

Temporality in Link Prediction: Understanding Social Complexity

Abstract

This article describes research into the discovery and modelling of emergent temporal phenomena in social networks. It summarises experimental results that bring together two views in contemporary science: Bayesian analysis and link prediction, to enhance the current understanding of emergent temporal patterns in social network analysis (SNA), particularly in value creation through social connectedness – an important, and growing, discipline within management science.

Traditional link prediction methods use the values of metrics in a graph to determine where new links are likely to arise, and little work has been done on analysing long-term graph trends. We have found that existing graph generation models are unrealistic in their prediction, and can be complemented through the use of temporal metrics, in the study of some networks. To date, no temporal information has been used in link prediction research, thereby excluding valuable temporal trends that emerge in sociogram sequences and also lowering the accuracy of the link prediction. We extracted information from the *Pussokram* online dating network dataset, and 9,939 cases of each class were formed. Logistic regression in the Weka data mining system was used to perform link prediction. Our results show that temporal metrics are an extremely valuable new contribution to link prediction, and should be used in future applications.

In addition to using metrics to measure the local behaviours of participants in social networks, we used Bayesian networks to model the interrelationships between the metrics as local behaviours and links forming between individuals as emergent behaviours (social complexity). We also explored how the metrics evolve over time using Dynamic Bayesian Networks (DBN).

Keywords

Social networks, complex adaptive systems, link prediction, Bayesian analysis, temporal analysis, emergence.

Introduction

Social networks are complex systems that are characterised by high numbers of interconnected component entities, and a high degree of interaction between these entities. The interrelationships in such a network are dynamic and evolve over time. Temporal changes in social networks are difficult to understand and anticipate. The interrelationships between the component entities in a social network and its global behaviour can be so numerous and mostly hidden, and can affect so many different entities throughout the social network that it becomes extremely difficult to comprehend.

Complexity theory is ideally suited to study social networks. Complex adaptive systems theory is a branch of complexity theory that studies systems that consist of agents that are collectively able to evolve in response to environmental changes. The agents in such a system constantly act and react to the actions of other agents and events in the environment. A social network is a complex adaptive system, in which people are agents interacting with each other.

In order to understand social complexity, the local behaviours of the participants must be understood, as well as how they act together and interact with the environment to form the whole. To model this, we use Bayesian network techniques to mine and model emergent relationships between local behaviours and global behaviours in social networks over time.

Social Complexity

The complex structure of any social organisation can be thought of as a network of individuals/agents (Nohria & Eccles, 1992: 288; Lincoln, 1982), sometimes termed network actors, that operates and is operated on in an environment which itself is an environment of other distributed organisations (Van Wijk, Van Den Bosch & Volberda, 2003; Potgieter, April, Cooke & Lockett, 2006), and the actions of

agents within the network are shaped and constrained because of their position and embeddedness in the network (Nohria, 1992).

Complex adaptive systems theory, a branch of complexity theory that is well suited to study social networks, investigates systems that consist of agents that are collectively able to evolve in response to environmental changes (i.e., people who interact with each other in socially complex ways). Social networks are complex systems that are characterised by high numbers of dynamically interconnected component entities, and a high degree of time-based evolutionary interactions between these entities. Temporal changes in social networks are difficult to understand and anticipate, because of the structures of constraint and opportunity negotiated and reinforced between interacting individuals. The interrelationships between the component entities in a social network and its global behaviour can be so numerous and mostly hidden, and can affect so many different entities throughout the social network that it becomes extremely difficult to comprehend. Such changes, for modern organisations, can come in the form of globalisation, deregulation, competition, technology and political transformation. The agents in such a constantly changing system tend to step out of the traditional boundaries of their network (within-group identification), act and react to, and evolve in relation to, the actions of other agents and events in the environment to achieve a more desirable end, e.g., to learn faster in the changing environment, or to collaborate better, or to craft more relevant identities, or to seek help in making major decisions, or to negotiate and renegotiate relationships of power, etc.

In social network studies, the traditional archetype – which acts as an interpretive schema – of an organisation is that it is made up of a number of static departments, which themselves have specialised individuals (agents), and its strategy for competing, or winning, or succeeding, or sustaining itself is fixed for a time period, usually 1-3 years, and is clearly understood and acted upon by its agents. After the specified period, usually the organisational strategy is reviewed, altered, or rewritten or totally revoked in favour of a new strategy, and its agents then enact new behaviours against such a plan. We contend, however, that organisations are not static, atomistic agents, instead they are recurring in dynamic agent linkages and embedded in networks of increasing importance (Brock, 2006) that progressively influence competitive actions (Granovetter, 1985 & 1992; Burt, 1995) and may themselves challenge the interpretive schema and supportive structures, and ultimately even delegitimise the existing archetype.

Social Network Analysis and Link Prediction

Social Network Analysis (SNA) is a research area aimed at understanding social complexity by representing and analysing social networks using mathematical graphs. This was first done by researchers, such as Cartwright & Harary (1956) in the USA, who interpreted Kurt Lewin's social interaction theory into graph theory – thereby helping to transform the study of social networks from description to analysis. In this theory, a graph G is a structure consisting of a set of nodes V (also called vertices), and a set of links E (also called arcs or edges). In our discussion the terms graph, network, social network and sociogram are used interchangeably. A graph can be bidirected¹, meaning that we are not concerned about the order of the nodes in a link. Alternatively, a graph is directed, meaning that the nodes in a link form an ordered pair. Nodes represent people and links represent some type of relationship between people. For the purposes of this discussion a link exists between two people if they have exchanged a message in the past. Thus links are never removed once they have been formed.

There appears to be one standard textbook on SNA (Wassermann & Faust, 1994), cited by the large majority of researchers in this field. Although SNA has existed for over fifty years, most analysis

¹ Also called “undirected”, but bidirected is used in this paper to emphasise that links are directed and go in both directions.

techniques have been designed for static data, or at least a static archetype of social organisation. For example, Wassermann & Faust (1994) contains no mention of temporal metrics, even though it was written in 1994 when electronic networks were well established. With the increase in the use of computers, collecting enough data to create numerous graphs over fixed time intervals becomes possible. An example is creating a graph per week from email data, using a server's email log of 'to', 'from', and 'date' fields (Campbell, Maglio, Cozzi & Dom, 2003). This series of graphs can be used to study the evolution of the network and the change over time in various metrics. Predicting certain changes to a social network is called the link prediction problem. Liben-Nowell & Kleinberg (2003) explain it as:

Given a snapshot of a social network at time t , we seek to accurately predict the edges that will be added to the network during the interval from time t to a given future time t' .

Link prediction has many real world applications, especially in the fields of marketing and crime prevention. Examples include:

- Identifying the structure of a criminal network (i.e., predicting missing links in a criminal network using incomplete data).
- Overcoming the data-sparsity problem in recommender systems using collaborative filtering (Huang, Li & Chen, 2005).
- Accelerating a mutually beneficial professional- or academic connection that would have taken longer to form serendipitously (Farrell, Campbell & Myagmar, 2005).
- Improving hypertext analysis for information retrieval and search engines (Henziger, 2000).
- Monitoring and controlling computer viruses that use email as a vector (Lim, Negnevitsky & Hartnett, 2005).
- Predicting when webpage users will next visit, in order to improve the efficiency and effectiveness of a site's navigation (Zhu, 2003).
- Link prediction might also be useful in ecology, though interdisciplinary sharing between these two fields is still new (McMahon, Miller & Drake, 2001).

Metrics

Link prediction through topological analysis is performed by computing various metrics. A metric is a value calculated from a graph that describes the graph in some way. For instance, it is more likely that two nodes that both have a high degree (number of neighbours) are more likely to form a new link, than two nodes with a low degree (Liben-Nowell, 2005). Most traditional SNA metrics are described and defined in Wassermann & Faust (1994), and summarised in an online book by Hanneman (2001). We defined metrics to quantify local behaviours of agent resources in social resource combinations (Potgieter, April, Cooke & Lockett, 2006). The appearance of new links between individual agent resources is emergent behaviour of the social network. Table 1 lists the metrics that were used for link prediction in this research. They were chosen since prior research has found them to be useful. Huang, et al. (2005) found that the most useful dyadic metrics for link prediction, in descending order, were the Katz measure, preferential attachment, common neighbours and the Adamic\Adar number. Recency is a metric that, to the authors' knowledge, has not been used before. The table uses certain symbols that are now defined:

denotes the number of elements in a set. U_n is the set of bidirected links of the social network at time step n. A bi-directed link between nodes v_i and v_j is notated as $u_{i,j}$ and $u_{j,i}$. $\Gamma(v_i)$ is the set of neighbours of node v_i , i.e. the set $\{v_j : u_{i,j} \in U\}$.

$P(v_i, v_j)$ is the set of all shortest paths from v_i to v_j . $P(v_i, v_j, v_x)$ is the set of all shortest paths from v_i to v_j that pass through node v_x .

Name	Definition	Description
Degree	$\#\{u_{j,i} : u_{i,j} \in U_n\}$ <p style="text-align: center;">or</p> $\#\Gamma(v_i)$	The number of links from v_i to any node at time step n.
Recency	N/A	One plus the number of time steps elapsed since the node last communicated. Has a value of one if the node sent or received a message in the current time step.
Betweenness (Wassermann & Faust, 1994)	$\sum_{v_j \in V} \sum_{v_k \in V, v_k \neq v_j} \frac{\#P(v_j, v_k, v_i)}{\#P(v_j, v_k)}$	The sum of all shortest paths between all nodes that contain v_i as a percentage of all shortest paths between all nodes. It ranges from 0 to $\frac{(\#V - 1)(\#V - 2)}{2}$
Common neighbours (Huang <i>et al.</i> , 2005)	$\#\{v_k : u_{i,k} \in U_n, u_{k,j} \in U_n\}$ <p style="text-align: center;">or</p> $\#\{\Gamma(v_i) \cap \Gamma(v_j)\}$	The number of nodes linked to both focus nodes (i.e. mutual friends).
Adamic\ Adar similarity (Adamic & Adar, 2003)	<p style="text-align: center;">General case:</p> $\sum_{z: \text{shared feature}} \frac{1}{\log(\text{frequency}(z))}$ <p style="text-align: center;">Common neighbours case:</p> $\sum_{v_z \in \Gamma(v_i) \cap \Gamma(v_j)} \frac{1}{\log(\#\Gamma(v_z))}$	The number of features shared by the nodes, divided by the log of the frequency of the features. This metric rates rarer features more heavily.

Preferential attachment (Huang <i>et al.</i> , 2005)	$\#\{v_k : u_{i,k} \in U_n\} \cdot \#\{v_k : u_{j,k} \in U_n\}$ <p style="text-align: center;">or</p> $\#\Gamma(v_i) \cdot \#\Gamma(v_j)$	The product of the number of edges incident to the two nodes.
Katz measure (Liben-Nowell, 2005)	$\sum_{l=1}^{\infty} \beta^l \cdot \#(\text{paths}_{v_i, v_j}^{<l>})$, where $\text{paths}_{v_i, v_j}^{<l>}$ is the set of all paths of length l from v_i to v_j	The sum of all paths between the nodes exponentially damped by length to weight short paths more heavily. $0 < \beta < 1$.

Table 1. Metrics used in this link prediction research

Existing link prediction techniques use the values of metrics in a graph to determine where new links are likely to arise. Important contributions to the field include Popescul & Ungar (2003), Taskar, Wong, Abbeel & Koller (2004), Popescul & Ungar (2004) and Zhou & Scholkopf (2004). The classic paper on link prediction is by Liben-Nowell & Kleinberg (2003). They tested the predictive power of only proximity metrics, including common neighbours, the Katz measure and variants of PageRank. They found some of these measures had a predictive accuracy of up to 16% (compared to a random prediction's accuracy of less than a percent). A third of Liben-Nowell's doctoral thesis (Liben-Nowell, 2005) was a chapter on link prediction in social networks. A few other link prediction papers are summarised in Getoor & Diehl (2005).

Past Contributions to Temporal Analysis

Leskovec, Kleinberg & Faloutsos (2005) state that little work has been done on analysing long-term graph trends:

Many studies have discovered patterns in static graphs, identifying properties in a single snapshot of a large network, or in a very small number of snapshots; these include heavy tails for in- and out-degree distributions, communities, small-world phenomena, and others. However, given the lack of information about network evolution over long periods, it has been hard to convert these findings into statements about trends over time.

Their study of trends found that, over time, graphs densify and the average distance between nodes decreases. This was contrary to the existing beliefs that the average nodal degree remains constant and average distance slowly increases. They claimed that existing graph generation models are not realistic and proposed a new "forest-fire" generation model. Desikan & Srivastava (2004) studied the change in metrics of a set of webpages over time for the graph as a whole, and for single nodes (subgraphs are their current research). They found that temporal metrics, such as their Page Usage Popularity, can be effectively used to boost ranks of recently popular pages to those that are more obsolete. This seems to indicate that temporal metrics are a useful addition to traditional static metrics in the study of some networks.

Using Temporal Analysis to Aid Link Prediction

To date no temporal information has been used in link prediction research (excluding our earlier proposition of the idea in April, Potgieter & Cooke (2005) and Potgieter, April, Cooke & Lockett (2006). It is the authors' opinion that ignoring such information excludes valuable temporal trends that emerge in sociogram sequences that may greatly increase the accuracy of link prediction. We have three types of temporal metrics that have been defined, the first two of which are borrowed from finance:

- Return is the percentage increase or decrease of a value over a period of time (Ross, Westerfield, Jordan & Firer, 2001). It shows the rate of change of a given metric. For instance, if we are considering the degree of node v_i from time step one to time step fifty, the degree return would be
$$\frac{\text{degree}_{50}(v_i) - \text{degree}_1(v_i)}{\text{degree}_1(v_i)}$$
.
- Moving averages are used to extract long-term trends from short-term noise. They do not show trends or "movement", like return, but rather serve only to blur the values of the metrics around a point.
- Recency (Table 1), simply shows how much time has elapsed since a node has communicated.

Temporal Analysis Methodology

All the metrics listed in the Table 1 were calculated for a series of hundred networks constructed daily from information extracted from the *Pussokram* online dating network dataset. In addition to the computation standard metrics, their moving averages over twenty, ten and two days were computed, as well as their returns over the same period and the recency of each node. Up to one hundred sample dyads where new links formed on the day, and an equivalent number of unconnected dyads, were chosen for each network. There were 9,939 cases of each class in total. Logistic regression (Hosmer & Lemeshow, 1989) in the Weka data mining system (Witten & Frank, 2005) was used to perform link prediction. The dataset was split into a 70%/30% train/test division, and after the system was trained it predicted whether a given set of test instance metrics characterised a forming link or an unconnected link. The overall classification accuracy, the true positive rates for each class, and the kappa statistic (a measure of a system's accuracy improvement over a random guess) were calculated (Witten & Frank, 2005). The true positive rate of a class is the percentage of all the instances of that class that were predicted correctly. The kappa statistic is defined as:

$$\frac{P(A) - P(E)}{1 - P(E)}, \text{ where } P(A) \text{ is the total percentage accuracy of instances predicted by the learning system}$$

and $P(E)$ is the total percentage accuracy of instances predicted by random guessing. Since in this research an equal number of positive and negative instances were used we know that $P(E) = 0.5$ and $1 - P(E) = 0.5$, which we applied to the kappa statistic. A kappa value over 0.4 has been said to indicate 'good agreement beyond chance' (Fleiss, 1981).

Temporal Analysis Results and Conclusions

Table 2 shows the mean and standard deviation for each class, as well as the significance level of the mean difference using a Normal statistical difference of means test. The last column is the most important and shows the kappa value of each metric when used alone in a logistic regression against the class. The suffix 'F' stands for 'From' (the first node in a dyad instance), the suffix 'T' stands for 'To' (the second

node in a dyad instance). ‘A’ stands for ‘Average’, ‘R’ stands for ‘return’ and the numbers ‘20’, ‘10’ and ‘2’ stand for the number of time steps over which the average and return were calculated. The metric in the middle row of each group is the traditional static metric. The averages are shown above it and the returns are shown below. ‘Forming’ and ‘Unconnected’ are emergent properties of the social network. Forming means that a link does not exist in the current time step, but it will appear in the next time step. Unconnected means that link between two nodes does not exist in the current time step, and these nodes will remain unconnected in the next time step.

Metric	Unconnected mean	Forming mean	Unconnected standard deviation	Forming standard deviation	Mean Difference	Test statistic	Significance level	Kappa
DegreeFA20	1.9438	8.6818	3.5317	21.5584	-6.7379	-30.75	0.0001	35.48%
DegreeFA10	1.9152	8.3771	3.5181	22.4411	-6.4618	-28.36	0.0001	35.12%
DegreeFA2	1.8927	7.9833	3.5052	22.1150	-6.0906	-27.12	0.0001	29.85%
DegreeF	1.8876	7.8111	3.5006	21.9536	-5.9235	-26.56	0.0001	27.13%
DegreeFR20	0.1361	0.7550	0.5419	2.2471	-0.6189	-26.69	0.0001	10.26%
DegreeFR10	0.0626	0.4599	0.3011	1.8352	-0.3973	-21.3	0.0001	5.86%
DegreeFR2	0.0089	0.0505	0.0952	0.2291	-0.0417	-16.74	0.0001	-8.33%
DegreeTA20	1.9810	8.0030	3.4860	19.1849	-6.0220	-30.79	0.0001	42.46%
DegreeTA10	1.9606	8.0086	3.5178	20.3095	-6.0480	-29.25	0.0001	41.54%
DegreeTA2	1.9256	7.5180	3.5184	20.7846	-5.5923	-26.45	0.0001	32.30%
DegreeT	1.9223	7.0698	3.5223	20.2884	-5.1475	-24.92	0.0001	26.47%
DegreeTR20	0.1426	0.6271	0.5724	1.6312	-0.4844	-27.94	0.0001	15.68%
DegreeTR10	0.0658	0.4169	0.2848	1.5329	-0.3510	-22.45	0.0001	12.97%
DegreeTR2	0.0066	0.0795	0.0734	0.3285	-0.0729	-21.59	0.0001	-2.55%
KatzA20	0.0111	0.0427	0.0225	0.0473	-0.0316	-60.08	0.0001	41.78%
KatzA10	0.0110	0.0415	0.0229	0.0475	-0.0305	-57.7	0.0001	44.92%
KatzA2	0.0107	0.0373	0.0227	0.0471	-0.0267	-50.8	0.0001	37.32%
Katz	0.0106	0.0348	0.0227	0.0467	-0.0241	-46.31	0.0001	28.17%
KatzR20	0.2177	0.5243	0.6967	1.4813	-0.3065	-18.67	0.0001	-2%

Metric	Unconnected mean	Forming mean	Unconnected standard deviation	Forming standard deviation	Mean Difference	Test statistic	Significance level	Kappa
KatzR10	0.1344	0.3634	0.6071	1.2972	-0.2290	-15.94	0.0001	-7%
KatzR2	0.0094	0.0846	0.1194	0.6620	-0.0752	-11.15	0.0001	-23.1%
PAA20	3.8428	38.6423	17.7252	97.9615	-34.7995	-34.85	0.0001	45.46%
PAA10	3.7283	37.9504	17.0146	106.8425	-34.2221	-31.54	0.0001	48.36%
PAA2	3.6188	35.7246	16.5470	117.9152	-32.1059	-26.88	0.0001	41.61%
PA	3.6043	33.2047	16.4837	114.0258	-29.6005	-25.61	0.0001	33.27%
PAR20	0.2813	1.6457	0.8481	4.4977	-1.3645	-29.72	0.0001	5.31%
PAR10	0.1404	1.0239	0.5073	3.2094	-0.8835	-27.11	0.0001	1.36%
PAR2	0.0151	0.1359	0.1213	0.4344	-0.1209	-26.72	0.0001	-10.6%
CNA20	0.0015	0.0306	0.0477	0.2007	-0.0290	-14.02	0.0001	32.2%
CNA10	0.0019	0.0302	0.0507	0.2012	-0.0283	-13.58	0.0001	29.17%
CNA2	0.0019	0.0253	0.0497	0.1865	-0.0234	-12.1	0.0001	12.52%
CN	0.0019	0.0231	0.0501	0.1791	-0.0212	-11.38	0.0001	1.92%
CNR20	0.0000	0.0920	0.0000	0.3282	-0.0920	-27.93	0.0001	-0.87%
CNR10	0.0000	0.0794	0.0000	0.2639	-0.0794	-29.98	0.0001	-1.18%
CNR2	0.0000	0.0112	0.0000	0.1054	-0.0112	-10.57	0.0001	-0.00%
AAA20	0.0015	0.0442	0.0656	0.3151	-0.0427	-13.22	0.0001	32.2%
AAA10	0.0019	0.0417	0.0656	0.2989	-0.0399	-12.99	0.0001	29.17%
AAA2	0.0017	0.0357	0.0619	0.2785	-0.0340	-11.87	0.0001	12.52%
AA	0.0017	0.0329	0.0615	0.2701	-0.0311	-11.21	0.0001	1.92%
AAR20	-0.0393	-0.0091	0.0481	0.2333	-0.0301	-12.62	0.0001	-0.40%
AAR10	-0.0163	0.0539	0.0376	0.2920	-0.0702	-23.78	0.0001	-1.18%
AAR2	-0.0092	0.0066	0.0346	0.1109	-0.0158	-13.59	0.0001	0.07%

Metric	Unconnected mean	Forming mean	Unconnected standard deviation	Forming standard deviation	Mean Difference	Test statistic	Significance level	Kappa
RecFA20	38.1636	12.5866	30.5582	16.1022	25.5771	73.82	0.0001	28.38%
RecFA10	38.2149	11.9455	32.0390	16.3050	26.2694	72.85	0.0001	35.75%
RecFA2	38.7039	11.0723	33.3012	16.5778	27.6316	74.05	0.0001	44.52%
RecF	38.7928	10.8518	33.4656	16.6218	27.9410	74.55	0.0001	46.34%
RecFR20	1.5342	1.4653	3.1617	3.5274	0.0688	1.45	NS	-2.92%
RecFR10	0.7021	0.8002	1.5502	1.9719	-0.0981	-3.9	0.0001	1.82%
RecFR2	0.0814	0.1413	0.2480	0.4357	-0.0599	-11.91	0.0001	8.46%
RecTA20	37.3356	15.1974	30.3941	18.7683	22.1382	61.78	0.0001	14.80%
RecTA10	37.5869	14.0277	31.8041	18.7829	23.5592	63.59	0.0001	22.61%
RecTA2	37.7987	11.8050	33.1000	18.2863	25.9937	68.53	0.0001	37.67%
RecT	37.8514	11.1834	33.2870	18.0862	26.6681	70.18	0.0001	42.41%
RecTR20	1.5430	1.5128	3.1033	3.5408	0.0302	0.64	NS	0.00%
RecTR10	0.6878	0.8434	1.5252	2.0300	-0.1556	-6.11	0.0001	8.01%
RecTR2	0.0797	0.1669	0.2564	0.4539	-0.0872	-16.67	0.0001	12.87%

Table 2. Metric statistics, grouped by category

We can see that almost all metrics had significant differences between each class. Thus the significance levels become meaningless and we consider only the kappa statistic as a measure of worth. A kappa statistic of 40% corresponds roughly to an overall accuracy of 70%, and a kappa of 20% corresponds to an accuracy of 60%. We note that a metric's moving average is a better link predictor than the static metric (except in the case of recency, which is already a temporal metric). Furthermore, the increase in predictive accuracy of moving averages seems to level off somewhere between ten and twenty time steps prior to the current time step. Thus, it is recommended that taking averages further back than twenty time steps is unnecessary and extravagant. Finally, metric returns appear to be completely useless for link prediction and should not be used in the future.

Table 3 shows how the metrics perform in combination to predict links. This is the ultimate test of prediction accuracy – seeing how well we can classify links as forming or not (global behaviours) using any combination of metrics (local behaviours) at our disposal.

Metric subset	Kappa	Unconnected TP rate	Forming TP rate	Overall accuracy
<i>Static metrics:</i> DegreeF, DegreeT, Katz, PA, CN, AA	39.83%	79.8%	59.9%	69.9698%
<i>Static metrics with average 10:</i> DegreeF, DegreeT, Katz, PA, CN, AA, DegreeFA10, DegreeTA10, KatzA10, PAA10, CNA10, AAA10	51.13%	80.4%	70.1%	75.5869%
<i>Static metrics with average 10 and recency:</i> DegreeF, DegreeT, Katz, PA, CN, AA, RecF, RecT, DegreeFA10, DegreeTA10, KatzA10, PAA10, CNA10, AAA10	63.62%	81.3%	82.4%	81.8075%

Table 3. Metric sets prediction

The first row shows the accuracy of static metrics alone. The second row shows the accuracy of the same metrics using moving averages. It can be seen that the overall accuracy increases by 6% from the first row to the second. If we include the recency metrics as well the accuracy increases to 82%, with true positive rates above 80% for both classes. This shows that temporal metrics are an extremely valuable new contribution to link prediction and should be used in future applications. Further research to be performed includes suggesting new temporal variations of static metrics and determining exactly the optimum number of time steps over which to compute temporal metrics.

Using Bayesian Networks to Understand Social Complexity

In order to understand social complexity, the local behaviours of the participants/network actors must be understood, as well as how they act together and interact with the environment to form the whole. To model this, we used Bayesian network techniques.

According to Baas and Emmeche (1997), understanding is related to the notion of explanation. A complex adaptive system uses the hyperstructures in its internal model for explanation and understanding. It uses observation mechanisms to create and maintain these hyperstructures. The process of adaptation

relies heavily on the observation mechanisms and involves a progressive modification of the hyperstructures (Holland, 1995).

We measured local behaviours using metrics and we used Bayesian networks to model the interrelationships between the metrics as local behaviours and links forming between individuals as emergent behaviours. We also explored how the metrics evolve over time using Dynamic Bayesian Networks.

Link Prediction using a Bayesian Network

Bayesian Networks provide the ideal technology to reason about social resource combinations (Potgieter, April, Cooke & Lockett, 2006). A Bayesian network is a directed acyclic graph (DAG) that consists of a set of nodes that are linked together by directional links. Each node represents a random variable or uncertain quantity. Each variable has a finite set of mutually exclusive propositions, called *states*. The links represent informational or causal dependencies among the variables, where a parent node is the *cause* and a child node the *effect*. The dependencies are given in terms of conditional probabilities of states that a node can have, given the values of the parent nodes (Pearl, 1988). Each node has a conditional probability matrix to store these conditional probabilities, accumulated over time.

Bayesian learning can be described as ‘mining’ the structure of the network and calculating the conditional probability matrices from history data. The data may be incomplete and the structure of the Bayesian network can be unknown.

Figure 2 illustrates the Bayesian network that was mined from the same dataset used in the temporal analysis. It clearly illustrates the cause-effect relationships between the local behaviours (metrics ‘PA’, ‘DegreeT’, ‘DegreeF’, ‘Katz’, ‘AA’ and ‘CN’) and the emergent behaviour ‘Forming’.

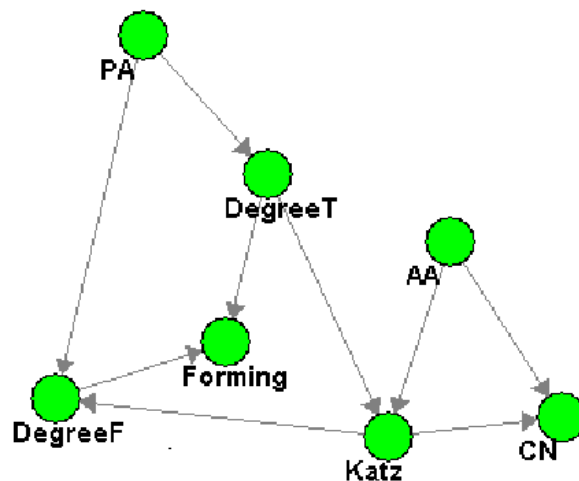


Figure 2. Link Prediction Bayesian Network

With a cross validation method, we extracted 30% records at random from the entire dataset of over 19,000 instances as used in the temporal analysis. The test set contained over 8,000 records. The remaining 70% (over 11,000) records served as the training set. As one of the pre-processing steps, we discretised the training and the test sets using the equal width algorithm (Osunmakinde, 2006). The pre-

processed training set was used by the Hybrid Genetic Algorithm (Osunmakinde, 2006) to mine a Bayesian Network model, shown in Figure 2. This algorithm is described later in this paper.

After training, we used *Bayesian Inference* in the model in Figure 2 to classify the 30% test set. Bayesian inference is the process of calculating the posterior probability of a hypothesis H (involving a set of query variables) given some observed event (assignments of values to a set of evidence variables e), $P(H|e)$.

For the link prediction, the hypothesis will be that ‘Forming’ is true or false, given the values of the evidence variables ‘DegreeT’, ‘DegreeF’, ‘Katz’, ‘PA’, ‘AA’, ‘CN’:

$$P(\text{Forming} = \text{true} \mid \text{DegreeT}, \text{DegreeF}, \text{Katz}, \text{PA}, \text{AA}, \text{CN}) ?$$

Or

$$P(\text{Forming} = \text{false} \mid \text{DegreeT}, \text{DegreeF}, \text{Katz}, \text{PA}, \text{AA}, \text{CN}) ?$$

The confusion matrix (Ron & Foster, 1998) in *Table 4* gives the link prediction and evaluation of the results.

Actual values	Predicted values	
	true	False
Forming = false (Total 4390)	4390	0
Forming = true (Total 4390)	39	4351

Table 4: Confusion Matrix

The implementation results were computed and generated from *Table 4* and presented in *Table 5*.

True positive rate of forming	99.1116173120729 %
True positive rate of not forming (unconnected)	100.0 %
Kappa	99.1116173120729 %
The overall accuracy	99.55580865603645 %

Table 5: Implementation results

Modelling Emergent Social Network Behaviour over time using a DBN

A Dynamic Bayesian network (DBN) is ideally suited to model changes in metrics and emergent behaviour over time. In Dynamic Bayesian networks, multiple copies of the variables are represented, one for each time step (Pearl & Russell, 2000). A DBN provides a compact representation of temporal probabilistic processes such as Hidden Markov Models (HMM). It contains slices of sub-models which facilitate its probabilistic reasoning over time. That is, we can interpret the present, understand the past and predict the future. We want to compute: $P(X_i(t) \mid E(t_1:t_2))$ from the complete joint distribution of variables over one or more time slices as defined as:

$$Pr(X_0, X_1, \dots, X_t, E_1, \dots, E_t) = P(X_0) \prod_{i=1}^n Pr(X_i | X_{i-1}) Pr(E_i | X_i)$$

where $(X_i(t))$ represents the i 'th hidden or unobservable variable at time t and $E(t_1:t_2)$ refers to the observable evidence from time step t_1 to t_2 . Observe that this follows a Markov assumption where the current state depends on only finite history of previous states (Russell & Norvig, 2002).

The construction of the DBN requires the following information components:

- 1) the prior distribution over the state variables, $Pr(X_0)$;
- 2) the transition model, $Pr(X_{t+1} | X_t)$; and
- 3) the sensor model, $Pr(E_t | X_t)$.

The learning of these components depends on the topological connections between the successive slices, the hidden states and the evidence variables. There are many algorithms for learning parameters in DBNs. Where sufficient training data is available, a maximum likelihood estimate (MLE) can be used. Many areas of science and engineering use approximate inference algorithms in DBNs because of the shortcomings of exact inference (Pavlovic, Rehg, Cham & Murphy, 1999).

In this paper, we incorporate a Dynamic Bayesian Network illustrated in *Figure 3* which models possible relationships between metrics and forming links. The observable nodes are 'Forming', 'DegreeT', 'DegreeF', 'Katz', 'PA', 'CN' and 'AA' which serve as the evidence variables 'E', to the unobservable variable called 'HiddenNode'. We tested this DBN using the 150 time slices of the dataset used in the temporal analysis, using the metrics defined in the previous sections. The primary objective of this DBN is to predict relationships between metrics (local behaviours) and links forming (global behaviours) in future time steps. Using this, possible links can be predicted between two or more people using the predicted relationships from observed nodes.

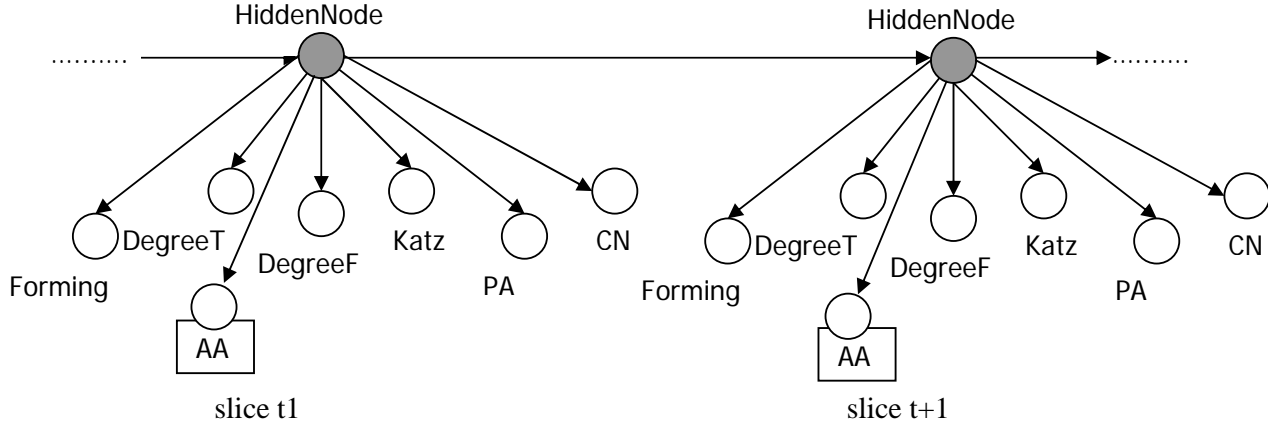


Figure 3: *Dynamic Bayesian Network*

Dynamic Bayesian Network Methodology

We used the same dataset as in the temporal metrics analysis, consisting of 150 time steps and used 70% as a training set, while 30% was used as a test set by a hybrid genetic algorithm (see next section). The average number of records in each time step was 201.

In our implementation, we trained the three components of DBN using the MLE (Russell & Norvig, 2002) and Mutual Information in information theory (Cheng, Bell & Liu, 2001). These 3 components were used during DBN inference using the Viterbi algorithm (see next section). We evaluated the relationships predicted by the DBN using the hybrid genetic algorithm. The DBN did not have access to the 30% dataset, but it was used by the genetic algorithm. That is, the genetic algorithm mined the *actual relationships* between the metrics and forming links at each time step, while the DBN *predicted relationships* between the metrics and forming links at each time step.

Dynamic Bayesian Network Prediction using the Viterbi Approximation

By definition (Pavlovic, Rehg, Cham & Murphy, 1999) the probability of the partial optimal path to a hidden state i at time t when evidence z_t is observed implies:

$$\delta_t(i) = \max_j (\delta_{t-1}(j) a_{ji} b_{iz_t})$$

Where ‘delta’ is the probability of reaching a particular intermediate hidden state on the DBN time slices, ‘ a ’ is the transition probability while ‘ b ’ is the corresponding observation probability. Recall the Markov assumption that says that the probability of a hidden state occurring, given the previous history, depends on the previous i -states. We wanted to compute the probability of the most probable path to any hidden node X . That is:

$$\Pr(X \text{ at time } t) = \max_i \Pr(i \text{ at time } (t-1)) * \Pr(X | a) * \Pr(\text{observation at time } t | X)$$

We used the knowledge of the previous states; the transition probabilities multiplied by the corresponding observation probabilities. This calculated the hidden node.

The HGA is a variant of classical genetic algorithm operators because it integrates information theoretic measures as learning components and mathematical concepts as population construction. Specifically, some of the information theoretic measures we used were: Mutual Information (MI) model, Shannon's information content, and Minimum Description Length (MDL) model. Also, the mathematical component part of the HGA was a PowerSet lattice which implemented crossover operator. Moreover, PowerSet lattice generated population space from which the best networks could emerge.

One of the interesting features of this methodology is that the components of the HGA were implemented as decomposable systems. The decomposability means the functionalities of every component such as MI perceives its inputs and produces its outputs differently from other HGA components. However, such decomposability of every HGA component is to level complexity into simplicity. More information can be found in Osunmakinde (2006).

Examples of the *predicted relationships* by the Dynamic Bayesian Network for time step 1 were as follows:

Probability of Relationship between: Forming and CN, at time step: 1 = 0.032570422535211266
Probability of Relationship between: Katz and PA, at time step: 1 = 7.435374363700479E-6
Probability of Relationship between: Katz and Forming, at time step: 1 = 0.002758431904558663
Probability of Relationship between: Katz and AA, at time step: 1 = 0.0016477018677722897
Probability of Relationship between: Forming and AA, at time step: 1 = 0.013710911301140175
Probability of Relationship between: DegreeF and Forming, at time step: 1 = 7.958115165232352E-4
Probability of Relationship between: Katz and CN, at time step: 1 = 0.002924745442351074

Table 6: Relationships predicted by DBN

We compared the *relationships* predicted by the DBN with the *actual relationships* mined using the hybrid genetic algorithm. We calculated the accuracy of each prediction by counting the number of predicted relationships and actual relationships that correspond in structure predicted by the DBN and the actual structure mined by the HGA. We obtained the following evaluation results:

<p>Time Step 0: Nr of relationships <i>predicted</i> by DBN: 21 Nr of <i>actual</i> relationships mined by HGA: 19.0 Accuracy: 90.47619047619048 %</p>
<p>Time Step 1: Nr of relationships <i>predicted</i> by DBN: 21 Nr of <i>actual</i> relationships mined by HGA: 10.0 Accuracy: 47.61904761904761 %</p>
<p>Time Step 2: Nr of relationships <i>predicted</i> by DBN: 21 Nr of <i>actual</i> relationships mined by HGA: 21.0 , Accuracy: 100.0 %</p>
<p>Time Step 3: Nr of relationships <i>predicted</i> by DBN: 21 Nr of <i>actual</i> relationships mined by HGA: 13.0 Accuracy: 61.904761904761905 %</p>
<p>Time Step 4: Nr of relationships <i>predicted</i> by DBN: 21 Nr of <i>actual</i> relationships mined by HGA: 21.0 Accuracy: 100.0 %</p>
<p>Time Step 5: Nr of relationships <i>predicted</i> by DBN: 21 Nr of <i>actual</i> relationships mined by HGA: 21.0 Accuracy: 100.0 %</p>
<p>Time Step 6: Nr of relationships <i>predicted</i> by DBN: 21 Nr of <i>actual</i> relationships mined by HGA: 20.0 Accuracy: 95.23809523809523 %</p>
<p>Time Step 7: Nr of relationships <i>predicted</i> by DBN: 21 Nr of <i>actual</i> relationships mined by HGA: 19.0 Accuracy: 90.47619047619048 %</p>
<p>” ” ” Overall DBN Prediction Accuracy: 79.52380952380953 %</p>

Table 7: DBN Prediction Accuracy

In future research, we will use these predicted relationships to study to see if time graphs densify over time and at what rate the average distance between nodes decreases.

Conclusion

In this paper, we recognise the role of complex adaptive systems theory in shedding light on an increasingly important aspect of organisational life (and survival), namely, social networks, and emergent temporal phenomena in these networks. It is our contention that the fluidity and temporality of the interrelationships, making up such networks, render it dependent on changing contexts – and therefore accurate prediction, going forward, becomes difficult (near impossible) using traditionally static methodologies.

The use of link prediction to improve collaborative filtering in recommender systems was investigated (Huang, Li and Chen, 2005), and the Katz measure was found to be the most useful, followed by preferential attachment, common neighbours and the Adamic\Adar measure. These path-based and neighbour-based measures outperformed simpler metrics, and so we used these in our experiments. We were able to show that when emergent temporal trends are ignored in sociogram sequences, the link prediction accuracy is significantly lower than when included.

An important contribution of this paper is that we have shown that temporal metrics are an extremely valuable new contribution to link prediction, and should be used in future research and applications. This is significant because, to date, no temporal information has been used in link prediction research, excluding our earlier proposition of the idea in April, Potgieter & Cooke (2005) and Potgieter, April, Cooke & Lockett (2006). Further research to be performed includes suggesting new temporal variations of static metrics and determining exactly the optimum number of time steps over which to compute temporal metrics.

We have used Bayesian networks for link prediction, and achieved extremely high prediction accuracies. Additionally, we trained the components of the Dynamic Bayesian Network in order to mine the *actual relationships* between the metrics (local behaviours) and forming links (emergent global behaviours) at each time step – of significance is the fact that the DBN was able to *predict relationships* between the metrics and forming links at each time step, as it evolved over time. In future research, we will use these predicted relationships to study more emergent temporal phenomena such as if time graphs densify over time and at what rate the average distance between nodes decreases.

References

- Adamic, L. A. and Adar, E. (2003). "Friends and neighbors on the web," *Social Networks*, ISSN 0378-8733, 25(3): 211-230.
- April, K. A. (2002). "Guidelines for developing a K-strategy," *Journal of Knowledge Management*, ISSN 1367-3270, 6(5): 445-456.
- Baas, N. A. and Emmeche, C. (1997). "On emergence and explanation," *Intellectica*, ISSN 0769-4113, 25: 67-83.
- Brock, D. (2006). "The changing professional organization: A review of competing archetypes," *International Journal of Management Reviews*, ISSN 1468-2370, 8(3): 157-174.
- Burt, R. S. (1995). *Structural Holes: The Social Structure of Competition*, ISBN 0674843711.
- Campbell, C. S., Maglio, P. P., Cozzi, A. C. and Dom, B. (2003). "Expertise identification using email communications," *Proceedings of the Twelfth International Conference on Information and Knowledge Management (CIKM)*, November: 528-531.
- Cartwright, D. and Harary, F. (1956). "Structural balance: A generalization of Heider's theory," *Psychological Review*, ISSN 0033-295X, 63: 277-292.

- Cheng, J., Bell, D. A. and Liu, W. (2001). "Learning belief networks from data: An information theory based approach," *Proceedings of the Sixth International Conference on Information and Knowledge Management*, Las Vegas, Nevada, ACM.
- Desikan, P. and Srivastava, J. (2004). "Mining temporally evolving graphs," <http://maya.cs.depaul.edu/webkdd04/final/desikan.pdf>.
- Farrell, S., Campbell, C. and Myagmar, S. (2005). "Relescope: An experiment in accelerating relationships," *Conference on Human Factors in Computing Systems*, 2-7 April 2005, Portland.
- Fleiss, J. L. (1981). *Statistical Methods for Rates and Proportions*. 2nd edition, ISBN 0471064289.
- Galaskiewicz, J. (1979). *Exchange Networks and Community Politics*, ISBN 0803911378.
- Getoor, L. and Diehl, C. (2005). "Link mining: A survey," *ACM SIGKDD Explorations Newsletter*, ISSN 1931-0145, 7(2): 3-12.
- Granovetter, M. S. (1985). "Economic action and social structure: The problem of embeddedness," *American Journal of Sociology*, ISSN 0002-9602, 91(3): 481-510.
- Granovetter, M. S. (1992). "Problems of explanation in economic sociology," in N. Nohria and R. G. Eccles (eds.), *Networks and Organizations: Structure, Form and Action*, ISBN 0875845789, pp. 25-56.
- Huang, Z, Li, X and Chen, H. (2005). "Link prediction approach to collaborative filtering," *Proceedings of the 5th ACM/IEEE-CS joint conference on Digital libraries*, 7-11 June 2005.
- Hanneman, R. (2001). "Introduction to social network methods," <http://faculty.ucr.edu/~hanneman/SOC157/NETTEXT.pdf>.
- Henzinger, M. (2000). "Link analysis in web information retrieval," *IEEE Data Engineering Bulletin*, 23(3): 3-8.
- Holland, J. H. (1995). *Hidden Order: How Adaptation Builds Complexity*, ISBN 0201407930.
- Hosmer, D. W. and Lemeshow, S. (1989). *Applied Logistic Regression*, ISBN 0471615536.
- Huang, Z., Li, X. and Chen, H. (2005). "Link prediction approach to collaborative filtering," *5th ACM/IEEE-CS Joint Conference on Digital Libraries*, 7-11 June 2005.
- Kautz, H., Selman, B. and Shah, M. (1997). "Referral web: Combining social networks and collaborative filtering," *Communications of the ACM*, ISSN 0001-0782, 40(3): 63-65.
- Leskovec, J., Kleinberg, J. and Faloutsos, C. (2005). "Graphs over time: Densification laws, shrinking diameters and possible explanations," *Proceedings of the 11th ACM SIGKDD International Conference on Knowledge Discovery in Data Mining*, 21-24 August 2005.
- Liben-Nowell, D. (2005). *An Algorithmic Approach to Social Networks*. PhD thesis at MIT Computer Science and Artificial Intelligence Laboratory.
- Liben-Nowell, D. and Kleinberg, J. (2003). "The link prediction problem for social networks," *Proceedings of the 12th International Conference on Information and Knowledge Management (CIKM)*, November: 556-559.
- Lincoln, J. R. (1982). "Intra- (and inter-) organizational networks," in S. B. Bacharach (ed.), *Research in the Sociology of Organizations*, Vol 1., Greenwich, CT: JAI Press, pp. 1-38.
- Lim, M., Negnevitsky, M. and Hartnett, J. (2005). "Artificial intelligence applications for analysis of email communication activities," *Proceedings International Conference On Artificial Intelligence In Science And Technology*, pp. 109-113.
- Loasby, B. J. (1999). "Capabilities," in B. J. Loasby (ed.), *Knowledge, Institutions and Evolution in Economics*, ISBN 0415205379, pp. 49-68.
- McMahon, S. M., Miller, K. H. and Drake, J. (2001). "Networking tips for social scientists and ecologists," *Science*, ISSN 0036-8075, 293:1604-1605.
- Nohria, N. (1992). "Is a network perspective a useful way of studying organizations?," in N. Nohria and R. G. Eccles (eds.), *Networks and Organizations: Structure, Form and Action*, ISBN 0875845789, pp. 1-22.

Nohria, N. and Eccles, R. G. (1992). "Face-to-face: Making network organizations work," in N. Nohria and R. G. Eccles (eds.), *Networks and Organizations: Structure, Form and Action*, ISBN 0875845789, pp. 288-308.

Osunmakinde, I. O. (2006). *Intelligent Detection of Anomalies in Telecommunications Customer Behaviour*, MSc. Thesis to be submitted, supervised by A. Potgieter, Agents Research Group, Computer Science Department, University of Cape Town, South Africa.

Pavlovic, V., Rehg, J., Cham, T. and Murphy, K. (1999). "A dynamic Bayesian network approach to figure tracking using learned dynamic models," *ICCV '99 – International Conference on Computer Vision*.

Pearl, J. (1988). *Probabilistic Reasoning in Intelligent Systems: Networks of Plausible Inference*, 2nd Edition, ISBN 0934613737.

Pearl, J. and Russell, S. (2000). "Bayesian networks," http://bayes.cs.ucla.edu/csl_papers.html.

Pepper, S. (2003). "ISO 13250:2002 – Topic maps: An international standard knowledge representation for humans and agents," [http://www.csc.liv.ac.uk/~valli/WG-Ontology/Helsinki/Id3-Ontopia.ppt#649,1,ISO 13250:2002 – Topic Maps](http://www.csc.liv.ac.uk/~valli/WG-Ontology/Helsinki/Id3-Ontopia.ppt#649,1,ISO%2013250:2002%20-%20Topic%20Maps).

Popescul, A and Ungar, L. H. (2004). "Cluster-based concept invention for statistical relational learning," <http://www.cis.upenn.edu/~popescul/Publications/popescul04clusterbased.pdf>.

Popescul, A. and Ungar, L. H. (2003). "Structural logistic regression for link analysis," <http://www.cis.upenn.edu/~popescul/Publications/popescul03mrdm.pdf>.

Potgieter, A., April, K. and Cooke, R. (2005). "Using Bayesian agents to enable distributed network knowledge: A critique," <http://www.mngt.waikato.ac.nz/ejrot/cmsconference/2005/abstracts/socialnetworks/Potgieter.pdf>.

Potgieter, A., April, K., Cooke, R. and Lockett, M. (2006). "Adaptive Bayesian agents: Enabling distributed social networks," *South African Journal of Business Management*, ISSN 0378-9098, 37(1): 41-55.

Potgieter, A., April, K. and Bishop, J. (2005). "Complex adaptive enterprises," in M. Khosrow-Pour (ed.), *Encyclopedia of Information Science and Technology*, Vol. 1, ISBN: 159140553X, pp. 475-480.

Pujol, J., Sanguesa, R. and Delgado, J. (2002). "Extracting reputation in multi-agent systems by means of social network topology," *AAMAS*, ISSN 0302-9743, 1:467-474.

Ron, K. and Foster, P. (1998). "Special issue on applications of machine learning and the knowledge discovery process," *Journal of Machine Learning*, ISSN 1532-4435, 30: 271-274.

Ross, S., Westerfield, R., Jordan, B. and Firer, C. (2001). *Fundamentals of Corporate Finance*, ISBN 007470799X.

Russell, S. J. and Norvig, P. (2002). *Artificial Intelligence: A Modern Approach*, 2nd edition, ISBN 0137903952.

Ryle, G. (1949). *The Concept of Mind*, ISBN 0226732959.

Taskar, B., Wong, M-F., Abbeel, P. and Koller, D. (2004). "Link prediction in relational data," *Proceedings of Neural Information Processing Systems*.

Van Wijk, R., Van Den Bosch, A. J. and Volberda, H. W. (2003). "Knowledge and networks," in M. Easterby-Smith. and M. A. Lyles (eds.), *The Blackwell Handbook of Organizational Learning and Knowledge Management*, ISBN 140513304X, pp. 428-453.

Wassermann, S. and Faust, S. (1994). *Social Network Analysis: Methods and Applications*, ISBN 0521382696.

Witten, I. and Frank, E. (2005). *Data Mining: Practical Machine Learning Tools and Techniques*, 2nd edition, ISBN 0120884070.

Zhou, D. and Scholkopf, B. (2004). "A regularization framework for learning from graph data," <http://www.cs.umd.edu/projects/srl2004/Papers/zhou.pdf>.

Zhu, J. (2003). "Mining web site link structures for adaptive web site navigation and search,"
<http://kmi.open.ac.uk/people/jianhan/thesis.pdf>.