs are often filtered on weights. This assumes a linear relationship between weights and confidence. Robust ML solutions, such as Namata's Coupled Collective Classifiers (C3)², can perform network inference (extract the "real" graph from the observed). These approaches use semi-supervised learning and inference models², which are computationally expensive. While some semi-supervised models process noisy graphs, they are usually limited to clustering, (do not assess components individually¹²). Others use ensembles (generative stochastic graph models)³. These rarely have closed-form solutions, requiring simulations. A common goal with noisy networks is community detection,^{15, 16} which also does not assess individual components). A method developed by Dr. Du (Beijing U of T) processes patent networks to find influential inventors¹⁸. His is more efficient than prior approaches, but does not distinguish ground truth. If non-noisy patent data is acquired at scale, comparisons between Du's and Newman's methods would be insightful. In 1999, US patents contained 3.8 million nodes. This data is rich, with attributes for inventors, assignee, etc.⁵ This richness results in noise, with innovators hidden by peripheral individuals. Older patent records can suffer from duplicate inventors, which present as false discoveries⁶. This data has historically been used to describe and forecast economic conditions, with common algorithms including shortest paths, and maximum flow.⁵ Modern patent analysis has matured in several areas, including patent infringement, hotspots, and competitors¹⁷. Traditionally the talent industry relied on passive acquisition⁸. Today, "head-hunting" is limited to small groups. Patent data is rarely used, and only as a reference. NGTA brings patents to the forefront, enabling a data-driven recruiting strategy.

What's new in your approach? Why will it be successful?

In 2000, research on patent networks was primarily trend analysis of global structures, which fueled macroeconomic decisions¹³. Recently, patent analysis has been utilized to find impactful inventions and highly innovative companies¹⁴ - largely fueling stock speculations. Never before has NGTA been executed on a patent network. NGTA will enable us to quantify confidence in each edge's existence in the network (where nodes represent inventors/ businesses, edges represent patents), then filter the network dynamically by confidence levels. After filtering, remaining connected nodes constitute significant inventors. Thus, NGTA extracts inventors who are "real" contributors, and filters out those who are less important to a company's patent portfolio. NGTA produces three network-level descriptive statistics. The true positive rate (low TP implies missing edges¹), the false positive rate (low FP implies an edge is rarely observed where none exists¹), and false discovery rate (low FDR implies observed edges are part of the network¹). This implementation has a strong chance of success due to the algorithm being lightweight, with rapid convergence. NGTA can be run on tabular data, and does not require a graph database. The network will be instantiated using NetworkX, and visualized with Graphviz⁴. These open source packages offer robust, effective, and lightweight processing and interactive visualization capabilities.

Who cares?

Most empirical studies of networks take a naive view of structural data, where one assumes that the data are the network¹. Accurate analysis and understanding of networked systems requires a way of estimating the true structure of networks from such rich but noisy data¹. NGTA quantifies this discrepancy, and level-sets confidence in downstream outputs. Human capital has a profound impact on a firm's strategy, outcomes, and performance. Furthermore, it is the most challenging resource for competitors to replicate.⁹ In an ever-competitive war for talent, companies will increasingly turn to data-driven recruiting strategies. Patent networks are directly related statistically to a firm's acquisition and dispersion of new technology⁷. Identifying true innovators allows firms to poach those people for competitive advantage. NGTA allows recruiters to confidently identify the best individuals to pursue.

If you're successful, what difference and impact will it make, and how do you measure them (e.g., via user studies, experiments, ground truth data, etc.)?

Successful implementation of NGTA will enable organizations to greatly increase their confidence in any business intelligence based upon the network. Recruiters can focus their efforts where the effort is worthwhile - the pursuit of valuable candidates. NGTA can be tested using a test set. If the resulting statistics are similar between the two, we can be confident that the algorithm reached true convergence.

What are the risks and payoffs?

This is a low-risk, high-reward implementation. One risk is that the network exhibits heteroskedasticity across industries, prohibiting NGTA from converging. If so, the network will be segregated by industry, with NGTA run upon sub-networks. Successful NGTA implementation upon patents results in business intelligence for the talent industry - extracting performant individuals, upon whom recruiters can focus.

How much will it cost?

Nothing. We can operate on the patent data in GCP within the limits of the free tier. While weekly updates to the US patent database can be hundreds of megabytes in size, the *patentpy* and *patentr* libraries allow for tidy and small records¹⁰. This greatly reduces the required data volume. Established conversion models between relational data and target graph data will further reduce the data volume¹¹.

How long will it take?

TASK TITLE	TASK OWNER	START DATE	DUE DATE	DURATION	% COMPLETE
Data Acquisition and Cloud Storage	Sang Yoon, Erik S	10/15	10/22	7	0%
Data Engineering	Xingpeng, Erik W	10/22	10/29	7	0%
Algorithm Development & Cloud Execution	Tylor	10/29	11/5	6	0%
Front End Development	Sandro, Xingpeng	11/5	11/12	7	0%
Front End Testing	Everyone	11/12	11/19	7	0%

* All team members will contribute a similar volume of effort.

What are the midterm and final "exams" to check for success? How will progress be measured?

The midterm consists of a small scale experimental run of NGTA on a subset of the network, with limited visualization. The final consists of a full run of NGTA on the full network, with interactive visualization.

References:

1. Newman, Mark EJ. "Network structure from rich but noisy data." *Nature Physics* 14.6 (2018): 542-545.
2. Namata, Galileo Mark, Stanley Kok, and Lise Getoor. "Collective graph identification." *Proceedings of the 17th ACM SIGKDD international conference on Knowledge discovery and data mining*. 2011.
3. Casiraghi, Giona, et al. "From relational data to graphs: Inferring significant links using generalized hypergeometric ensembles." *International conference on social informatics*. Springer, Cham, 2017.
4. Faysal, Md Abdul Motaleb, and Shaikh Arifuzzaman. "A comparative analysis of large-scale network visualization tools." *2018 IEEE International Conference on Big Data (Big Data)*. IEEE, 2018.
5. Batagelj, Vladimir, et al. "Analyzing the structure of US patents network." *Data science and classification*. Springer, Berlin, Heidelberg, 2006. 141-148.
6. Li, Guan-Cheng, et al. "Disambiguation and co-authorship networks of the US patent inventor database (1975–2010)." *Research Policy* 43.6 (2014): 941-955.
7. Duguet, Emmanuel, and Megan MacGarvie. "How well do patent citations measure flows of technology? Evidence from French innovation surveys." *Economics of innovation and new technology* 14.5 (2005): 375-393.
8. Rose, Jacqueline A. "Building An Internet Recruiting Strategy For A Big Five Professional Services Firm." (1999).
9. Valenti, Alix, and Stephen V. Horner. "Leveraging board talent for innovation strategy." *Journal of Business Strategy* (2019).
10. Yu, James, et al. "Accessing United States Bulk Patent Data with patentpy and patentr." *arXiv preprint arXiv:2107.08481* (2021).
11. De Virgilio, Roberto, Antonio Maccioni, and Riccardo Torlone. "Converting relational to graph databases." *First International Workshop on Graph Data Management Experiences and Systems*. 2013.
12. Chang, Jui-Hung, and Hsiu-Chen Weng. "Fully used reliable data and attention consistency for semi-supervised learning." *Knowledge-Based Systems* 249 (2022): 108837.
13. Wu, Chao-Chan, and Ching-Bang Yao. "Constructing an intelligent patent network analysis method." *Data Science Journal* (2012): 011-003.
14. Chakraborty, Manajit, Maksym Byshkin, and Fabio Crestani. "Patent citation network analysis: A perspective from descriptive statistics and ERGMs." *Plos one* 15.12 (2020): e0241797.
15. He, Kun, et al. "Hidden community detection in social networks." *Information Sciences* 425 (2018): 92-106.
16. Newman, Mark EJ, and Michelle Girvan. "Finding and evaluating community structure in networks." *Physical review E* 69.2 (2004): 026113.
17. Abbas, Assad, Limin Zhang, and Samee U. Khan. "A literature review on the state-of-the-art in patent analysis." *World Patent Information* 37 (2014): 3-13.
18. Du, Yong-ping, Chang-qing Yao, and Nan Li. "Using heterogeneous patent network features to rank and discover influential inventors." *Frontiers of Information Technology & Electronic Engineering* 16.7 (2015): 568-578.