Distribution and Dependence of Extremes in Network Sampling Processes

Jithin K. Sreedharan^{*} with Konstantine Avrachenkov^{*} and Natalia M. Markovich[†]

*INRIA Sophia Antipolis, France [†]Institute of Control Sciences, Russian Academy of Sciences, Moscow

March 30, 2015



No complete picture a priori !

All we have: X_1, X_2, \dots, X_n

Samples: any stationary (most likely dependent) sequence e.g. node ID's, degrees, number of followers or income of the nodes in OSN etc

CORRELATIONS IN GRAPHS AND SAMPLING

- Correlations in graph properties exist in real networks e.g: correlation in Coauthorship network
- Usually neglected in analysis of sampling algorithms



Effect of neglecting correlations:

- Assuming i.i.d. degrees, largest degree $\approx KN^{1/\gamma}$, N no. of nodes, γ tail index of Pareto distribution (N. Litvak et al, LNCS'12)
- Twitter graph (2012): N = 537M, $\gamma = 1.124$ for out-degree.
- Largest out-degree predicted is 59M. Actual largest out-degree is 22M!

QUESTIONS WE ADDRESS HERE...



Kth largest value of samples and many more extremal properties

Is there a simple way to get information about many extremal properties? Ans: Extremal Index

Relation to Extreme Value Theory

Extremal Index (θ) :

Definition: If $\lim_{n \to \infty} E[\text{no of exceedances}] = \tau$,



Point process of exceedances \rightarrow Compound poisson process (rate $\theta \tau$) Tendency to form clusters

$$P\{\max\{X_1,\ldots,X_n\} \le x\} = F^{n\theta}(x) + o(1), n \to \infty$$

Gives maxima of the degree sequence with certain probability

Pareto case revisited:

- i.i.d. degrees, largest degree $\approx KN^{1/\gamma}$, N no. of nodes, γ tail index of Pareto distribution (N. Litvak, LNCS'12)
- Stationary degree samples with EI, largest degree $\approx K(N\theta)^{1/\gamma}$

EXTREMAL INDEX: APPLICATIONS



Lower the value of EI, more time to hit extreme levels

e.g. Pareto
$$P(X_i > u_n) = u_n^{-\alpha}, u_n = (n)^{1/\alpha}$$
 for $\tau = 1$
 $\implies E[T_n] \approx \frac{u_n^{\alpha}}{\theta}$

EXTREMAL INDEX: APPLICATIONS

Relation to Mean Cluster Size:



 $\lim_{n\to\infty} E[\text{cluster size with } n \text{ samples}] = \frac{1}{\text{Extremal Index}}$

CALCULATION OF EXTREMAL INDEX

Two mixing conditions on the samples **Cond-1**: Limits long range dependence $|P(\mathcal{AB}) - P(\mathcal{A})P(\mathcal{B})| \leq \alpha_n \quad \mathcal{A} \text{ and } \mathcal{B}: \text{ events } \subset \{X_i \leq u_n\}, l_n \text{ seperated}$ $l_n = o(n), \alpha_n \to 0$

Stationary Markov samples or its measurable functions satisfy this



PROPOSITION

If the sampled sequence is stationary and satisfies mixing conditions, then Extremal Index

$$\theta = <\mathbf{1}, \nabla C > \big|_{u=v=1} - 1,$$

 $0 \leq \theta \leq 1$ and $C(u, v) = P(X_1 \leq F^{-1}(u), X_2 \leq F^{-1}(v))$ is the Copula.

DEGREE CORRELATIONS

- Undirected and correlated
- $f(d_1, d_2)$ is enough to construct graph



- Crawling via Random Walks on vertices
- Degree sequence is a Hidden Markov chain
- What is the joint stationary distribution on degree state space?



Standard Random Walk



 $f_{RW}(d_{t+1}|d_t) \approx \frac{1}{d_t} \cdot \frac{E[D]f(d_t, d_{t+1})}{f_d(d_t)}$

$$f_{RW}(d_{t+1}, d_t) \approx f(d_{t+1}, d_t)$$

Page Rank



with c, follow RW with 1 - c, uniform node sampling

$$f_{PR}(d_{t+1}|d_t) \approx cf_{RW}(d_{t+1}|d_t) + (1-c)f_d(d_{t+1})$$

$$P(\text{head}) = c$$

$$c = \frac{d_t}{d_t + \alpha}$$

$$f_{RWJ}(d_{t+1}, d_t) \approx \frac{E[D]}{E[D] + \alpha} f(d_{t+1}, d_t) + \frac{\alpha}{E[D] + \alpha} f_d(d_{t+1}) f_d(d_t)$$

CHECK OF MEANFIELD MODEL IN RANDOM WALKS



Degree correlation among neighbours, bivariate Pareto distributed

$$\bar{F}(d_1, d_2) = \left(1 + \frac{d_1 - \mu}{\sigma} + \frac{d_2 - \mu}{\sigma}\right)^{-\gamma} \mu = 10, \sigma = 15, \gamma = 1.2$$

Extremal Index for Bivariate Pareto Model

$$\bar{F}(d_1, d_2) \sim \left(1 + \frac{d_1 - \mu}{\sigma} + \frac{d_2 - \mu}{\sigma}\right)^{-\gamma}$$

Random Walk:
$$EI = 1 - 1/2^{\gamma}$$

Random Walk with Jumps: $EI = 1 - \frac{E[D]}{E[D] + \alpha} 2^{\gamma}$
PageRank: $EI \ge (1 - c)$
(for any kind of degree correlations)

ESTIMATION OF EXTREMAL INDEX

Empirical Copula based estimator:

$$C_n(u,v) = \frac{1}{n} \sum_{i=1}^n \mathbb{I}\left(\frac{R_i^X}{n+1} \le u, \frac{R_i^Y}{n+1} \le v\right)$$

EI: slope at (1; 1),Linear least square fitting & numerical differentiation



Use $E(T_{\theta}^2) = 2/\theta$ to obtain estimates

NUMERICAL RESULTS: SYNTHETIC GRAPHS

Degree correlation between neighbours

$$\bar{F}(d_1, d_2) = \left(1 + \frac{d_1 - \mu}{\sigma} + \frac{d_2 - \mu}{\sigma}\right)^{-\gamma} \mu = 10, \sigma = 15, \gamma = 1.2$$

EI	Analysis	Copula based estimator	Intervals Estimator
Synthetic graph (5K Nodes)	0.56	0.53	0.58



NUMERICAL RESULTS: REAL GRAPHS

EI	Copula based estimator	Intervals Estimator
DBLP (32K Nodes, 1.1M Edges)	0.29	0.25
Enron Email (37K Nodes,368K Edges)	0.61	0.62

CONCLUSIONS

- Associated Extremal Value Theory of stationary sequence to sampling of large graphs
- For any general stationary samples meeting two mixing conditions, knowledge of bivariate distribution or bivariate copula is sufficient to derive many extremal properties
- Extremal Index (EI) encapsulates this relation
- Applications of EI to many relevant extrems:
 - First hitting time
 - Order statistics
 - Mean cluster size
- Modeled correlation in degrees of adjacent nodes and random walk in degree state space
- Estimates of EI for synthetic graph with degree correlations and find a good match with theory
- Estimated EI for two real world networks

Thank You!