#### Using Partial Probes to Infer Network States

#### Pavan Rangudu°, Bijaya Adhikari\*, B. Aditya Prakash\*, Anil Vullikanti\* °

\*Department of Computer Science, Virginia Tech °NDSSL, Biocomplexity Institute, Virginia Tech

Contact: badityap@cs.vt.edu



イロト 不得 トイヨト イヨト ヨー ろくぐ

# Motivation

- Network nodes and links fail dynamically
- Networks not known fully because of privacy constraints
- Our focus: if some failed nodes are known, can we infer the states of the remaining nodes?



Node failures in internet



Traffic jam in road network

Prior works fail to the address the problem directly.

# Our model

- Graph G(V, E) with set  $I \subseteq V$  which have failed
- Goegraphically correlated failure model [Agarwal et al., 2013]
  - Single seed of the failure, with probability  $p_s(v)$  of node v being the seed
  - Correlated failure model: F(u|v) denotes the probability that node u fails given that v has failed
    - Assume independence, i.e.,  $F(u_1, u_2|v) = F(u_1|v) \cdot F(u_2|v)$
  - Motivation: attacks or natural disasters in infrastructure networks
- Probes: subset  $Q \subseteq I$  of failed nodes is known
- Objective: find the set I Q



Figure: A toy road network with node failures

イロト 不得 トイヨト イヨト ヨー ろくぐ

## Our approach: Minimum Description Length

• Model cost  $\mathcal{L}(|\mathcal{Q}|, |I|, I)$  has three components

$$\mathcal{L}(|\mathcal{Q}|,|I|,I) = \mathcal{L}(|\mathcal{Q}|) + \mathcal{L}(|I| \mid |\mathcal{Q}|) + \mathcal{L}(I \mid |\mathcal{Q}|,|I|).$$

• 
$$\mathcal{L}(|\mathcal{Q}|) = -\log\left(Pr(|\mathcal{Q}|)\right)$$
 by using the Shannon-Fano code  
•  $\mathcal{L}(|I| \mid |\mathcal{Q}|) = -\log\left(\frac{Pr(|\mathcal{Q}| \mid |I|)Pr(|I|)}{Pr(|\mathcal{Q}|)}\right)$   
•  $\mathcal{L}(I \mid |\mathcal{Q}|, |I|) = -\log\left(Pr(I \mid |\mathcal{Q}|, |I|)\right) = -\log\left(Pr(I \mid |I|)\right)$ 

◆□▶ ◆□▶ ◆三▶ ◆三▶ 三三 のへぐ

Data cost: description of Q<sup>+</sup> = I \ Q (assuming no observation errors)

• 
$$\mathcal{L}(\mathcal{Q}^+|I) = -\log\left(\gamma^{|\mathcal{Q}|}(1-\gamma)^{|\mathcal{Q}^+|}\right) = -|\mathcal{Q}|\log(\gamma) - (|I| - |\mathcal{Q}|)\log(1-\gamma)$$

### Problem Description

Model Cost

$$\begin{split} \mathcal{L}(|\mathcal{Q}|,|I|,I) = & \mathcal{L}(|\mathcal{Q}|) + \mathcal{L}(|I| \mid |\mathcal{Q}|) + \mathcal{L}(I \mid |\mathcal{Q}|,|I|) \\ = & -\log \binom{|I|}{|\mathcal{Q}|} - |\mathcal{Q}|\log(\gamma) - (|I| - |\mathcal{Q}|)\log(1 - \gamma) \\ & -\log \Big(\sum_{s \in V} p_s(s) \prod_{v \in I} F(v \mid s) \prod_{v' \notin I} \Big(1 - F(v' \mid s)\Big)\Big) \\ & * \text{after algebra} \end{split}$$

#### **Problem Formulation**

Given G,  $p_s$ ,  $F(\cdot)$ , Q, find I that minimizes the total MDL cost:

$$\mathcal{L}(|\mathcal{Q}|, |I|, I, \mathcal{Q}) = -\log \binom{|I|}{|\mathcal{Q}|} - \log \left( \sum_{s \in V} p_s(s) \prod_{v \in I} F(v \mid s) \prod_{v' \notin I} \left( 1 - F(v' \mid s) \right) \right)$$
$$-2|\mathcal{Q}|\log(\gamma) - 2(|I| - |\mathcal{Q}|)\log(1 - \gamma)$$

◆□▶ ◆□▶ ◆三▶ ◆三▶ 三三 のへぐ

## Algorithm GREEDY

**Input:** Instance  $(V, Q, p, P, \gamma)$ **Output:** Solution  $\hat{I}$  that minimizes  $\mathcal{L}(|\mathcal{Q}|, |\hat{I}|, \hat{I}, \mathcal{Q})$ 1: for each  $s \in V$  do for each  $k \in [|\mathcal{Q}|, |V|]$  do 2:  $I_{s}(k) \leftarrow \text{Top } k - |\mathcal{Q}| \text{ nodes in } V \setminus \mathcal{Q} \text{ with highest weight}$ 3: f(s, v) $I_{s}(k) \leftarrow I_{s}(k) \cup \mathcal{Q}$ 4: end for 5: 6: end for 7:  $\mathcal{S} \leftarrow \{I_s(k) : \forall s \in V \& k \in [|\mathcal{Q}|, |V|]\}$ 8:  $\hat{I} \leftarrow \arg \min \mathcal{L}(|\mathcal{Q}|, |I|, I, \mathcal{Q})$  $I \in S$ 9: Return

◆□▶ ◆□▶ ◆ □▶ ◆ □▶ ○ □ ○ ○ ○ ○

## Analysis of $\operatorname{GREEDY}$

#### Theorem: (Additive Approximation)

Let  $I^*$  be the set minimizing the MDL cost, and let I denote the solution computed by Algorithm GREEDY. Then,  $\mathcal{L}(|\mathcal{Q}|, |I|, I, \mathcal{Q}) \leq \mathcal{L}(|\mathcal{Q}|, |I^*|, I^*, \mathcal{Q}) + \log(n)$ , where n is the number of seed nodes.

イロト 不得 トイヨト イヨト ヨー ろくぐ

#### Running time

Algorithm GREEDY runs in  $O(|V|^3)$  time

# Experiments

- Baseline: local improvement algorithm LOCALSEARCH
- Datasets
  - Synthetic grid
    - $60 \times 60$  grid
    - Uniform seed probability  $p_s(\cdot)$
    - Conditional failure probability distribution using model of [Agarwal et al., 2013]: F(v | s) = 1 d(s, v), where  $d(\cdot)$  is (normalized) distance
  - Real datasets: Seed and conditional failure probability distributions computed from data
    - JAM data from WAZE for Boston: road network with 2650 nodes.
    - WEATHER data from WAZE for Boston: road network with 1520 nodes.
    - POWER-GRID: network of 24 nodes from Electric disturbance events

WAZE dataset



Visualization of Waze dataset. Partitions in the  $119 \times 78$  grid represent nodes in our network.

# Takeaways



- Our MDL based approach helps identify missing failures
- Promising approach for other problems with missing information



◆□▶ ◆□▶ ◆□▶ ◆□▶ □ のQ@