

Patently True

Introduction

Patent data is inherently noisy, with many inventors for each patent. This makes it difficult to determine those who actually contribute to meaningful work. Furthermore, the data is quite large, making manual filters difficult - if not impossible - to implement efficiently or effectively. We propose using Newman's Ground Truth Algorithm (NGTA) to efficiently filter the noisy patent data and identify prolific inventors. NGTA's outputs can then be used for recruiting purposes, to hire top talent.

Problem Definition

Using 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.

Survey

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.

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.

Proposed Method

NGTA is the result of academic research by Dr. Mark Newman, published as *Network Structure from Rich but Noisy Data* in 2018. Newman discusses how data representing a network contains a hidden, "true" network structure that is obscured by noisiness. When a noisy network is not scrutinized, the true network structure remains hidden. NGTA is a lightweight algorithm that enables filtering of the network to edges that are statistically significantly likely to be a part of the true network structure. This likelihood is computed for and attributed to each edge in the network.

$$Q_{ij} = \frac{\rho \alpha^{E_{ij}} (1 - \alpha)^{N - E_{ij}}}{\rho \alpha^{E_{ij}} (1 - \alpha)^{N - E_{ij}} + (1 - \rho) \beta^{E_{ij}} (1 - \beta)^{N - E_{ij}}}$$
$$\alpha = \frac{\sum_{i < j} E_{ij} Q_{ij}}{N \sum_{i < j} Q_{ij}}$$

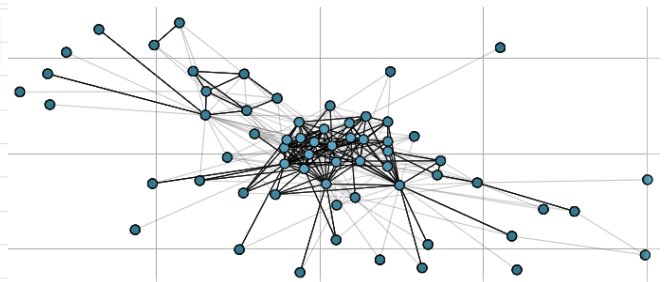
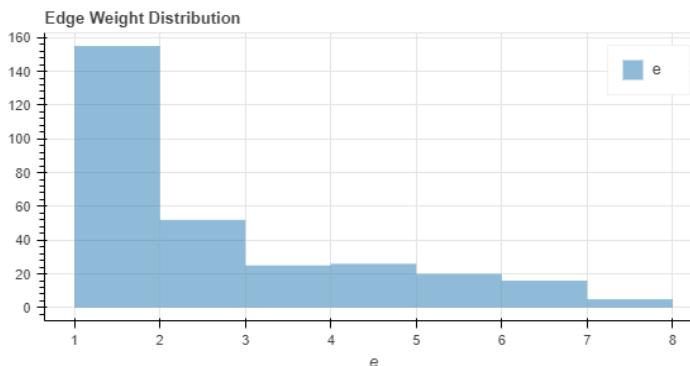
$$\beta = \frac{\sum_{i < j} E_{ij} (1 - Q_{ij})}{N \sum_{i < j} (1 - Q_{ij})}$$

$$\rho = \frac{1}{\binom{n}{2}} \sum_{i < j} Q_{ij}$$

Each edge's likelihood (Q_{ij}) is based upon three network-level statistics: Alpha (True Positive Rate), Beta (False Positive Rate), and Rho (Prior Edge Probability). Each Q_{ij} value also depends on that edge's weight (e) and an overall measurement count (m). Alpha, Beta, and Rho are instantiated randomly. Then, an iteration begins. Within each iteration, each row's Q_{ij} , Alpha, Beta, and Rho are calculated. Iterations continue until Alpha, Beta, and Rho converge to stable values.

Alpha, Beta, and Rho *usually* converge rapidly to stable values. When convergence is not stable, this suggests one of two things: that the network exhibits significant heteroskedasticity across its surface, or that the distribution of the edge weights contains hills and valleys that create local minima and maxima. Fortunately, both of these issues are addressable. Heteroskedasticity is addressed by dividing the network into cohorts, and running NGTA upon each cohort. The patent network contains attribution for each patent, including industry. These attributes will serve as the means to subdivide the network.

If the distribution of edge weights presents hills and valleys, a transformation can be imposed upon that distribution, resulting in smoothing. Should this situation present itself, the specific and most suitable transformation will have to be discovered through trials.



Newman's final network. Edge darkness corresponds to Q_{ij} values. Observe a core of nodes that are very likely connected to each other, and several peripheral nodes far less likely to be a part of the true network. Q_{ij} values range from above 0.999 to below 0.1.

NGTA is innovative in many ways. First, it is lightweight and unsupervised. It can operate on a large network (e.g. the patent network) without incurring large computational expenses. Secondly, the patent network is constantly changing, almost on a daily basis. Organizations tend to tackle noisy networks occasionally, due to the cost of conventional algorithms. NGTA can be re-run cheaply. Lastly, NGTA outputs network-level true positive and false positive rates. This makes NGTA predictive and diagnostic.

Experiments/Evaluation

To measure progress, two main experiments will be performed, a small proof-of-concept, and a larger production version. For the proof-of-concept, we will randomly sample 10% of the nodes and edges from the original set, and perform NGTA on it. If NGTA converges, we will visualize the results using d3 via a graph network, with filtering capabilities based on the values of NGTA. If this scales appropriately, we will then conduct our second production experiment which will use the full dataset. To conduct this experiment, a holdout set of the patent data will be created which contains a uniformly randomly sampled set of all the nodes and edges that is 10% the size of the original data. The rest of the data will be put into the training set. Then, we will run the NGTA algorithm on both the holdout set and the training set. Once NGTA converges, we will display the results on our proof-of-concept network graph visualization.

In order to evaluate the success of NGTA, we will compare the results between the holdout and training set. If NGTA converges on both sets and the alpha and beta statistics calculated by NGTA are within 5% between each set, we can be confident that the results are due to NGTA actually converging and correctly identifying outliers in our graph, and are not caused by random chance. A second test for performance is to test the if NGTA network confidence (γ) vs the number of measurements (x) is not a straight line. This means the algorithm has learned a trend in the network and not random noise.

To evaluate the visualization, we will test for performance and readability. If the graph loads without excessive loading times (greater than 10 seconds), with filtering functionality working properly, the graph will be deemed performant enough for our project. If the graph can display all of the nodes with overlay capability to display the results of the NGTA as well as important meta-data for each node and edge including the patent name, inventors, and company names, the visualization will be deemed a success.

Plan of Activities: 10/14 and 11/1

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 have contributed a similar volume of effort.

TASK TITLE	TASK OWNER	START DATE	DUE DATE	DURATION	% COMPLETE
Data Acquisition and Cloud Storage	Sang Yoon, Erik S	10/15	10/22	7	100%
Data Engineering	Rick, Erik W	10/22	11/2	10	75%
Algorithm Development & Cloud Execution	Tylor	11/2	11/9	7	25%
Front End Development	Sandro, Rick	11/9	11/16	7	0%
Front End Testing	Everyone	11/16	11/23	7	0%

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.