

# Innovation Networks: Review

Alun Meredith



## 1 INTRODUCTION

Patent systems have been studied in many ways, from natural language processing, graph theoretical approaches and exploring the analogy to natural evolution.

The USPTO dataset provides one of the most comprehensive patent databases offering full text documents from 1976 and images since 1790. In this paper we review the contributions of some of the most important works in this field.

There have been a number of papers exploring this dataset but one of the more recent and comprehensive was Bernard Gress (2009) [?]. Here Bernard investigates the ratios citations given and received for patents in different technology groups. High number and diversity of citations given was treated as an indication of generality, which number of citations received is a measure of productivity and originality. He then compared these measures and how they varied over time for different technology categories, assessing some areas as sinks. He primarily concludes that these categories are fundamentally different and therefore future research needs to take this into account.

Beyond bibliographic networks the patent documents themselves have been analysed using natural language techniques by many researchers [?]. Focusing on one, Byungun Yoon [?] conducted a study to produce a network based weighted term similarities of patent documents using standard bag of words methods and comparing these measures to analogies in citation networks. Through inspection they argue that the centrality of their network yields a more relevant approximation of impact because it is less biased by age and preferential

attachment mechanisms.

## 2 GRAPH THEORY

The field of scale-free networks is relatively young with the bulk of research since 1999 where the web was shown to be a scale-free network rather than a random one the field has grown significantly[?]. A defining trait of scale-free networks is to have simple interactions creating emergent behaviour in a way in which it is hard to predict from observing the interactions, because of this simulations and modelling has proven a powerful tool in understanding the links between the two.

The foundation of the models used is Price's model which [?] explains the Mathew effect [?] of the rich getting richer by modelling the attachment of new nodes as a linear function of the number of edges of that node.

Since then there has been a wealth of research using different scale-free networks to build and test adjustments to Price's model with concepts of non-linear growth, local events, growth constraints and initial attractiveness[?].

The patent system has only recently been studied in detail from a graph theoretical perspective. In 2007 Csardi et. al published the first paper examining the USPTO dataset from this perspective [?]. They did this by applying a basic model assuming the attractiveness of a node (rate at which new nodes attach to it) is a function of the age and number of citations of the node. Normalising for the growth of the patent network over time they found the total attractiveness of the system over time could only be replicated with a super-linear preferential attachment model.

They also explore the idea that the increase in number of citations per patent increasing has been coupled with a fundamental change in the structure of the network. Here they find that the level of stratification starts to increase in alignment with the higher citation rates.

Valverde et. al. [?] builds on this work by suggesting an extended power-law form for the attractiveness function and shows with correct parameters this form can describe both the scale-invariant and exponential tail of the citation count distribution. They also explore the clustering and modularity of the system. The clustering coefficient being approximately inversely proportional to the number of citations, suggesting a hierarchical structure.

There is an increasing body of research in academic bibliographic networks which aims to explain more of the mechanics behind citation formation in their graph network models. This includes the propagation of errors which suggests 70 - 80% of citations are copy and pasted from a secondary source [?] or the study of redundant edges to show that the majority of references ( 70%) are secondary [?].

### 3 INNOVATION EVOLUTION

There is a body of research exploring the analogy between the evolution of innovation and biological evolution. The premise is that each invention is built from the recombination of previous inventions. The two have their differences however, there is a limited concept of 'death' in innovation as very old patents may still be cited by new ones and it is hard to think of bibliometric patent networks as direct lineages to name a few.

Yeoun et. al [?] explores this idea of invention as a recombination process by looking at the use of technology codes in patents as a proxy for novelty. Technology codes map the technological niche of a patent consisting of categories and subcategories. Patents can and predominantly do have combinations of technology codes. They show that as the number of patents increases the number of new codes being generated falls off concluding new technologies has a minimal role relative to recombination.

They also show that 40% of patents use existing combinations vs. new ones suggesting these are incremental improvements.

The study of technology codes isn't directly related to a network but they do note that the ageing of codes is different from a bibliometric sense, codes appear not to age with 99% of codes being used at least once every 7 years.

They also look at the dissimilarity of the codes as a proxy for novelty. If a patent is used in a very different field from its parents it is argued that it is more likely to be a bigger leap in novelty. This is done in a binary way, categorising patents as either narrow or broad leaps in novelty and only uses the count as a metric, as such we only get a sense how novelty has changed with time and not any of the network factors which may be present here, such as the distribution of novelty could be a power law or the degree of novelty can be a measure of linkage between clusters like N-K landscapes.

The limitations of the evolutionary analogy are loosely addressed warning that citations in patents aren't directly related to lineage but about carving a legal niche and there being no good metric of fitness for patents.

Buchanan et. al [?] glosses over some of these limitations, using the number of citations a patent has as an "impact" metric, a proxy for fitness. Prior art citations also function as a proxy for combinatorial lineage. They tell the story of the most cited patents in the network over the past 30 years and show that such a distribution of citations cannot come from random natural selection and therefore must be due to adaptive selection in an evolutionary model of innovation. This argument is the random network vs. preferential attachment network argument from an evolutionary perspective.

They focus on showing this idea more robustly incorporating a multitude of normalisation techniques and simulating a null hypothesis random network model by sampling existing data, rather than building a clean model from scratch. Observing the familiar hallmarks of a fat-tailed distribution they conclude that these "superstars" high impact is due to adaptive features, however they do not address the

role of preferential attachment here, how many of the citations received are due to 'rich getting richer' mechanics or due to the intrinsic quality of those innovations.

Their paper also makes investigates the dissimilarity of technology codes as a proxy for novelty making the claim that large leaps in novelty are responsible for the largest "impact" patents, however because its scope is only looking at the 20 most cited patents falls short of being able to make such a general argument about the network as a whole.

Finally Arthur et. al [?] makes the most direct contrast between the search process of a genetic algorithm and the evolution of technology. In their paper they simulate an evolving population of logical circuits starting from simple logic in order to meet a selection of logical needs. Analysing the evolution of the population and resulting network.

They find many of the typical evolutionary features present such as building blocks being formed as intermediary steps to complex solutions i.e. the building block hypothesis, sub-optimal solutions slowly being replaced by better ones and solutions with redundancy simplifying over time.

Complex features are also observed such as a loose power law distribution of edges in the network and avalanches of redundancy as new technologies replace old ones and their dependencies, the size of these redundancies follows a power law showing self-organised criticality.

The paper incorporates a standard genetic algorithm, ignoring many of the observed differences between natural selection and the evolution of innovation, such as holding a finite population, therefore incorporating the "death" of patents and use of random selection shown earlier to not be descriptive of the patent network. Despite this they achieve results similar to observed patent networks, further research could conclude that many of the differences observed naturally arrive from a simple model for example patent ageing could be a result from the saturation of combinatorial space around older technologies.

## 4 CONCLUSION

In conclusion there are three main directions in which Patent Networks are studied, as a branch of already existing Price model networks, through natural language processing techniques and through an analogy to evolutionary processes.

We have seen how efforts have been made to link the evolution of innovations to a Darwinian evolutionary models, that despite some success a lot more effort needs to be done to incorporate the differences between natural and technological evolution into models. linking these more strongly with networks models can be a key to understanding the mechanics of the innovation of evolution.

There have only been a few studies of this dataset from a graph theoretical perspective. There is room for more work in this area especially when paralleling the advances of the academic citation literature to include more subtlety to the models. Although there is still debate as to the efficacy of number of citations as a measure of value or citations as a proxy for lineage.

The greatest challenge to studying patent networks lies in the efficacies of citations as a measure of impact, there has been great debate in the papers referenced here to what extent impact can be measured this way and to what degree these measures can be improved.

## REFERENCES