## Algorithmic Problems in Epidemiology

#### V.S. Anil Kumar

#### Dept. of Computer Science and Virginia Bioinformatics Institute Virginia Tech

Joint work with Stephen Eubank and Madhav Marathe (Virginia Tech), Aravind Srinivasan and Nan Wang (U. Maryland, College Park)

> Email: akumar@vbi.vt.edu Web: http://ndssl.vbi.vt.edu/





## The first Network Approach to Epidemiology



"Riddle of the Cholera Outbreak" solved by Dr. John Snow Made an implicit network between people and water pumps
People who drank water from Broad St. pump died
also proved that cholera spreads

through water





## *EpiSims/Simdemics*<sup>1</sup> - a network approach to epidemiology

Main steps:

1. Construct a realistic social network

Not readily available

2. Develop an efficient disease simulation

- Must scale to millions of people
- Static analysis to guide simulation
- 3. Help in policy questions
  - Abstract out a simpler combinatorial model for the disease dynamics
  - whom to vaccinate/quarantine
  - where to place sensors





Main steps:

- Synthetic population from census data
- Location Information from census, land use, GIS, etc.
- Individual activities from surveys, traffic data, behavioral models, etc.



Main Difficulty: Slow and expensive Can we construct graphs that look "like" real social networks? Common approach: generative random graph models





### The Social Network







#### Generative Graph Models: Chung-Lu model locations people $Pr[(p_2, l_1)] = d(p_2)d(l_1)/\sigma$ p<sub>1</sub> $d(I_1)$ $d(p_1)$ $d(l_2)$ $d(p_2)$ Degree distribution Degree distribution of people of locations 0.25 Portland data $10^{-1}$ Portland data Chung-Lu model $\sigma = \sum_{p \in P} d(p) = \sum_{l \in L} d(l)$ FastGen-1 Chung-Lu model Fractions of people 10<sup>-2</sup> Fraction of locations Degree distributions of people Degree distribution of locations (log–log) $\max_{p \in P, l \in L} d(p) d(l) \le \sigma$ $10^{-4}$ 10<sup>-5</sup> <sup>\_⊥</sup> 10<sup>−6</sup> 10<sup>4</sup> Location degrees 10<sup>0</sup> 10<sup>1</sup> 10<sup>3</sup> People degrees 25 30



Network Dynamics and Simulation Science Laboratory



## Connectivity in the people-people graph

Assume: d(p)=c, for each  $p \in P$ 10 Fractions of people  $d_0 \leq d(I) \leq d_1$ , for each  $d(I) \in L$  $n_i$ : #locations with degree i ~ c'|L|/i^{\beta} For  $p_1, p_2 \in P$ :  $\Pr[(p_1, p_2) \notin G_P] = \prod_{\ell \in L} (1 - \frac{d(p_1)d(p_2)d(\ell)^2}{r^2})$ 10  $\leq exp(-\sum_{\ell \in I} \frac{c^2 d(\ell)^2}{\sigma^2})$  $= exp(-\frac{1}{|P|^2}\sum \frac{i^2n_i}{i^\beta})$  $\simeq exp(-rac{c'|L|}{|P|^2} \sum_{i=1}^{d_1} rac{1}{i^{\beta-2}})$  $\simeq exp(-\frac{1}{|P|}\frac{c(\beta-2)d_0^{\beta-2}d_1^{3-\beta}}{3-\beta})$  $\ll exp(-\frac{2}{|P|}) \le 1 - \frac{1}{|P|}$ 









## The Configuration Model







## The Configuration Model







## Chung-Lu vs Configuration

![](_page_9_Figure_1.jpeg)

Candidates 
$$\begin{aligned} f(dist(p_2,l_1)) &= dist(p_2,l_1)^{\beta} \\ f(dist(p_2,l_1)) &= e^{a \cdot dist(p_2,l_1)} \end{aligned}$$

![](_page_9_Picture_3.jpeg)

![](_page_9_Picture_5.jpeg)

### Fast Generation of CL

- Speed up (naive CL takes ~ 10 hours) random generation
- Preserve degrees
- Be "similar" to CL

![](_page_10_Picture_4.jpeg)

![](_page_10_Picture_5.jpeg)

- Speed up (naive CL takes ~ 10 hours) random generation
- Preserve degrees
- Be "similar" to CL

$$\begin{split} \forall p \in P, \forall \ell \in L, \Pr[X_{p,\ell} = 1] = \frac{d(p) \cdot d(\ell)}{\sigma} \\ \forall \ell \in L, \Pr[|\{p : X_{p,\ell} = 1\}| = d(\ell)] = 1 \\ \forall \ell \in L, P' \subseteq P : \\ \\ \mathsf{Negative}_{\mbox{correlation}} \\ \begin{cases} \Pr[\bigwedge_{p \in P'} (X_{p,\ell} = 0)] \leq \prod_{p \in P'} \Pr[X_{p,\ell} = 0], \\ \Pr[\bigwedge_{p \in P'} (X_{p,\ell} = 1)] \leq \prod_{p \in P'} \Pr[X_{p,\ell} = 1]. \end{cases} \end{split}$$

![](_page_11_Picture_5.jpeg)

![](_page_11_Picture_7.jpeg)

## Dependent Rounding<sup>1</sup>

![](_page_12_Figure_1.jpeg)

Time taken: O( # fractional variables x(p,l) ) $\Rightarrow O(|P| |L|)$  time overall

<sup>1</sup>R. Gandhi, S. Khuller, S. Parthasarathy, A. Srinivasan, Dependent Rounding, FOCS 2004

![](_page_12_Picture_4.jpeg)

![](_page_12_Picture_6.jpeg)

#### Graph generation by Dependent Rounding $x(p,l) + \alpha$ β $\alpha + \beta$ р *x(p,l)* $x(p',l) - \alpha$ / **x(p,l) -** β **)** p' $\alpha$ *x(p',l)* $\alpha + \beta$ $x(p',l) + \beta$

 $\alpha = \min\{ 1- x(p,l), x(p',l) \}$  $\beta = \min\{ x(p,l), 1- x(p',l) \}$ 

![](_page_13_Picture_2.jpeg)

![](_page_13_Picture_3.jpeg)

![](_page_14_Figure_0.jpeg)

![](_page_14_Picture_1.jpeg)

![](_page_14_Picture_2.jpeg)

![](_page_14_Picture_3.jpeg)

![](_page_15_Figure_0.jpeg)

![](_page_15_Picture_1.jpeg)

![](_page_15_Picture_2.jpeg)

![](_page_15_Picture_3.jpeg)

![](_page_16_Figure_0.jpeg)

![](_page_16_Picture_1.jpeg)

![](_page_16_Picture_3.jpeg)

## Fast Generation of CL

![](_page_17_Figure_1.jpeg)

- 1. Partition into maximal blocks  $B_1, B_2, ..., B_q$ 
  - $B_k = \{p(j_k), \dots, p(j_{k+1})-1\}$
  - $w(B_k) = \sum_{j \in Bk} d(p_j) d(l) / \sigma$
  - $1 d(p(j_{k+1}))d(l) / \sigma < w(B_k) \le 1$
- 2. Reduce to smaller graph H on  $\{I, B_1, B_2, ..., B_q\}$
- 3. Run Dependent rounding on H and choose subset  $B' \subseteq \{B_1, B_2, ..., B_q\}$
- 4. For each block  $B_i \in B'$ , choose  $p \in B_i$  with prob.  $x(p,l) / w(B_i)$

![](_page_17_Picture_9.jpeg)

![](_page_17_Picture_11.jpeg)

Observation:  $q \le 2(d(l) + 1)$ Proof:  $\forall k:$  $w(B_k) + w(B_{k+1}) > 1 - d(p(j_{k+1}))d(\ell)/\sigma + w(B_{k+1})$  $= 1 + w(B_{k+1}) - d(p(j_{k+1}))d(\ell)/\sigma$ > 1 $\Rightarrow \lfloor \frac{q}{2} \rfloor < \sum^{q} w(B_k)$ k=1 $\Rightarrow q \leq 2(d(\ell) + 1)$ 

![](_page_18_Picture_2.jpeg)

![](_page_18_Picture_4.jpeg)

# Fast Generation of CL: Running

#### Time

![](_page_19_Figure_2.jpeg)

 $\Rightarrow O(\Sigma q(I) \log |P|) \text{ time overall}$ =  $O(\Sigma d(I) \log |P|) = O(|E| \log |P|)$ 

![](_page_19_Picture_4.jpeg)

![](_page_19_Picture_5.jpeg)

## Fast Generation Properties

- Equivalent to running Dependent Rounding on  $G(\{I\},P, \{(p,I): p \in P\})$ with fractional solution:  $x(p,I) = d(p)d(I) / \sigma$
- 1. Partition into maximal blocks  $B_1$ ,  $B_2$ , ...,  $B_q$ 
  - $B_k = \{p(j_k), \dots, p(j_{k+1})-1\}$
  - 1 d(p(j<sub>k+1</sub>))d(l)/  $\sigma < w(B_k) \le 1$
- 2. Reduce to smaller graph H on  $\{I, B_1, B_2, ..., B_q\}$
- 3. Run Dependent rounding on H and choose subset  $B' \subseteq \{B_1, B_2, ..., B_q\}$
- 4. For each block  $B_i \in B$ , choose  $p \in B_i$ with prob.  $x(p,l) / w(B_i)$
- ⇒ 1. Pr[edge (p,l)] = x(p,l) = d(p)d(l)/σ
   2. Rounded-degree(l) = Σ<sub>p∈P</sub> x(p,l) = d(l)
   3. Negative correlation property

![](_page_20_Picture_9.jpeg)

![](_page_20_Picture_10.jpeg)

#### FastGen

	Portland data	CL	FastGen-1	FastGen-2
number of locations	181230	178746~(98.63%)	$181230\;(100\%)$	178668~(98.59%)
size of giant(locs.)	181192	178571	181088	178611
number of people	1615860	1507234~(93.28%)	$1507291 \ (93.28\%)$	$1615860\ (100\%)$
size of giant(ppl.)	1615813	1507054	1507148	1615803
number of edges	6060679	6065637~(100.08%)	6060679~(100%)	6060679~(100%)
average deg. (ppl.)	3.7507	4.0227	4.0209	3.7507
time of generating a graph		> 10  hours	< 40 seconds	< 30 seconds

![](_page_21_Figure_2.jpeg)

	Portland data	$\operatorname{CL}$	FastGen-1	FastGen-2
c.c. mean	0.6376	0.6161	0.6235	0.7021
c.c. std.	0.0167	0.0315	0.0236	0.0201
c.c l.b.	0.6286	0.6220	0.6241	0.7061
c.c. a.l.b.	0.6262	0.6220	0.6240	0.7061

![](_page_21_Picture_4.jpeg)

![](_page_21_Picture_6.jpeg)

#### Main steps:

1. Construct a realistic social network

- Not readily available
- 2. Develop an efficient disease simulation
  - Must scale to millions of people
  - Static analysis to guide simulation
- 3. Help in policy questions
  - Abstract out a simpler combinatorial model for the disease dynamics
  - whom to vaccinate/quarantine
  - where to place sensors

![](_page_22_Picture_11.jpeg)

![](_page_22_Picture_13.jpeg)

## Policy Planning problems

- Disease detection
  - Fast detection + response, instead of mass vaccination
  - Sensor Location: dominating set problem
- Quarantining
  - Remove some edges so that disease is contained
- Vaccination
  - Remove some nodes so that disease is contained

![](_page_23_Picture_8.jpeg)

![](_page_23_Picture_9.jpeg)

## The Sensor Placement Problem

![](_page_24_Figure_1.jpeg)

Dominating set problem: choose  $L' \subseteq L \text{ s.t. } N(L') = P$ 

(1- $\varepsilon$ )-Dominating set problem: choose L'  $\subseteq$  L s.t.  $|N(L')| \ge (1-\varepsilon)|P|$ 

![](_page_24_Picture_4.jpeg)

![](_page_24_Picture_5.jpeg)

## Temporal Dominating Set

![](_page_25_Figure_1.jpeg)

I(p, l): interval during which p visits l

The Temporal Dominating set problem: choose  $L' \subseteq L$ , and time interval  $[t_1, t_2]$  for each  $l \in L'$ s.t.  $\forall p \in P, \exists l \in L'$  s.t.  $I(p,l) \cap [t_1, t_2] \neq \emptyset$ 

![](_page_25_Picture_4.jpeg)

![](_page_25_Picture_6.jpeg)

- FastGreedy: choose large locations in non-increasing order of degrees sum of degrees is  $(1-\epsilon')|P|$
- FastGreedy gives (1+o(1))-approximation to optimum dominating set in Chung-Lu model
  - Same approximation in FastGeneration model
- FastGreedy works well in practice
  - 10% of the locations can dominate ~90% of people in Portland network
  - Very close to Greedy
  - Takes ~15 sec
- Temporal version hard to approximate within  $\Omega(n^{\epsilon})$

![](_page_26_Picture_9.jpeg)

![](_page_26_Picture_10.jpeg)

## Dominating Set in the Chung-Lu model

Theorem: FastGreedy is a (1+o(1))-approximation to the (1- $\epsilon$ )-Dominating Set in the Chung-Lu and Fast Generation models

![](_page_27_Figure_2.jpeg)

Non-increasing order of degrees

![](_page_27_Picture_4.jpeg)

![](_page_27_Picture_5.jpeg)

![](_page_27_Picture_6.jpeg)

## Dominating Set in the Chung-Lu model

Theorem: FastGreedy is a (1+o(1))-approximation to the (1- $\epsilon$ )-Dominating Set in the Chung-Lu and Fast Generation models

![](_page_28_Figure_2.jpeg)

- $L(d) = \{ \text{locations } | \text{ with } d(l) \ge d \}$   $S(L') = \sum_{l \in L'} d(l)$   $d_2: \text{ largest degree s.t. } S(L(d_2)) \ge (1-\epsilon')|P|$ To prove:  $1. |N(L(d_2))| \ge (1-\epsilon)|P|$
- 2. Approximation ratio =  $|L(d_2)|/|OPT| = 1+o(1)$

![](_page_28_Picture_5.jpeg)

![](_page_28_Picture_7.jpeg)

## Property 1

![](_page_29_Picture_1.jpeg)

Chernoff bound, also holds for Fast Generation

Lemma: 
$$S(L') \ge |P| \ln \frac{1 - \delta_1}{\epsilon_1} \Rightarrow |N(L')| \ge (1 - \epsilon_1 - \delta_1)|P|, w.h.p.$$
  
Proof:  
 $\forall p \in P, \Pr[p \notin N(L')] = \prod_{\ell \in L'} (1 - \frac{d(p)d(\ell)}{\sigma})$   
 $< e^{-S(L')/|P|}$   
 $\Rightarrow E[|P \setminus N(L')|] < |P|e^{-S(L')/|P|}$ 

$$S(L') \ge |P| \ln \frac{1 - \delta_1}{\epsilon_1} \Rightarrow E[|P \setminus N(L')|] < |P| \frac{\epsilon_1}{1 - \delta_1}$$
  
$$\Rightarrow E[|N(L')|] > |P|(1 - \frac{\epsilon_1}{1 - \delta_1})$$
  
$$\Rightarrow |N(L')| > (1 - \epsilon_1 - \delta_1)|P| \quad w.h.p.$$

 $\Rightarrow$  Fast Greedy gives a (1- $\varepsilon$ )-Dominating Set

![](_page_29_Picture_6.jpeg)

![](_page_29_Picture_7.jpeg)

![](_page_30_Figure_1.jpeg)

$$\text{Lemma: } N(L') \geq (1-\epsilon)|P| \ \Rightarrow \ S(L') \geq |P| \frac{\ln \frac{2-\epsilon}{2\epsilon}}{1+\gamma}, \ w.h.p.$$

∴ S(OPT) is large d<sub>3</sub> = largest degree s.t. S(L(d<sub>3</sub>)) ≤  $|P| \frac{\ln \frac{2-\epsilon}{2\epsilon}}{1+\gamma}$ 

 $\Rightarrow |\mathsf{OPT}| \geqq |\mathsf{L}(\mathsf{d}_3)|$ 

![](_page_30_Picture_5.jpeg)

![](_page_30_Picture_7.jpeg)

## Approximation factor of FastGreedy

Approximation factor 
$$\leq rac{|L(d_2)|}{|L(d_3)|}$$

Power law 
$$\Rightarrow S(L(d)) \simeq \frac{c'|L|}{(\beta - 2)d^{\beta - 2}}$$

 $\therefore$  Approximation factor

$$\leq \left(\frac{S(L(d_2))}{S(L(d_3))}\right)^{\frac{\beta-1}{\beta-2}}$$

≈ 1+o(1), w.h.p.

![](_page_31_Picture_6.jpeg)

![](_page_31_Picture_7.jpeg)

## Policy Planning problems

- Disease detection
  - Fast detection + response, instead of mass vaccination
  - Sensor Location: dominating set problem
- Quarantining
  - Remove some edges so that disease is contained
- Vaccination
  - Remove some nodes so that disease is contained

![](_page_32_Picture_8.jpeg)

![](_page_32_Picture_9.jpeg)

### Policy problems: whom to vaccinate

- Given a social network G(P,E)
  - $\boldsymbol{\cdot}$  initial infected set A
  - budget B on # vaccinations

![](_page_33_Figure_4.jpeg)

![](_page_33_Picture_5.jpeg)

![](_page_33_Picture_6.jpeg)

#### Policy problems: whom to vaccinate

![](_page_34_Figure_1.jpeg)

![](_page_34_Picture_2.jpeg)

### Policy problems: whom to vaccinate

- Given a social network G(P,E)
  - $\boldsymbol{\cdot}$  initial infected set A
  - budget B on # vaccinations
- Goal choose  $S \subseteq P$  to vaccinate so that  $|S| \le B$ , and # nodes reachable from A in  $G[P \setminus S]$  is minimized
- **Result** Bicriteria approximation: Vaccinate  $(1+\epsilon)$ B nodes so that at most  $(1+1/\epsilon)$  OPT infected people<sup>1,2</sup>

![](_page_35_Figure_6.jpeg)

<sup>1</sup>S. Eubank, V.S. Anil Kumar, M. Marathe, A. Srinivasan and N. Wang, AMS-DIMACS special volume, 2005 <sup>2</sup>A. Hyrapetyan, D. Kempe, M. Pal and Z. Svitkina, ESA 2005

![](_page_35_Picture_8.jpeg)

![](_page_35_Picture_10.jpeg)

In reality: initial infected set is not known

Goal: Choose  $S \subseteq P$ , s.t. every component in  $G[P \setminus S]$  is small

 $\rho\mbox{-separator}$  problem: polylog approximations known, but algorithms are not scalable

![](_page_36_Picture_4.jpeg)

![](_page_36_Picture_5.jpeg)

## The "High Degree" Vaccination Policy

Vaccinate people of high degrees

![](_page_37_Figure_2.jpeg)

Giant component remains after deleting a large fraction of nodes

Also true for  $G_P$  in the Chung-Lu and FastGen models

![](_page_37_Picture_5.jpeg)

![](_page_37_Picture_6.jpeg)

![](_page_37_Picture_7.jpeg)

- Social Networks
  - Models for temporal and demographic aspects
  - Distance function
- Other network design problems for policy planning
  - Remove nodes/edges so that expansion is reduced
  - Remove nodes/edges so that average distance increases
  - Faster algorithms on large graphs
- More realistic disease models (SIR, SIS)
  - Dynamics on social networks
  - Policy planning problems
- Problems on temporal graphs

![](_page_38_Picture_12.jpeg)

![](_page_38_Picture_13.jpeg)

## Thank You

![](_page_39_Picture_1.jpeg)

![](_page_39_Picture_2.jpeg)