



Introduction

- Given a large social graph, like a scientific collaboration network, what can we say about its robustness?
- Can we estimate a robustness index for a graph quickly?
- If the graph evolves over time, how these properties change?
- Robust Graph:** Capable to retain its structure and its connectivity properties after the loss of a portion of its nodes and edges
- The property of robustness in real-world graphs is closely related to the notion of **community structure**
- We tackle the problem of estimating the robustness of a graph quickly, studying the **expansion properties**
- Contributions:**
 - Fast robustness index
 - Patterns of real static and time-evolving social graphs
 - Anomaly detection

Preliminaries: Expansion Properties

- Good expander: Simultaneously sparse and highly connected
- Given a graph $G = (V, E)$, the *expansion* of any subset of nodes $S \subset V$, with size at most $\frac{|V|}{2}$, is defined as $\frac{|N(S)|}{|S|}$
- A graph is considered to have good expansion properties if every subset of nodes has good expansion (i.e., many neighbors)

Expansion, Robustness and Community Structure

- Why the expansion properties of a graph are important?
 - They offer crucial insights about the structure of a graph
 - They can act as a natural measure of the graph's robustness
 - Information about the presence or not of edges which can operate as bottlenecks inside the network
- Good expansion properties \rightarrow high robustness, while poor expansibility reflects exactly the opposite behavior
- Connections with the community structure: good expansibility requires cuts with large size (i.e., large number of edges crossing the cut) \rightarrow poor community structure
- The expansion properties of a graph can be approximated by the *spectral gap* $\Delta\lambda = \lambda_1 - \lambda_2$ of the adjacency matrix A
- Large $\Delta\lambda$ implies high robustness
 - However**, it is not clear how large the spectral gap should be

Spectral Gap + Subgraph Centrality

- Combine the spectral gap with the subgraph centrality [Estrada, Eur. Phys. J. B '06]
- Subgraph Centrality:** # of closed walks that a node participates

$$SC(i) = \sum_{j=1}^{|V|} u_{ij}^2 \sinh(\lambda_j), \forall i \in V$$
- Good expansion properties \rightarrow High robustness $\rightarrow \lambda_1 \gg \lambda_2 \rightarrow u_{i1}^2 \sinh(\lambda_1) \gg \sum_{j=2}^{|V|} u_{ij}^2 \sinh(\lambda_j)$
 - $SC(i) \approx u_{i1}^2 \sinh(\lambda_1), \forall i \in V \rightsquigarrow u_{i1} \propto \sinh^{-1/2}(\lambda_1) SC(i)^{1/2}$
 - Deviation from this behavior \rightarrow existence (or lack thereof) of high robustness properties
- Shortcoming:**
 - Scalability issues (it requires all the pairs $(\lambda_i, u_i), \forall i \in V$)
 - It cannot be applied directly to bipartite graphs

Proposed Metric: Generalized Robustness Index r_k

Q: Can we efficiently approximate the *SC* of every node in the graph?

- The eigenvalues of A follow a power-law distribution [Faloutsos et al., SIGCOMM '99]
- The eigenvalues are almost symmetric around zero (except from the first few) [Tsourakakis, ICDM '08]
- $\sinh(\cdot)$: odd function (i.e., $\sinh(-x) = -\sinh(x)$)

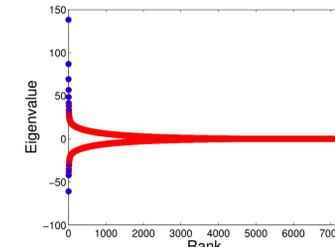


Figure: Skewed spectrum (WIKI-VOTE)

\rightsquigarrow Approximate the *SC* using **only the first top k** eigenvalues and their corresponding eigenvectors

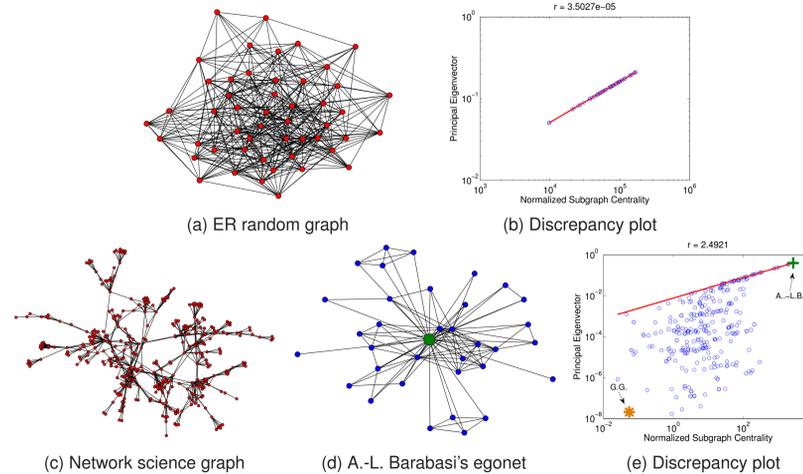
$$NSC_k(i) = \sum_{j=1}^k u_{ij}^2 \sinh(\lambda_j), \forall i \in V$$

□ Generally, $k \ll |V|$ for real-world graphs

$$\rightsquigarrow r_k = \left(\frac{1}{|V|} \sum_{i=1}^{|V|} \left\{ \log(u_{i1}) - \left(\log(\sinh^{-1/2}(\lambda_1)) + \frac{1}{2} \log(NSC_k(i)) \right) \right\}^2 \right)^{1/2}$$

□ Summarizes the robustness of a graph in a single number (smaller $r_k \rightarrow$ better robustness)

An Illustrative Example: Random vs. Real Graphs



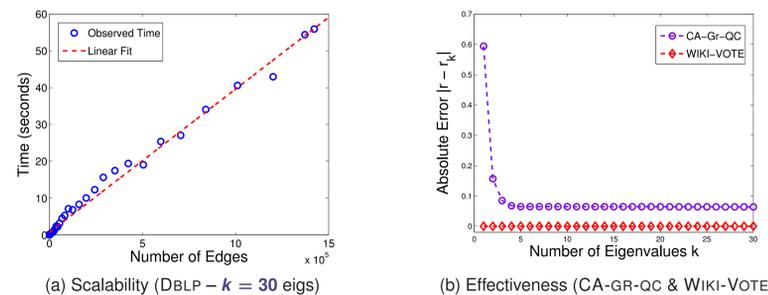
Datasets

Graph	# Nodes	# Edges	Graph	# Nodes	# Edges
EPINIONS	75,877	405,739	CA-ASTRO-PH	17,903	197,031
EMAIL-EUALL	224,832	340,795	CA-GR-QC	4,158	13,428
SLASHDOT	77,360	546,487	CA-HEP-TH	8,638	24,827
WIKI-VOTE	7,066	100,736	DBLP	404,892	1,422,263
FACEBOOK	63,392	816,886	CIT-HEP-TH	26,084	334,091
YOUTUBE	1,134,890	2,987,624			



Effectiveness and Scalability of r_k index

Q: How effective and scalable (efficient) is the proposed r_k index?



- The r_k index scales linearly with respect to the number of edges
- Only a few eigenpairs are enough to achieve a very good approximation of the robustness index

Robustness of Large Static Graphs: High Robustness Pattern

Q: What can we say about the robustness of large social graphs?

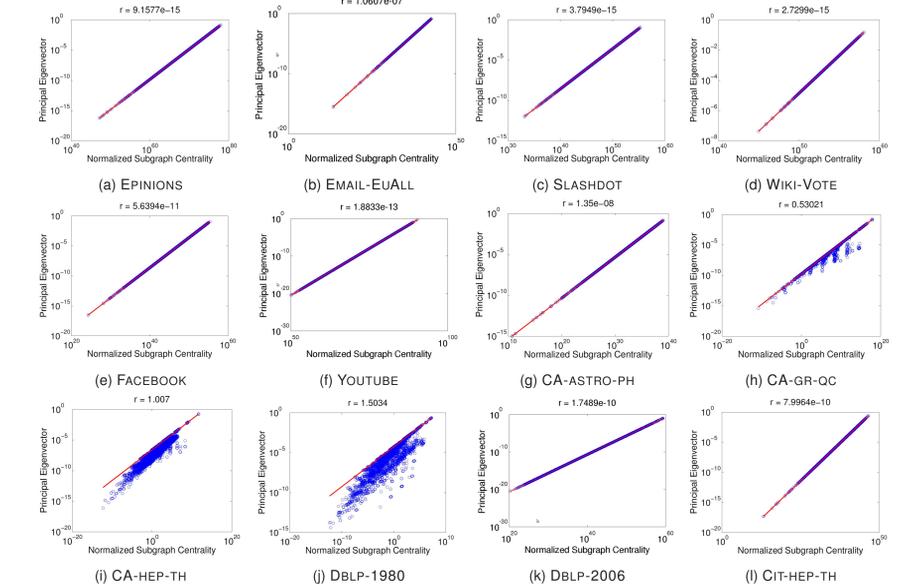


Figure: **Discrepancy plots:** Principal eigenvector vs. *NSC* in log-log scales. Almost all graphs tend to be extremely robust (linearity, r_k close to zero)

Large Social Graphs: Good expansion properties \rightarrow High robustness \rightarrow Not a clear modular structure

Time Evolving Graphs: Fragility Evolution Pattern

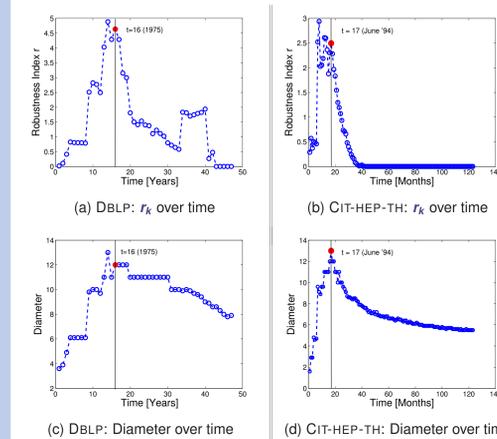


Figure: Fragility evolution pattern: The spike of the r_k index aligns with the diameter's spike

Q: How the robustness index r_k of a graph changes over time?

- Study the *fragility evolution* of a graph

General Observation:

- At the first time points $r_k \nearrow$ gradually \rightarrow Low robustness \rightarrow Good community structure
- After a specific time point, r_k starts \searrow gradually \rightarrow The graphs tend to become more robust
- The time point that r_k changes, corresponds to the **gelling point** [McGlohon et al., KDD '08]

The fragility evolution pattern can be considered as a natural explanation for the structural differences (regarding robustness and community structure) between different scale graphs

Anomaly Detection

Q: Can we spot anomalies over time using the r_k index?

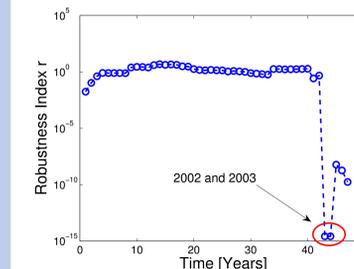


Figure: Fragility evolution of the DBLP graph (lin-log scales)

- Examine the r_k index over time, trying to identify and track abrupt changes and deviations
- Sudden deviations from the fragility evolution pattern can possibly correspond to anomalies
- DBLP graph: Strange behavior of the r_k index for 2002 and 2003
- These two time graphs are outliers
 - After 2001 a large number of new publications were introduced \rightarrow Robustness $\nearrow \rightarrow r_k \searrow$
 - After 2002-2003 new research fields are covered from DBLP \rightarrow New fields formed new communities $\rightarrow r_k \nearrow$

Acknowledgements

Data: SNAP, MPI-SWS, NEU, DBLP

Funding: NSF (Grant No. IIS1017415), ARL (No. W911NF-09-2-0053)