

A framework for imperfectly observed networks

David Aldous

2 May 2016

A math model of a real-world network typically starts as a graph. This is weird, because almost all real networks are better represented as *edge-weighted* graphs. The reason this isn't the default (I guess) is that there are several conceptually different interpretations of edge-weight:

- flow capacity (road network, water network)
- distance or cost (TSP)
- strength of association (close friend or acquaintance or Facebook friend).

I'll consider the last class and think of *social networks* – collaboration networks, corporate directorships, Senators' voting record, etc (note many biological networks are also in this class). Even within this class of social networks there are different interpretations of *strength of association* , but (envisaging *friends*) I abstract this as *frequency of interaction*.

Introduce randomness by saying:

for each edge $e = (vy)$, individuals v and y interact at the times of a rate- w_e Poisson process.

So this is the meaning of the edge-weights $w_e \geq 0$.

Aside. As discussed in my 2013 paper *Interacting Particle Systems (IPS) as Stochastic Social Dynamics* this setup underlies what probabilists call IPS: each individual is in some “state” and some update rule changes the states when individuals interact. This covers numerous models like the voter model or SIS epidemic – a line of research going back to statistical physics study of the Ising model on \mathbb{Z}^d .

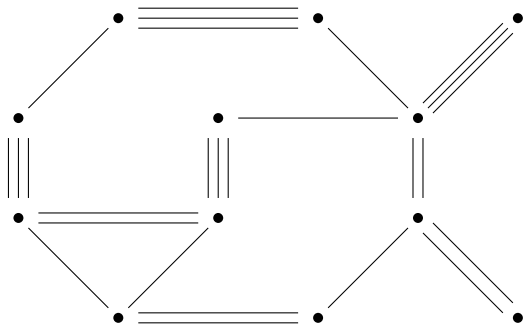
This talk goes in a different direction: Suppose we are interested in some quantitative feature of a network which we could calculate if we knew exactly what the network is.

But suppose we don't know it then what can we do?

I'll call this the **imperfectly-observed network** problem. I will talk about one particular formalization – not claimed to be useful for real-world data but (I do claim) interesting as math theory.

A **network** is a finite edge-weighted graph. We are concerned with some “statistic” Γ , a functional $G \rightarrow \Gamma(G)$ on finite edge-weighted graphs G . There is a network G^{true} with known vertices but unknown edges and edge-weights w_e . What we observe is the interaction process described above. That is, what we observe over time $[0, t]$ is the $\text{Poisson}(tw_e)$ number of interactions $N_e(t)$ over edges e . We can represent our observations in two equivalent ways: either as the random multigraph with $N_e(t)$ copies of edge e , or as the random weighted graph $G^{\text{obs}}(t)$ in which edge e has weight $t^{-1}N_e(t)$.

How do we use these observations to estimate $\Gamma(G^{\text{true}})$, and how accurate is the estimate?



Some general comments.

- For any problem about networks where you assumed the network is known, you could ask this “imperfectly-observed” variation.
- There are many other ways to think about “imperfectly-observed networks” [one popular way will be shown later].
- We always have the naive frequentist estimator $\Gamma(G^{\text{obs}}(t))$. It's natural to study, but there is no reason to think it is optimal.
- We always have the naive Bayes estimator (flat prior on each w_e) but
- “Computation is free” – not concerned with computational complexity – instead we regard observation time as the “cost”.

Any estimator like $\Gamma(G^{\text{obs}}(t))$ for fixed t will have error depending on the unknown G^{true} . The “elegant” formulation of a mathematical problem is:

Program

*Given a statistic Γ , define a (“universal”) stopping rule T and an estimator such that the relative error of the estimator, say $\Gamma(G^{\text{obs}}(T))/\Gamma(G^{\text{true}}) - 1$, is small **uniformly** over all networks G^{true} .*

Program

*Given a statistic Γ , define a (“universal”) stopping rule T and an estimator such that the relative error of the estimator, say $\Gamma(G^{\text{obs}}(T))/\Gamma(G^{\text{true}}) - 1$, is small **uniformly** over all networks G^{true} .*

The bottom line of this talk. We have no idea how to do this for most interesting/natural statistics, but we can do this for a few statistics which are less interesting/natural.

This is ongoing joint work with grad student Lisha Li.

Given G , write n for the number of vertices and $w_v = \sum_y w_{vy}$ for the total interaction rate of vertex v . We are thinking of results for large networks, formalized as $n \rightarrow \infty$ limits. **For discussion purposes here** (not as assumptions in theorems) assume $w_v \equiv 1$, so in time t we see on average t edges per vertex.

Qualitatively there are 3 time regimes.

- For $t = o(1)$ can only estimate statistics like (weighted) degree distributions (cf. birthday problem).
- To make the observed graph connected we need $t = \Theta(\log n)$ (cf. coupon collector problem) at which time we see $\Theta(\log n)$ edges per vertex and (intuitively) “we can estimate anything well”.
- The interesting/challenging regime is where t is a (large-ish) constant; what can we infer when we have seen 13 interactions per individual?

Here is a fundamental, albeit vague, open problem.

if we observe $G^{\text{obs}}(t)$ has a “highly connected” (in some sense) giant vertex set of size αn , then we can infer that G^{true} has a similarly “highly connected” giant vertex set of size $\beta(\alpha)n$?

There are many ways to quantify connectedness by a statistic Γ in this context, for instance via spectral gap of the (restricted) graph Laplacian. We conjecture that our program (repeated below) can be done in this setting. The *intuition* is that randomness makes G^{obs} less well connected than G^{true} – but we have no idea how to prove any reasonable version.

Program

*Given a statistic Γ , define a (“universal”) stopping rule T and an estimator such that the relative error of the estimator, say $\Gamma(G^{\text{obs}}(T))/\Gamma(G^{\text{true}}) - 1$, is small **uniformly** over all networks G^{true} .*

On the positive side, here is a “sideways” approach to our program.
Consider

$$T_k^{tria} = \inf\{t : \text{observed multigraph contains } k \text{ edge-disjoint triangles}\}.$$

$$T_k^{span} = \inf\{t : \text{observed multigraph contains } k \text{ edge-disjoint spanning trees}\}.$$

Proposition

$$\frac{\text{s.d.}(T_k^{tria})}{\mathbb{E}T_k^{tria}} \leq \left(\frac{e}{e-1}\right)^{1/2} k^{-1/6}, \quad k \geq 1.$$

$$\frac{\text{s.d.}(T_k^{span})}{\mathbb{E}T_k^{span}} \leq k^{-1/2}, \quad k \geq 1.$$

So here the bounds are independent of \mathbf{w} , meaning that we can estimate the statistics $\mathbb{E}T_k$ without assumptions on \mathbf{w} .

So the “sideways” approach is to seek some observable quantity which is concentrated around its mean, independent of \mathbf{w} , which therefore provides an estimator of the statistic defined by the expectation.

Proposition from arXiv preprint *Weak Concentration for First Passage Percolation Times on Graphs and General Increasing Set-valued Processes* and the title give a hint of the proof method. Our observation process, considered as a growing multigraph, is an increasing set-valued process, for which there is a simple general bound on $\frac{\text{s.d.}(T)}{\mathbb{E}T}$ for the first time T that some “increasing” property holds. In our context, we have

$$T_k = \inf\{t : \text{observed multigraph contains } k \text{ edge-disjoint **objects**\}$$

and the argument for the bound uses only one object-specific calculation, which I will outline as a game, which is trivial in the two cases (triangles and spanning trees) above.

The game. I choose a multigraph with the given “contains k edge-disjoint **objects**” property, and I then delete an edge, and then show you. Can you always find many different ways to restore the property by creating a few new edges?

Spanning trees; deleting edge creates a split $(A, \mathcal{V} \setminus A)$ of vertex-set \mathcal{V} ; sufficient for you to create any edge between A and $\mathcal{V} \setminus A$.

Triangles: sufficient for you to create one new triangle.

The bound in the general inequality involves (worst-case) mean “restore” time in the observation process.

Open problem; Can we do this for the “ k -edge connected” property? (Menger’s theorem doesn’t seem to help).

Here is a first example of a “natural” statistic. Identify a graph with its matrix \mathbf{w} of edge-weights.

Maximum matching. Take n even. A *matching* is a set π of $n/2$ edges such that each vertex is in exactly one edge.

The *weight* of the matching is $\text{weight}(\pi, \mathbf{w}) := \sum_{e \in \pi} w_e$.

The *maximum-weight* is $\Gamma_1(\mathbf{w}) := \max_{\pi} \text{weight}(\pi, \mathbf{w})$.

Can we estimate $\Gamma_1(\mathbf{w})$ from the observed $G^{\text{obs}}(t)$ at (large) times $t = O(1)$?

The naive frequentist estimator $\Gamma_1(G^{\text{obs}}(t))$ does not work – consider the “dense” case of the complete graph with edge-weights $w_e = 1/(n-1)$.

We will finesse this issue by reformulating the problem. Because real-world networks are typically sparse, we can say that, although we require our estimates to be **valid** for all G^{true} , we only require them to be **informative** for sparse G^{true} .

Informally, we regard a weighted graph as *sparse* if the vertex-weight sums $w_v = \sum_y w_{vy}$ are dominated by the largest $O(1)$ terms.

For discussion, assume $\max_v w_v = 1$. For a sparse graph we will have $\Gamma_1(\mathbf{w}) = \Theta(n)$, so we reformulate the problem as

can we estimate $n^{-1}\Gamma_1(\mathbf{w})$ up to small additive error?

Such an estimator will be informative in the sparse case, but not for *dense* graphs like the complete graph above, for which $\Gamma_1 = \Theta(1)$.

A moment's thought says that to know anything about the weight of some specific edge we must observe at least two interactions (cf. unseen species problem).

This suggests making an estimator using only edges for which we have observed at least two “interactions”. That is, we define

$$\text{weight}_2(\pi, G^{\text{obs}}(t)) := t^{-1} \sum_{e \in \pi} N_e(t) \mathbf{1}_{\{N_e(t) \geq 2\}}$$

$$\Gamma_2(G^{\text{obs}}(t)) := \max_{\pi} \text{weight}_2(\pi, G^{\text{obs}}(t))$$

and our goal is to obtain a bound of the form

$$\mathbb{E} n^{-1} |\Gamma_2(G^{\text{obs}}(t)) - \Gamma_1(\mathbf{w})| \leq \psi(t) \quad \forall \mathbf{w}. \quad (1)$$

The best we can hope for is a $\psi(t) = O(t^{-1/2})$ bound: consider the graph with only one edge. And a conceptually straightforward argument (large deviations and counting) shows (1) is true for some

$$\psi(t) = O(t^{-1/2} \log t).$$

[Also a factor $\max_v w_v$, but we can estimate this more quickly].

Observed and true community structure.

For a subset A of vertices write A^* for the set of edges with both end-vertices in A . Write

$$\bar{\mathbf{w}}_m = m^{-2} \max \left\{ \sum_{e \in A^*} w_e : |A| = m \right\}$$

– essentially the maximum edge-density in a size- m community. Ignoring computational complexity, suppose we can compute the analogous observable quantity

$$\bar{W}_m(t) = m^{-2} \max \left\{ \sum_{e \in A^*} N_e(t)/t : |A| = m \right\}.$$

To make inferences from the observed $G^{\text{obs}}(t)$ to G^{true} we need $m \sim \gamma \log n$. Then (as in previous example, just using large deviations and counting) we can be confident that $\bar{\mathbf{w}}_m$ is in a certain interval, roughly

$$\left[\bar{W}_m(t) - \sqrt{\frac{2\bar{W}_m(t)}{\gamma t}}, \bar{W}_m(t) \right].$$

Here is one case where it seems **impossible** to carry out this program. It is a basic example of a process built over a weighted graph.

First passage percolation (FPP).

Given an edge-weighted graph (G, \mathbf{w}) with distinguished vertices (v^*, v^{**}) , create independent random variables ξ_e with $\text{Exponential}(w_e)$ distributions, and view ξ_e as the “traversal time” of edge e . Let $X(\mathbf{w})$ be the (random) FPP time from v^* to v^{**} , that is the minimum value of $\sum_{e \in \pi} \xi_e$ over all paths π from v^* to v^{**} . Take the expectation of this FPP time as our statistic

$$\Gamma(\mathbf{w}) = \mathbb{E}X(\mathbf{w}).$$

We will argue informally that there is no “good” general stopping rule T for which $\Gamma(G^{\text{obs}}(T))$ is a good estimate of $\Gamma(G^{\text{true}})$. That is, one cannot improve on using the observation process to simulate the FPP process itself.

First consider the linear graph on vertex set $\{v^* = 0, 1, 2, \dots, n = v^{**}\}$ with unknown edge-weights $\mathbf{w} = (w_i, 1 \leq i \leq n)$. Clearly

$$\Gamma(\mathbf{w}) = \sum_{i=1}^n 1/w_i.$$

Fix $k \geq 1$ and consider the first time that we have observed k interactions across each edge:

$$T_k = \min\{t : N_i(t) \geq k, 1 \leq i \leq n\}.$$

It is intuitively clear (and true) that $\Gamma(G^{\text{obs}}(T_k))$ should be a good estimator for $\Gamma(\mathbf{w})$ for the linear graph. By analogy with earlier “weak concentration” results one might hope this holds generally for

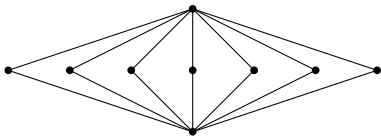
$T_k =$ time until observe k edge-disjoint paths from v^* to v^{**} .

But this is false, as explained below.

For the linear graph with weights $w_i \equiv 1$ we have

“observation time needed” $T_k = \Theta(\log n)$; actual FPP time $\Gamma(G^{\text{true}}) = n$.

Consider instead the graph with n 2-edge routes from v^* to v^{**} , and with edge weights w .



Here both “observation time needed” and actual FPP time are the same order, $\Theta(w^{-1}n^{-1/2})$. So choose $w_n = 1/n$ to make this $n^{1/2}$.

Now superimpose this graph and the linear graph. Then at time $\Theta(\log n)$ we observe the presence of the linear route, for which FPP time is n , but we do not observe the presence of the shorter-time 2-edge routes until time $n^{1/2}$.

Bottom line: a “universal” algorithm for this statistic cannot stop before the actual FPP time, because there might be unobserved analogs of such 2-edge paths “in the gap”.

Comments on example above.

- The example is artificial; perhaps $\Gamma(G^{\text{obs}}(T_k))$ is indeed a good estimator of $\Gamma(\mathbf{w})$ under some weak assumptions on \mathbf{w} .
- A simple argument shows that for this FPP statistic Γ the natural estimator is an overestimate. Precisely, the unconditional distribution of $X(G^{\text{obs}}(t))$ stochastically dominates $X(\mathbf{w})$.

As mentioned earlier, a major open problem is to prove some version of the latter for a statistic measuring “connectivity of giant component”.

Note the **weird logic**: usually with a random structure we are interested in proving some desirable property holds. In our framework we want to reach a conclusion of the format

if the observed graph has a given desirable property then we can be confident that the true graph has a similar property.

So we want the observed graph to have (slightly) worse behavior, as regards the given property.

A very different framework for “imperfectly-observed networks”.

[from a 2011 survey *Link prediction in complex networks* by Linyuan Lü and Tao Zhou, cited 683 times.]

Consider unweighted graphs, and only the possibility of unobserved edges – this is called *link prediction*. In this literature, the goal is to define an algorithm that takes the observed edges as input, and outputs an ordering e_1, e_2, \dots of all the other possible edges, intended as decreasing order of assessed “likelihood” of the edge being present. This is done by defining, for each possible edge (v_1, v_2) , some statistic based on (typically) the local structure of the observed graph near v_1 and v_2 , for instance

$$s(v_1, v_2) = \frac{|\mathcal{N}(v_1) \cap \mathcal{N}(v_2)|}{|\mathcal{N}(v_1)| |\mathcal{N}(v_2)|}$$

where $\mathcal{N}(v)$ is the set of neighbors of v . Then list edges in decreasing order of $s(v_1, v_2)$.

In this framework there is no probability model involved; different algorithms are compared empirically by taking a real-world network, randomly deleting a proportion of edges to create a synthetic “observed graph”, and comparing the algorithms’ effectiveness in predicting the deleted edges.

Returning to our framework – unknown G^{true} and an observed $G^{\text{obs}}(t)$ – it is conceptually simpler to take the Bayesian view. Put a prior on G^{true} , compute the posterior distribution of G^{true} given $G^{\text{obs}}(t)$, then any given statistic has a posterior distribution.

In particular, if we assume G^{true} is connected and wish to estimate the spectral gap of the graph Laplacian, in our previous setup we need $t = \Theta(\log n)$ to make $G^{\text{obs}}(t)$ connected and get a non-trivial estimate, where in the Bayes setup we can put a prior on connected graphs.

But not so easy in practice – how do you choose a plausible prior?

To play with the mathematics, consider the “naive Bayes” procedure – take as prior the uniform law on $[0, \infty)$ for each w_{ij} – for which the posterior distribution on \mathbf{w} given observed interactions (n_{ij}) is that the w_{ij} are independent with densities

$$\nu \rightarrow p(n_{ij}; \nu t) \tag{2}$$

where $p(k; \lambda)$ denotes the Poisson probability function.

Informally, this “flat” prior lives on highly connected graphs, and for small t the posterior distribution on \mathbf{w} will concentrate on too-highly-connected graphs, with spectral gap around ne^{-t} . So we will not get a good estimate of true spectral gap before time $\Theta(\log n)$.