

Invariant Theory for Hypothesis Testing on Graphs

Priebe, Carey

Johns Hopkins University, Applied Mathematics and Statistics

3400 North Charles Street

Baltimore, Maryland 21218-2682, USA

E-mail: cep@jhu.edu

Rukhin, Andrey

Naval Surface Warfare Center, Sensor Fusion Department

18444 Frontage Road

Dahlgren, Virginia 22448, USA

E-mail: andrey.rukhin@navy.mil

1 Introduction

Following the setting outlined in Priebe et al. [2011] we aim to detect anomalies within attributed graphs. In particular, let $\mathcal{V} = \{1, \dots, n\}$ be the fixed set of vertices and $\phi : \binom{\mathcal{V}}{2} \rightarrow \{0, 1, \dots, K\}$ be an *edge-attribution* function. The graph on \mathcal{V} is defined to be $G = (\mathcal{V}, \mathcal{E}_\phi)$ where

$$(u, v) \in \mathcal{E}_\phi \iff \phi(u, v) > 0.$$

We say that the edge (u, v) has attribute $c \in \{1, \dots, K\}$ if $\phi(u, v) = c$. One can view the categorical edge attributes as some mode of the communication event between actors u and v (e.g., a topic label derived from the content of the communication).

The specific anomaly we aim to detect is the “chatter” alternative – a small (unspecified) subset of vertices with altered communication behavior in an otherwise homogeneous setting. Our inference task is to determine whether or not a graph $(\mathcal{V}, \mathcal{E}_\phi)$ includes a subset of vertices $\mathcal{M} = \{v_1, v_2, \dots, v_m\}$ whose edge-connectivity within the subset exhibits a different behaviour than that found among the remaining vertices in the graph.

To this end we consider the problem of detecting chatter anomalies in a graph using hypothesis testing on a *fusion* of attributed graph invariants. In particular, the focus of this paper is analyzing and comparing the inferential power of the linear attribute fusion of the attributed q -clique invariant

$$T_q^W(G) = \sum_{c_1, \dots, c_K \in P(\binom{[q]}{2}, K)} w_{c_1, \dots, c_K} \sum_{(u_1, \dots, u_q) \in \binom{\mathcal{V}}{q}} h(u_1, \dots, u_q; c_1, \dots, c_K),$$

where the sum is over the collection of partitions $P(\binom{[q]}{2}, K)$ of $\binom{[q]}{2}$ into K non-negative parts, $W = \{w_i\}_{i \in P(\binom{[q]}{2}, K)}$ are the fusion weights, and the summand $h(u_1, \dots, u_q; c_1, \dots, c_K)$ indicates the event that the vertices u_1, \dots, u_q are elements of a q -clique with c_r edges of color r . Specifically, we consider the cases $q = 2$ which yields the *size* fusion T_2^W and $q = 3$ which yields the *triangle* fusion T_3^W .

Our random graph model is motivated by the time series model found in Lee and Priebe [forthcoming]: for each vertex $v \in \mathcal{V}$ we assign a latent variable $\mathbf{X}^v = (X_1^v, \dots, X_d^v)$ drawn independently of all other vertices from some d -dimensional distribution. The edge-attribution function will be a random variable where the probability of an edge (u, v) having attribute c is defined to be a some predetermined function of the inner product of the latent variables. We assume that the edge attributes, conditioned on the latent variables, are independent. In this paper, we will assume that $\mathbf{X}^v \sim \text{Dirichlet}(\lambda_0^v, \dots, \lambda_K^v)$ and

$$\mathbb{P}\{\phi(u, v) = c\} = X_c^u X_c^v$$

for all $(u, v) \in \binom{\mathcal{V}}{2}$ and all $c \in \{1, \dots, K\}$. This model choice is analogous to the first and second approximations found in Lee and Priebe [forthcoming]: if we write $\lambda^v = (\lambda_0^v, \dots, \lambda_K^v) = (1 + rx_0^v, \dots, 1 + rx_K^v)$ for some fixed (x_0^v, \dots, x_K^v) in the unit simplex and non-negative real r , then $r \rightarrow \infty$ yields the first approximation model (i.e., the “independent edge model”). We mention that our approach herein differs from the second approximation in Lee and Priebe [forthcoming]; their second approximation yields a inner product model with truncated Gaussian latent variables.

Related work may be found in Bollobas et al. [2007] Section 16.4 and references therein. We also direct the interested reader to Priebe et al. [2011] in which the authors study other linear attribute fusion invariants; in particular, the authors consider

$$\text{maxd}^W(G) = \max_v \sum_{c=1}^K w_c \sum_{u \in N[v]} I\{\phi(u, v) = c\}$$

and

$$\text{scan}^W(G) = \max_v \sum_{c=1}^K w_c \sum_{u, x \in N[v]} I\{\phi(u, x) = c\},$$

where $N[v] = \{u \mid (u, v) \in \mathcal{E}\} \cup \{v\}$ is the closed neighborhood of vertex v in the graph.

Finally, we add that we will restrict ourselves to simple undirected graphs. We will not consider hyper-graphs (hyper-edges consisting of more than two vertices), multi-graphs (more than one edge between any two vertices), self-loops (an edge from a vertex to itself), or weighted edges.

2 Notation

For each positive integer l we use the notation $[l] = \{1, \dots, l\}$.

For each $v \in \mathcal{V}$ we assign a *latent position* vector $\mathbf{X}^v = (X_0^v, \dots, X_K^v) \sim \text{Dirichlet}(\lambda^v)$ for some fixed parameter vector $\lambda^v \in \mathbb{R}_+^{K+1}$. We also assume that the latent positions are independent.

Our null hypothesis assumes a version of homogeneity among the vertices; specifically,

$$\mathbb{H}_0 : \mathbf{X}^v = (X_0^v, \dots, X_K^v) \sim \text{Dirichlet}(\lambda) \text{ for all } v \in \mathcal{V}$$

for some Dirichlet parameter vector $\lambda = (\lambda_0, \dots, \lambda_K)$. Our alternative hypothesis incorporates the anomaly feature described in the preceding section as follows. Assume $m = m(n) < n$ satisfies the following two conditions: $\lim_{n \rightarrow \infty} m(n) = \infty$ and $\lim_{n \rightarrow \infty} \frac{m(n)}{n} = 0$. Our alternative hypothesis is defined to be

$$\mathbb{H}_1 : \mathbf{X}^v = \begin{cases} (Y_0^v, \dots, Y_K^v) \stackrel{iid}{\sim} \text{Dir}(\eta) & i \in [m], \\ (X_0^v, \dots, X_K^v) \stackrel{iid}{\sim} \text{Dir}(\lambda) & i \in [n] - [m]. \end{cases}$$

for some fixed Dirichlet parameter vector $\eta = (\eta_0, \dots, \eta_K)$ and the same $\lambda = (\lambda_0, \dots, \lambda_K)$ from the null hypothesis. For convenience, we also define $\Lambda = \sum_{0 \leq c \leq K} \lambda_c$ and $H = \sum_{0 \leq c \leq K} \eta_c$.

We define

$$\varepsilon_c = \sum_{(u,v) \in \binom{\mathcal{V}}{2}} \mathbb{I}\{\phi(u, v) = c\}$$

to be the *size* (i.e., number of 2-cliques) of attribute c in the graph. Similarly, for the number of *triangles* (i.e., 3-cliques) we write τ_c , $\tau_{b,c}$, and $\tau_{b,c,d}$ to denote the number of 3-cliques with three c -colored edges, two b -colored and one c -colored edge, and one edge of each of three edge-colors b, c, d , respectively.

Before proceeding, we highlight a relevant property of the mixed moments of the Dirichlet distribution (see Johnson and Kotz [1972]): if r_1, \dots, r_s are non-negative and $(X_1, \dots, X_s) \sim \text{Dirichlet}(\theta_1, \dots, \theta_s)$ then

$$(1) \quad E \left[\prod_{i=1}^s X^{r_i} \right] = \frac{\Gamma(\sum_{i=1}^s \theta_i) \prod_{j=1}^s \Gamma(\theta_j + r_j)}{\prod_{i=1}^s \Gamma(\theta_i) \Gamma(\sum_{j=1}^s (\theta_j + r_j))}$$

where Γ denotes Euler’s standard Gamma function. With this property one can compute the exact moments of the Hajek projection of T_2^W and T_3^W under either hypothesis. We write $\nu_{(c)}^{(i)}$ to denote the i -th moment of X_c in the null latent vector and $\nu_{(b,c)}^{(i,j)}$ to denote the joint (i, j) -moment of (X_b, X_c) . Similarly, we’ll write $\mu_{(c)}^{(i)}$ to denote the i -th moment of Y_c in the anomalous latent vector and $\mu_{(b,c)}^{(i,j)}$ to denote the joint (i, j) -moment of (Y_b, Y_c) .

3 Analysis

We will appeal to Hajek’s Projection method, detailed in Nowicki and Wierman [1988], in order to demonstrate the asymptotic normality of the fusion invariants in this article. This approach is outlined as follows: We define the *projection* of the fusion T to be the centered sum of independent random variables

$$T^* = \sum_{v \in \mathcal{V}} E[T | \mathbf{X}^v] - (n - 1)E[T].$$

For both the size and triangle fusion we aim to show that

$$\frac{T - E[T]}{\sqrt{\text{Var}(T^*)}} = \frac{T - T^*}{\sqrt{\text{Var}(T^*)}} + \frac{T^* - E[T]}{\sqrt{\text{Var}(T^*)}} \xrightarrow{\mathcal{D}} N(0, 1).$$

To this end, one appeals to Chebyshev’s Inequality to show that $\text{Var}(T - T^*) = o(\text{Var}(T^*))$ (see Nowicki and Wierman [1988] for the detailed argument). Specifically, if

$$\mathbb{P} \left\{ \frac{|T - T^*|}{\sqrt{\text{Var}(T^*)}} \geq \epsilon \right\} \leq \frac{\text{Var}(T - T^*)}{\epsilon^2 \text{Var}(T^*)} \rightarrow 0$$

for any positive ϵ , then

$$\frac{T - T^*}{\sqrt{\text{Var}(T^*)}} + \frac{T^* - E[T]}{\sqrt{\text{Var}(T^*)}} \xrightarrow{\mathcal{D}} N(0, 1)$$

by applying the Central Limit Theorem to the normalized sum of independent random variables in the second term of the left-hand side.

3.1 The Attributed Size Fusion

For each $c \in K$ define

$$\epsilon_c = \sum_{(u,v) \in \binom{\mathcal{V}}{2}} \mathbb{I}\{\phi(u, v) = c\}$$

to be the number of edges of color c in the graph. The linear attribute fusion with parameter $W = (w_1, \dots, w_K)$ is defined to be

$$T_2^W = \sum_{c=1}^K w_c \epsilon_c.$$

3.1.1 The Attributed Size Fusion under \mathbb{H}_0

We present all relevant terms within the Hajek Projection of the attributed size fusion under the null hypothesis.

The expectation of T_2^W under the null is given by

$$E_0 [T_2^W] = \binom{n}{2} \sum_{c=1}^K w_c \left[\nu_{(c)}^{(1)} \right]^2.$$

$E_0 [T_2^W | \mathbf{X}^a]$ for any fixed $a \in \mathcal{V}$ is given by

$$E_0 [T_2^W | \mathbf{X}^a] = (n - 1) \sum_{c=1}^K w_c \left[\nu_{(c)}^{(1)} \right] X_c^{(a)} + \binom{n - 1}{2} \sum_{c=1}^K w_c \left[\nu_{(c)}^{(1)} \right]^2.$$

We can now evaluate T^* under the null:

$$\begin{aligned} T^* &= \sum_{a \in [n]} E [T_2^W | \mathbf{X}^a] - (n - 1)E[T_2^W] \\ &= (n - 1) \sum_{a \in [n]} \sum_{1 \leq c \leq K} w_c \left[\nu_{(c)}^{(1)} \right] X_c^{(a)} - \binom{n}{2} \sum_{c=1}^K w_c \left[\nu_{(c)}^{(1)} \right]^2. \end{aligned}$$

The variance of this sum of independent and identically distributed random variables is

$$Var_0 (T^*) = \Theta (n^3).$$

As $Var_0 (T_2^W - T^*) \leq \binom{n}{2} (2+1)^2 E [\mathbb{I} \{ \phi(u, v) > 0 \}] = o(Var_0 (T^*))$ by the Cauchy-Schwarz Inequality (see Nowicki and Wierman [1988] for full details), we have the desired convergence to the standard normal distribution.

3.1.2 The Attributed Size Fusion under \mathbb{H}_1

For the alternative we write the edge attribution function as

$$\mathbb{P} \{ \phi(u, v) = c \} = \begin{cases} Y_c^u Y_c^v & u, v \in [m], \\ Y_c^u X_c^v & u \in [m], v \in [n] - [m], \\ X_c^u X_c^v & u, v \in [n] - [m]. \end{cases}$$

We perform a similar but more involved analysis to deduce the limiting distribution of the attributed size fusion of the graph under these conditions, yielding

$$Var_1 (T^*) = \Theta (n^3)$$

and

$$Var_1 (T_2^W - T^*) = \Theta (n^2) = o(Var_1 (T^*))$$

as desired.

3.1.3 Asymptotic Power Analysis of the Attributed Size Fusion

Returning to the context of hypothesis testing, assume we are interested in performing an α -level hypothesis test to determine whether or not the graph includes an anomalous set of m vertices whose underlying latent distribution differs from the the null component of the graph. We define $\beta_2^W = \lim_{n \rightarrow \infty} \mathbb{P}_1 \{ T_2^W > c_\alpha \}$ where $c_\alpha = c(\alpha, n)$ is the α -level critical value of the test.

Fix $c \in [K]$. The difference of the corresponding terms the hypotheses means can be written as

$$E_1 [\varepsilon_c] - E_0 [\varepsilon_c] = D_1^{(c)} + D_2^{(c)} = \binom{m}{1} \binom{n-m}{1} \nu_{(c)}^{(1)} (\mu_{(c)}^{(1)} - \nu_{(c)}^{(1)}) + \binom{m}{2} \left([\mu_{(c)}^{(1)}]^2 - [\nu_{(c)}^{(1)}]^2 \right)$$

(here $D_i^{(c)}$ corresponds to the edge-count that includes edges with exactly i anomalous vertices). The reader can verify that $\frac{Var_1(T^*)}{Var_0(T^*)} \rightarrow 1$. Moreover, given that the limiting distribution (under the null) is normal, we write

$$c_\alpha = z_\alpha \sqrt{Var_0(T^*)} + E_0 [T_2^W]$$

and thus

$$\beta_2^W = \mathbb{P} \left\{ Z > z_\alpha - \lim_{n \rightarrow \infty} \left(\frac{E_1 [T_2^W] - E_0 [T_2^W]}{\sqrt{Var_0(T^*)}} \right) \right\}.$$

Recall that $Var_0(T^*) = \Theta(n^3)$; thus, if

$$\sum_c w_c D_1^{(c)} \neq 0$$

(i.e. there is signal in the null-to-anomaly connectivity) then the limiting power $\beta_2^W > \alpha$ when $\frac{m(n-m)}{\sqrt{n^3}} \rightarrow 0$ or, equivalently, when $m = \Omega(\sqrt{n})$ (similarly, if $m = \omega(\sqrt{n})$ then $\beta_2^W \rightarrow 1$). Furthermore, if

$$\sum_c w_c D_1^{(c)} = 0$$

(i.e. there is no signal in the null-to-anomaly connectivity) the limiting power $\beta_2^W > \alpha$ when $\sum_c w_c D_2^{(c)} \neq 0$ and $\frac{m^2}{\sqrt{n^3}} \rightarrow 0$ (which is equivalent to $m = \Omega(\sqrt[4]{n^3})$). Moreover, if $m = \omega(\sqrt[4]{n^3})$ under these conditions then $\beta_2^W \rightarrow 1$.

It follows that the optimal choice of weights (w_1, \dots, w_K) is the one which maximizes the expression

$$\lim_{n \rightarrow \infty} \left(\frac{E_1 [T_2^W] - E_0 [T_2^W]}{\sqrt{Var_0(T^*)}} \right)$$

in either of the two above-mentioned cases.

3.2 The Attributed Number of Triangles Fusion

We begin by writing

$$\tau = \sum_{c \in [K]} \tau_c + \sum_{b \neq c} \tau_{b,c} + \sum_{d \neq b,c} \tau_{b,c,d}.$$

We denote the number-of-triangles fusion invariant to be

$$T_3^W = \sum_{c \in [K]} w_c \tau_c + \sum_{b \neq c} w_{b,c} \tau_{b,c} + \sum_{d \neq b,c} w_{b,c,d} \tau_{b,c,d}.$$

Similar to what was done in the previous section, we obtain

$$E_0 [T^*] = \binom{n}{3} \left[\sum_{c \in [K]} w_c [\nu_{(c)}^{(2)}]^3 + \sum_{b \neq c} w_{b,c} 3 \nu_{(b)}^{(2)} [\nu_{(b,c)}^{(1,1)}]^2 + \sum_{d \neq b,c} w_{b,c,d} 3 \nu_{(b,c)}^{(1,1)} \nu_{(b,d)}^{(1,1)} \nu_{(c,d)}^{(1,1)} \right]$$

and

$$Var_0(T^*) = \Theta \left(n \binom{n-1}{2}^2 \right)$$

under the null and

$$E_1[T^*] = \Theta \left(\sum_{i=0}^3 \binom{m}{i} \binom{n-m}{3-i} \right)$$

and

$$Var_1(T^*) = \Theta \left(n \binom{n-1}{2}^2 \right)$$

under the alternative. Again, $Var_0(T_3^W - T^*) = o(Var_0(T^*))$ and $Var_1(T_3^W - T^*) = o(Var_1(T^*))$.

As in the case with the attributed size fusion, we are interested in performing an α -level hypothesis test.

The terms within the difference in means can be expressed as

$$\begin{aligned} E_1[T^*] - E_0[T^*] &= \sum_{c \in [K]} D^{(c)} + \sum_{b \neq c} D^{(b,c)} + \sum_{d \neq b,c} D^{(b,c,d)} \\ &= \Theta \left(\sum_{i=1}^3 \binom{m}{i} \binom{n-m}{3-i} \delta_i(H, \Lambda) \right) \end{aligned}$$

where $\delta_i(H, \Lambda)$ is the mixed-moments difference when there are i anomalous vertices in a 3-clique.

The reader can verify that $\frac{Var_1(T^*)}{Var_0(T^*)} \rightarrow 1$. Since $Var_0(T^*) = \Theta(n^5)$, we have that the limiting power $\beta_3^W > \alpha$ when when $m = \Omega(\sqrt[2i]{n^{2i-1}})$ and the corresponding mixed-moments expression $\delta_i(H, \Lambda)$ is non-zero.

4 Conclusion

We have presented preliminary results for linear attribute fusion for clique sizes $q = 2$ and 3 in terms of inferential power when detecting the prescribed anomaly within our model. In general, the most powerful choice of q depends on m as a function of n and on the Dirichlet parameter vectors λ and η through the mixed moments.

References

- B. Bollobas, S. Janson, and O. Riordan. The phase transition in inhomogeneous random graphs. *Random Structures and Algorithm*, 31:3–122, 2007.
- N. Johnson and S. Kotz. *Distributions in Statistics: Continuous Multivariate Distributions*. Wiley, New York, 1972.
- N. Lee and C. E. Priebe. A Latent Process Model for Time Series of Attributed Random Graphs. *Statistical Inference for Stochastic Processes*, forthcoming.
- K. Nowicki and J. Wierman. Subgraph counts in random graphs using incomplete u -statistics methods. *Discrete Mathematics*, 72:299–310, 1988.
- C. E. Priebe, N. Lee, Y. Park, and M. Tang. Attribute fusion in a latent process model for time series of graphs. *The 2011 IEEE Workshop on Statistical Signal Processing (SSP2011)*, 2011.