

# Probabilistic Contagion and Models of Influence

CS224W: Machine Learning with Graphs  
Jure Leskovec, Stanford University  
<http://cs224w.stanford.edu>

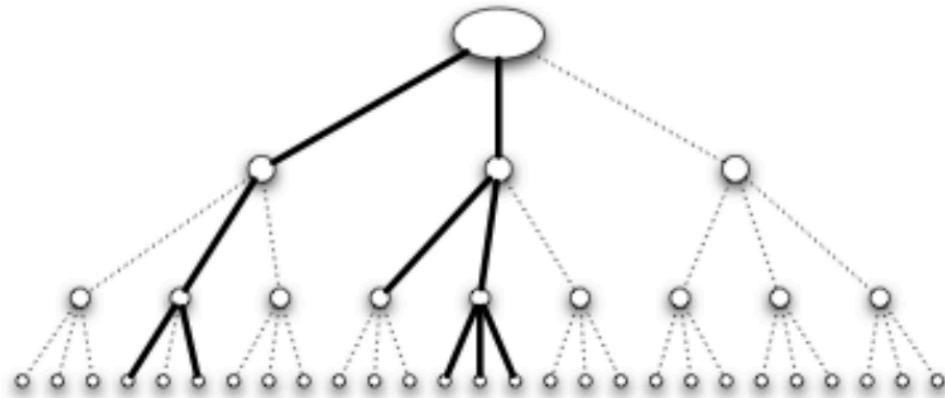


# Epidemics vs Cascade Spreading

- In decision-based models nodes make decisions based on pay-off benefits of adopting one strategy or the other
- In epidemic spreading:
  - Lack of decision making
  - Process of contagion is complex and unobservable
    - In some cases it involves (or can be modeled as) randomness

Recap

## Example with k=3



High contagion probability:  
The disease spreads

## Low contagion probability: The disease dies out



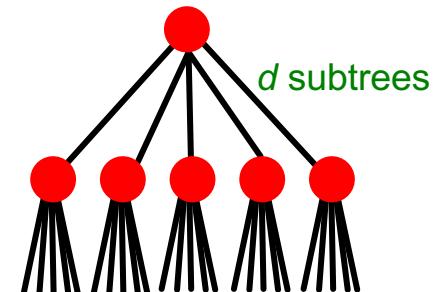
# Probabilistic Spreading Models

## ■ Epidemic Model based on Random Trees

- (a variant of branching processes)
- A patient meets  $d$  new people
- With probability  $q > 0$  she infects each of them

■ Q: For which values of  $d$  and  $q$  does the epidemic run forever?

Root node,  
“patient 0”  
Start of epidemic



■ Run forever:  $\lim_{h \rightarrow \infty} P \left[ \begin{array}{c} \text{a node at depth } h \\ \text{is infected} \end{array} \right] > 0$

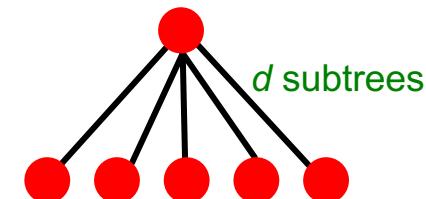
■ Die out:  $\lim_{h \rightarrow \infty} P \left[ \begin{array}{c} \text{a node at depth } h \\ \text{is infected} \end{array} \right] = 0$

# Probabilistic Spreading Models

- $p_h$  = prob. a node at depth  $h$  is infected
- We need:  $\lim_{h \rightarrow \infty} p_h = ?$  (based on  $q$  and  $d$ )
  - We are reasoning about a behavior at the root of the tree. Once we get a level out, we are left with identical problem of depth  $h - 1$ .

- Need recurrence for  $p_h$

$$p_h = 1 - \underbrace{(1 - q \cdot p_{h-1})^d}_{\text{No infected node at depth } h \text{ from the root}}$$



- $\lim_{h \rightarrow \infty} p_h$  = result of iterating

$$f(x) = 1 - (1 - q \cdot x)^d$$

- Starting at the root:  $x = 1$  (since  $p_1 = 1$ )

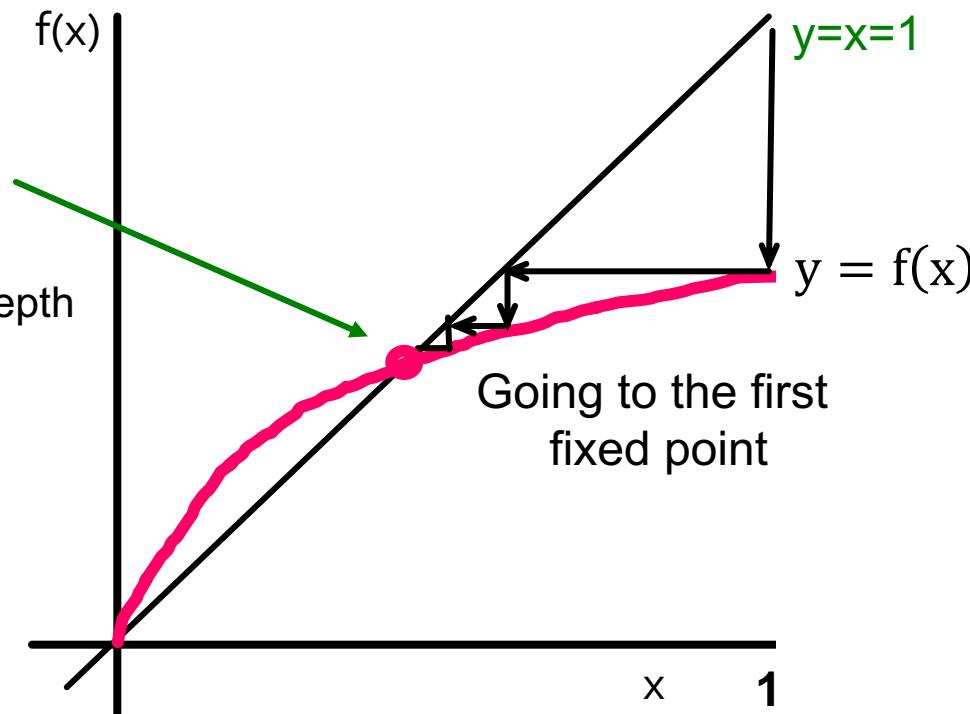
We iterate:  
 $x_1 = f(1)$   
 $x_2 = f(x_1)$   
 $x_3 = f(x_2)$

# Fixed Point: $f(x) = 1 - (1 - qx)^d$

**Fixed point:**

$$f(x) = x$$

This means that prob. there is an infected node at depth  $h$  is constant ( $>0$ )



$x$  ... prob. a node at level  $h-1$  is infected.  
We start at  $x=1$  because  $p_1=1$ .  
 $f(x)$  ... prob. a node at level  $h$  is infected  
 $q$  ... infection prob.  
 $d$  ... degree

**We iterate:**

$$x_1 = f(1)$$

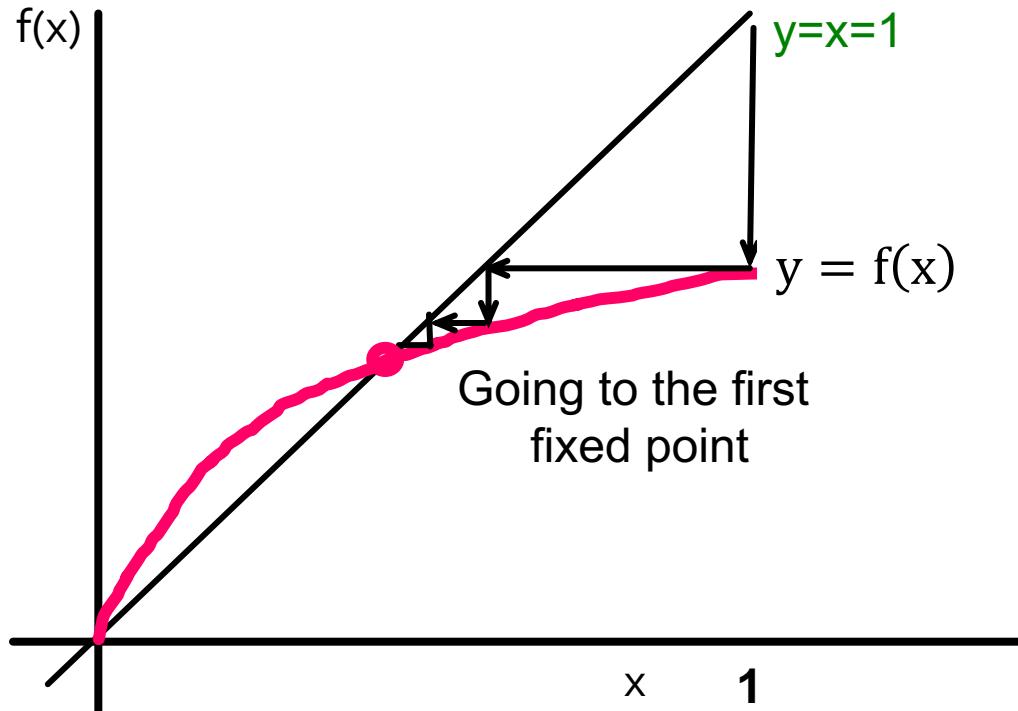
$$x_2 = f(x_1)$$

$$x_3 = f(x_2)$$

If we want to epidemic to die out, then iterating  $f(x)$  must go to zero. So,  $f(x)$  must be **below**  $y = x$ .

- What's the shape of  $f(x)$ ?

# Fixed Point: $f(x) = 1 - (1 - qx)^d$



$x$  ... prob. a node at level  $h-1$  is infected.  
We start at  $x=1$  because  $p_1=1$ .  
 $f(x)$  ... prob. a node at level  $h$  is infected  
 $q$  ... infection prob.  
 $d$  ... degree

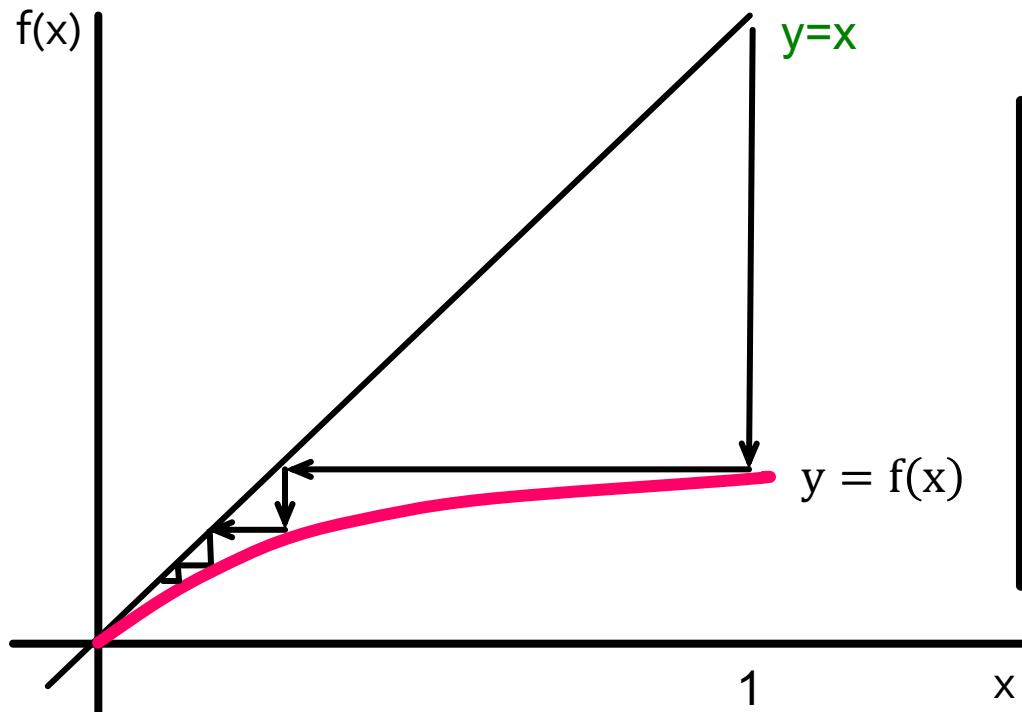
## What do we know about the shape of $f(x)$ ?

- $f(0) = 0$
- $f(1) = 1 - (1 - q)^d < 1$
- $f'(x) = q \cdot d(1 - qx)^{d-1}$
- $f'(0) = q \cdot d$

**$f'(x)$  is monotone:** If  $g'(y)>0$  for all  $y$  then  $g(y)$  is monotone. In our case,  $0 \leq x, q \leq 1$ ,  $d > 1$  so  $f'(x) > 0$ , so  $f(x)$  is monotone.  
 **$f'(x)$  non-increasing:** since term  $(1-qx)^{d-1}$  in  $f'(x)$  is decreasing as  $x$  decreases.

$f'(x)$  is monotone non-increasing on  $[0,1]!$

# Fixed Point: When is this zero?



Reproductive number  $R_0 = q \cdot d$ :  
There is an epidemic if  $R_0 \geq 1$

For the epidemic to die out  
we need  $f(x)$  to be below  $y = x$ !

$$\text{So: } f'(0) = q \cdot d < 1$$

$$\lim_{h \rightarrow \infty} p_h = 0 \text{ when } q \cdot d < 1$$

$q \cdot d$  = expected # of people that get infected

# Important Points

- Reproductive number  $R_0 = q \cdot d$ :
  - It determines if the disease will spread or die out.
- There is an epidemic if  $R_0 \geq 1$
  
- Only  $R_0$  matters:
  - $R_0 \geq 1$ : epidemic never dies and the number of infected people increases exponentially
  - $R_0 < 1$ : Epidemic dies out exponentially quickly

# Measures to Limit the Spreading

- When  $R_0$  is close 1, slightly changing  $q$  or  $d$  can result in epidemics dying out or happening
  - Quarantining people/nodes [reducing  $d$ ]
  - Encouraging better sanitary practices reduces germs spreading [reducing  $q$ ]
  - HIV has an  $R_0$  between 2 and 5
  - Measles has an  $R_0$  between 12 and 18
  - Ebola has an  $R_0$  between 1.5 and 2

# **Application: Social cascades on Flickr and estimating $R_o$ from real data**

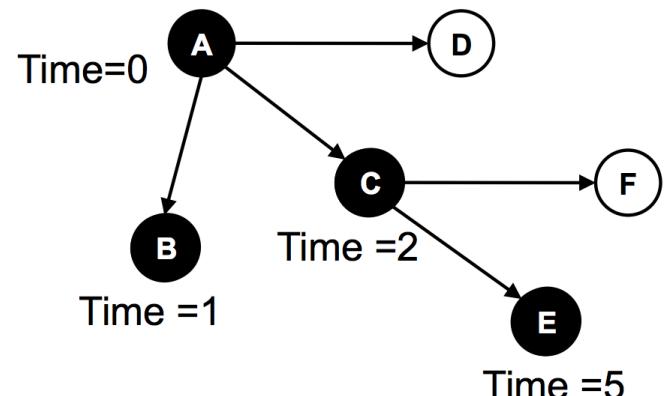
[Characterizing social cascades in Flickr](#). Cha et al. ACM WOSN 2008

# Dataset

- **Flickr social network:**
  - Users are connected to other users via friend links
  - A user can “like/favorite” a photo
- **Data:**
  - **100 days of photo likes**
  - Number of users: 2 million
  - 34,734,221 likes on 11,267,320 photos

# Cascades on Flickr

- Users can be exposed to a photo via social influence (cascade) or external links
- Did a particular like spread through social links?
  - No, if a user likes a photo and if none of his friends have previously liked the photo
  - Yes, if a user likes a photo after at least one of her friends liked the photo → **Social cascade**
- Example social cascade:  
 $A \rightarrow B$  and  $A \rightarrow C \rightarrow E$

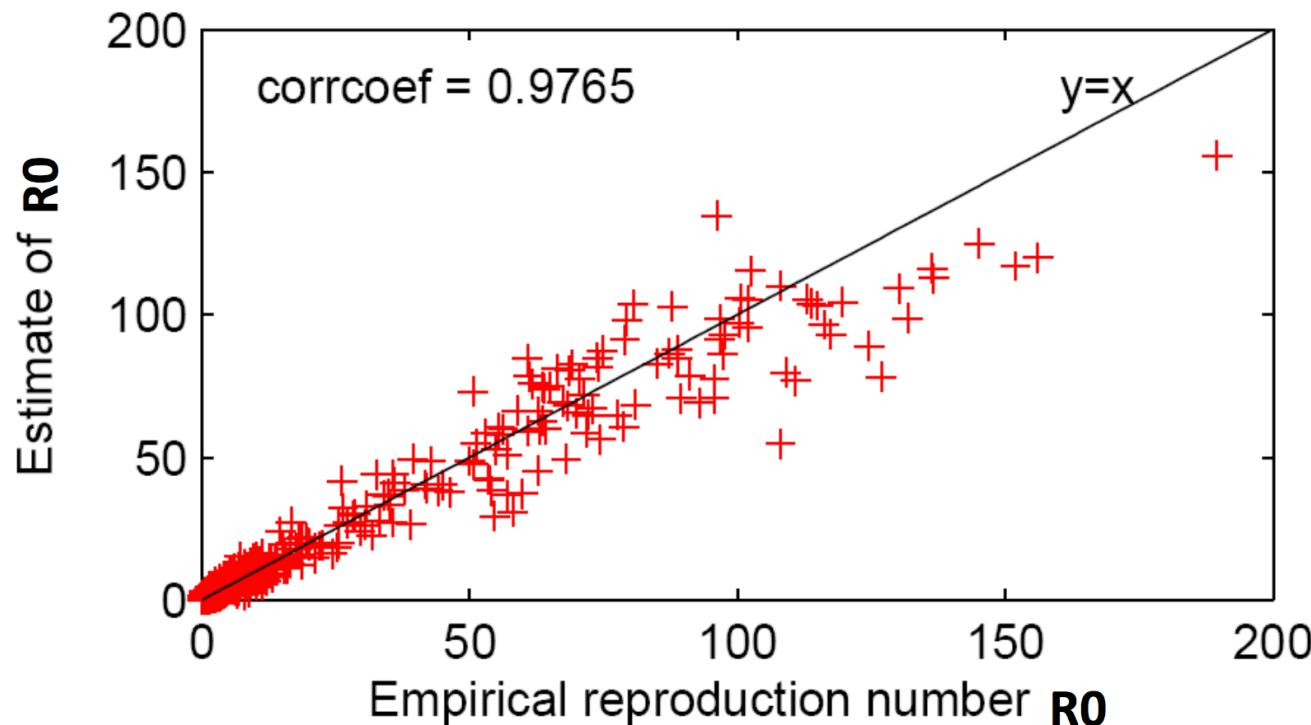


# How to estimate $R_0$ from real data?

- Recall:  $R_0 = q * d$
  - Estimate of  $R_0$ :
    - Estimating  $q$ : Given an infected node count the proportion of its neighbors subsequently infected and average
      - Then:  $R_0 = q * d * \frac{\text{avg}(d_i^2)}{(\text{avg } d_i)^2}$
  - Empirical  $R_0$ :
    - Given start node of a cascade, count the fraction of directly infected nodes and proclaim that to be  $R_0$
- $d$  ... avg degree  
 $d_i$  ... degree of node  $i$
- Correction factor due to skewed degree distribution of the network

# $R_0$ correlation across all photos

- Data from top 1,000 photo cascades
- Each + is one cascade



# Discussion

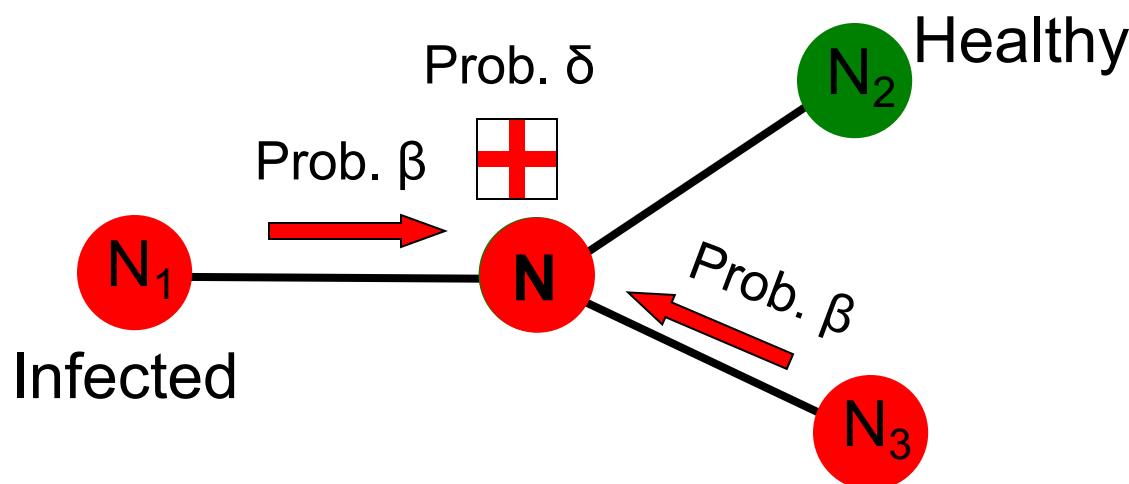
- The basic reproduction number of popular photos is between **1 and 190**
- This is much higher than very infectious diseases like measles, indicating that social networks are efficient transmission media and online content can be very infectious.

# Epidemic models

# Spreading Models of Viruses

## Virus Propagation: 2 Parameters:

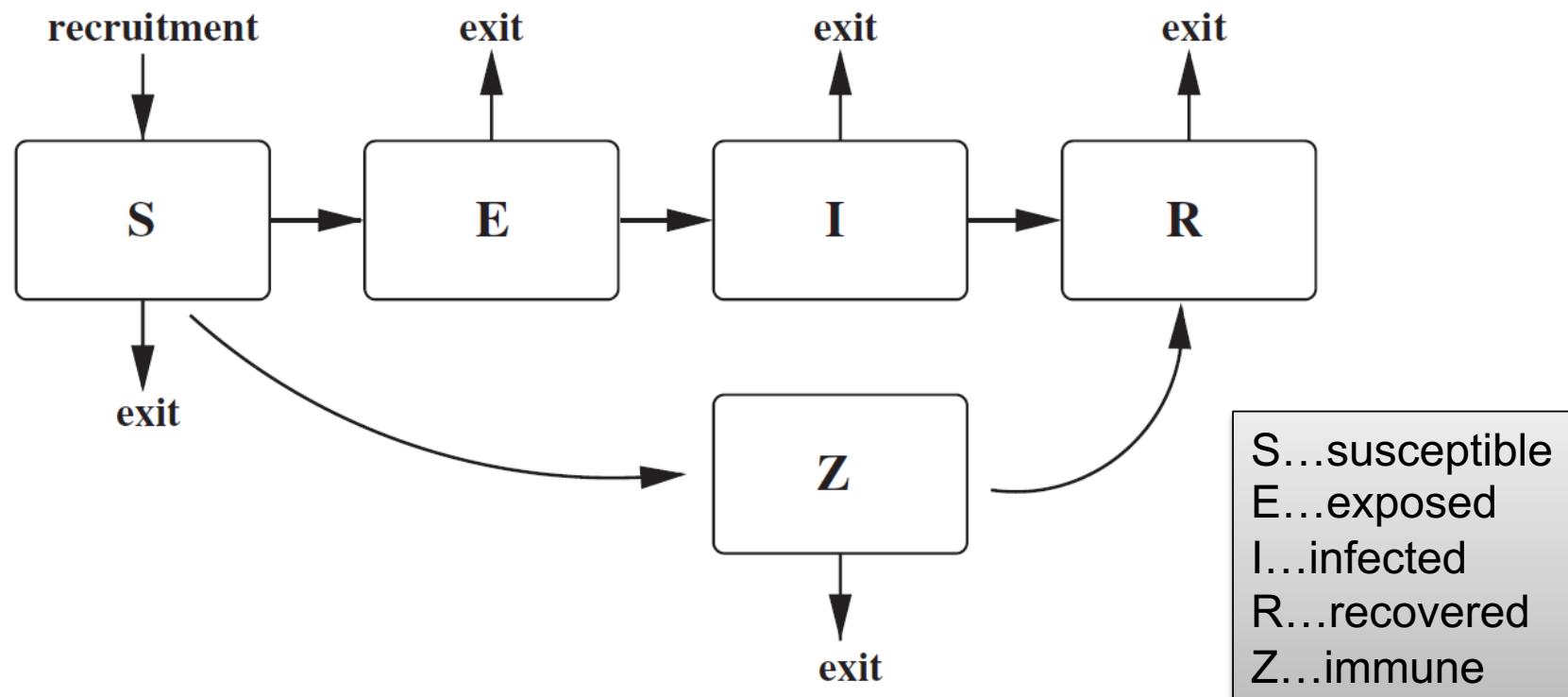
- **(Virus) Birth rate  $\beta$ :**
  - probability that an infected neighbor attacks
- **(Virus) Death rate  $\delta$ :**
  - Probability that an infected node heals



# More Generally: S+E+I+R Models

- General scheme for epidemic models:

- Each node can go through phases:
  - Transition probs. are governed by the model parameters



# SIR Model

- **SIR model:** Node goes through phases

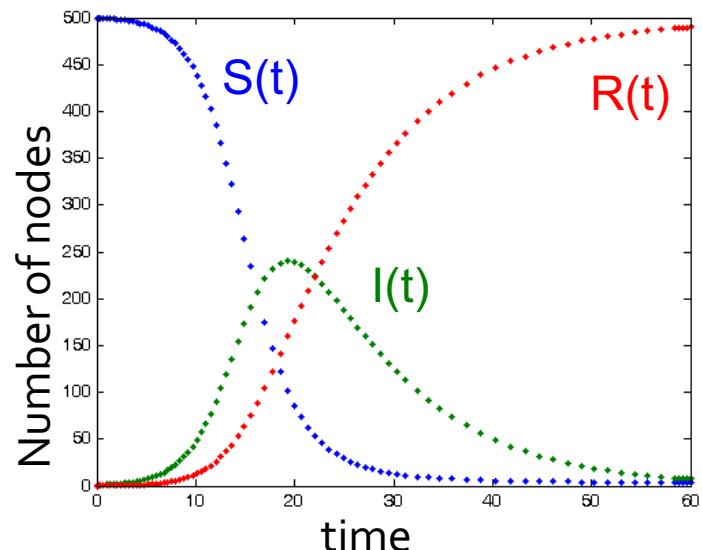


- Models chickenpox or plague:
  - Once you heal, you can never get infected again
- **Assuming perfect mixing (The network is a complete graph) the model dynamics are:**

$$\frac{dS}{dt} = -\beta SI$$

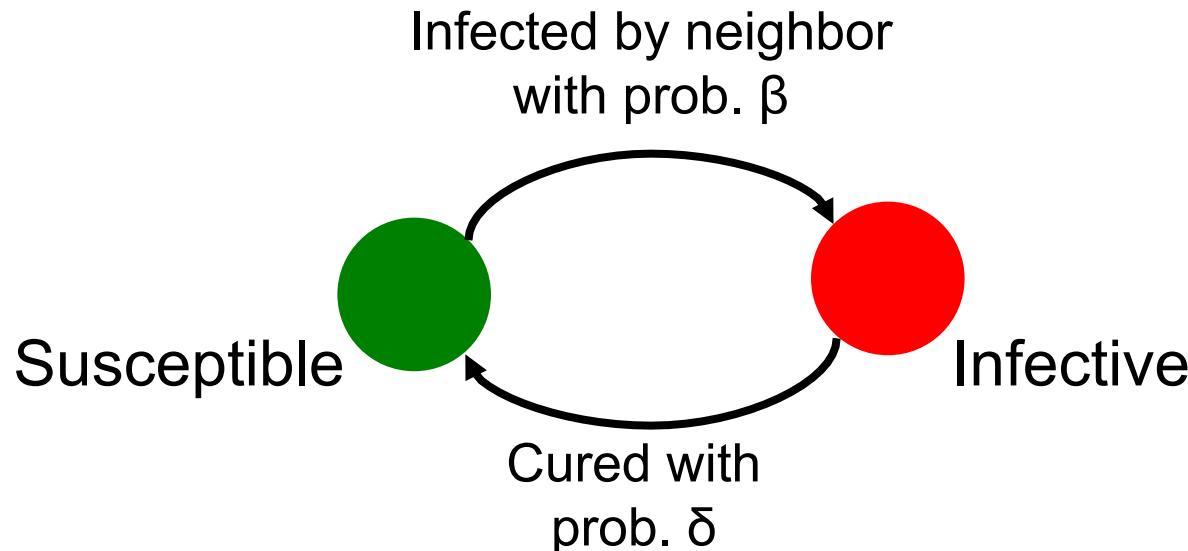
$$\frac{dR}{dt} = \delta I$$

$$\frac{dI}{dt} = \beta SI - \delta I$$

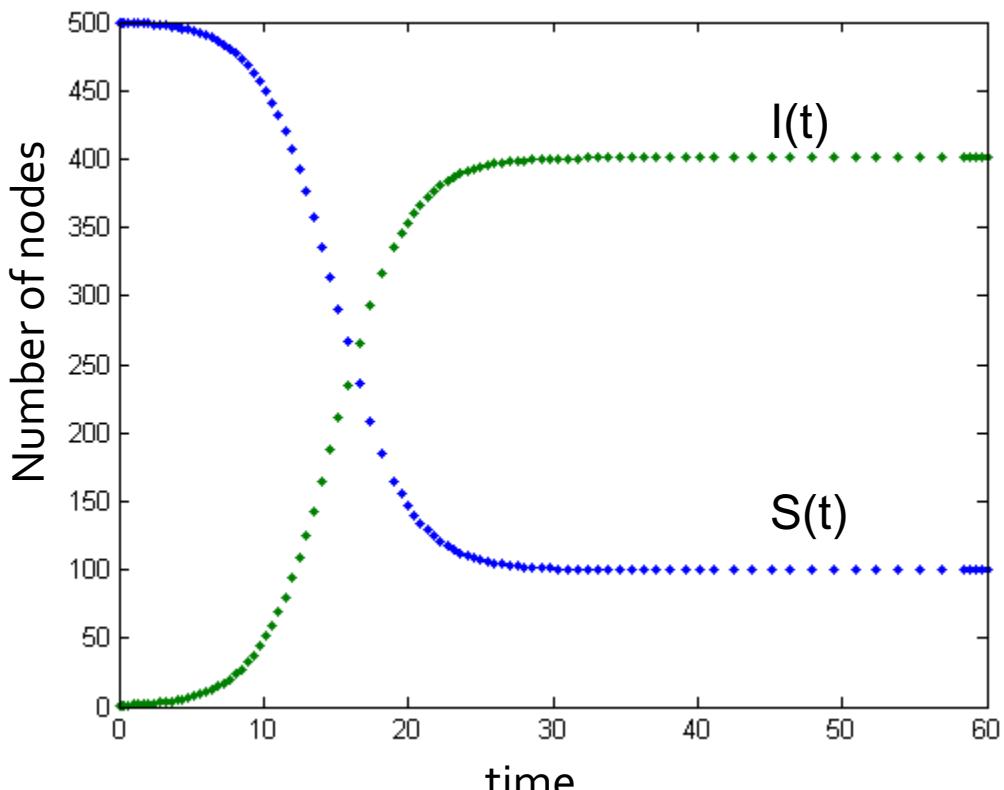


# SIS Model

- **Susceptible-Infective-Susceptible (SIS) model**
- Cured nodes immediately become susceptible
- Virus “strength”:  $s = \beta / \delta$
- **Node state transition diagram:**



# SIS Model



Susceptible



Infected

- **Models flu:**
  - Susceptible node becomes infected
  - The node then heals and become susceptible again
- **Assuming perfect mixing (a complete graph):**

$$\frac{dS}{dt} = -\beta SI + \delta I$$

$$\frac{dI}{dt} = \beta SI - \delta I$$

# Question: Epidemic threshold $\tau$

- **SIS Model:**  
Epidemic threshold of an arbitrary graph  $G$  is  $\tau$ , such that:
  - If virus “strength”  $s = \beta / \delta < \tau$  the epidemic can not happen (it eventually dies out)
- Given a graph what is its epidemic threshold?

# Epidemic Threshold in SIS Model

- Fact: We have no epidemic if:

$$\frac{\beta}{\delta} < \tau = \frac{1}{\lambda_{1,A}}$$

(Virus) Death rate

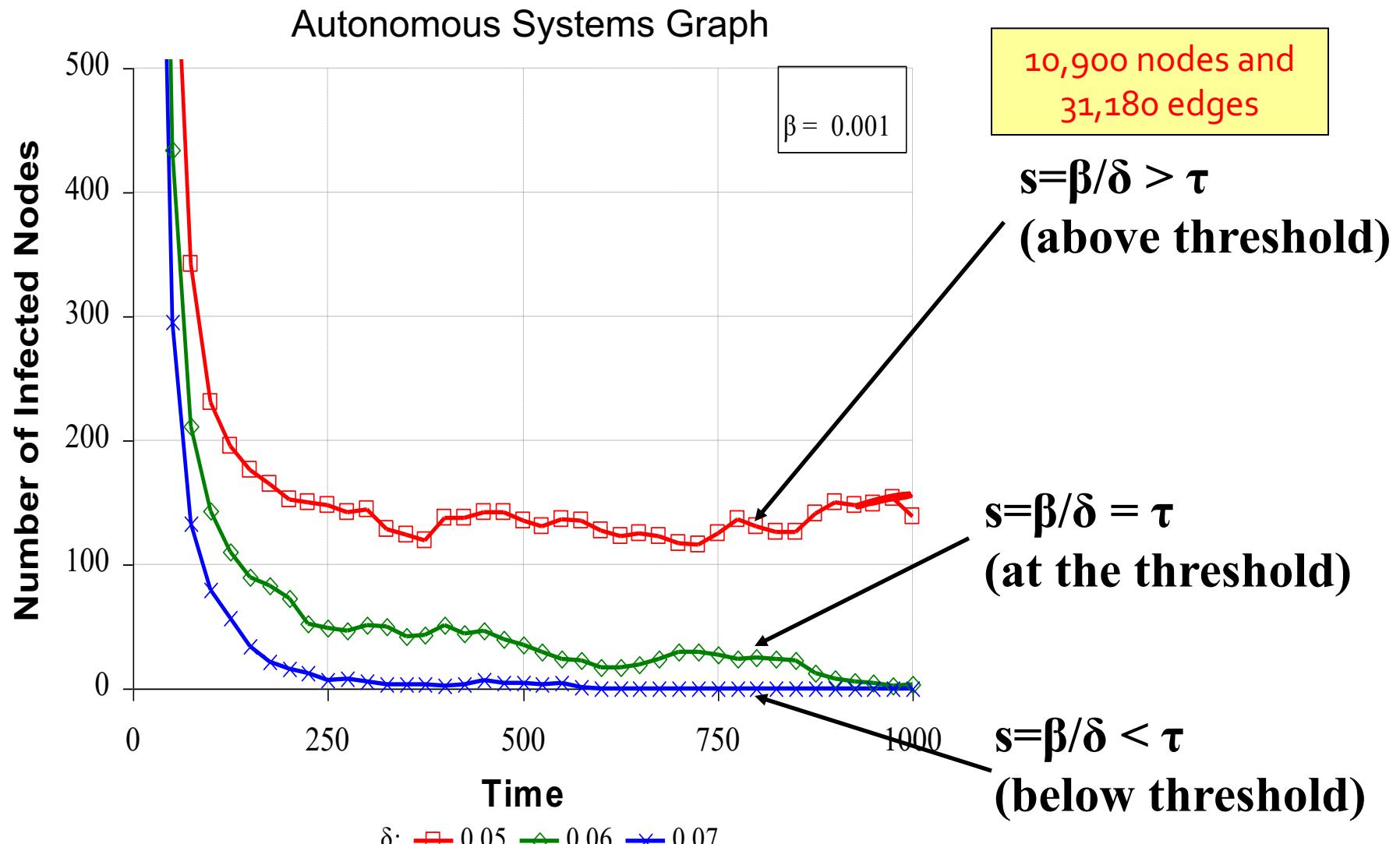
(Virus) Birth rate

Epidemic threshold

largest eigenvalue of adj. matrix **A** of **G**

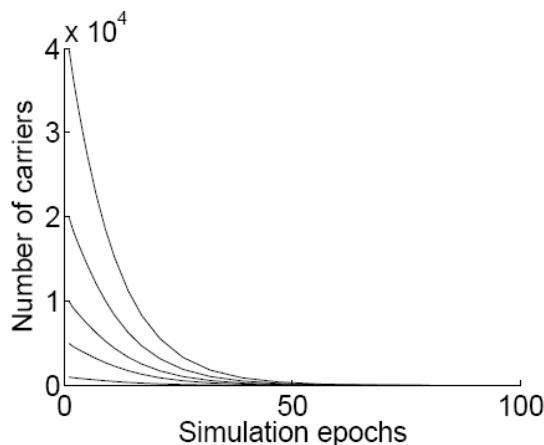
►  $\lambda_{1,A}$  alone captures the property of the graph!

# Experiments (AS graph)

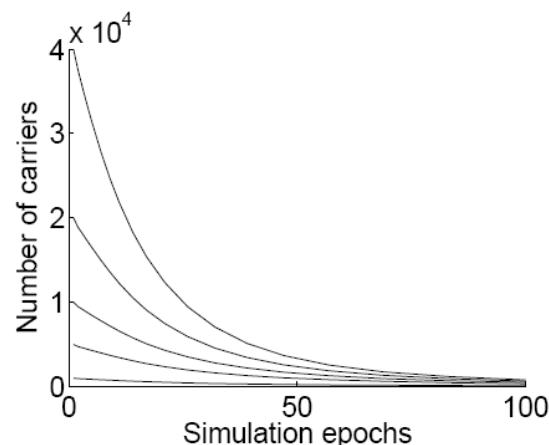


# Experiments

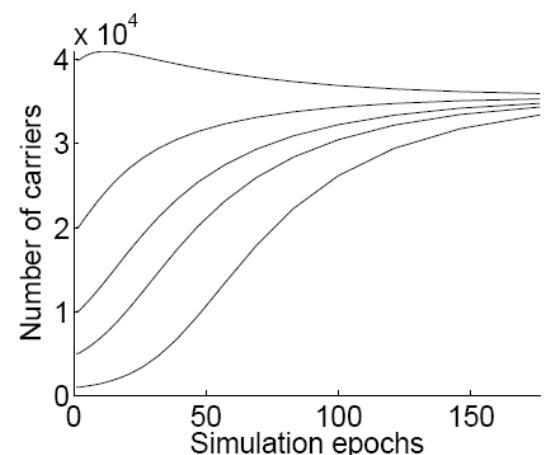
- Does it matter how many people are initially infected?



(a) Below the threshold,  
 $s=0.912$

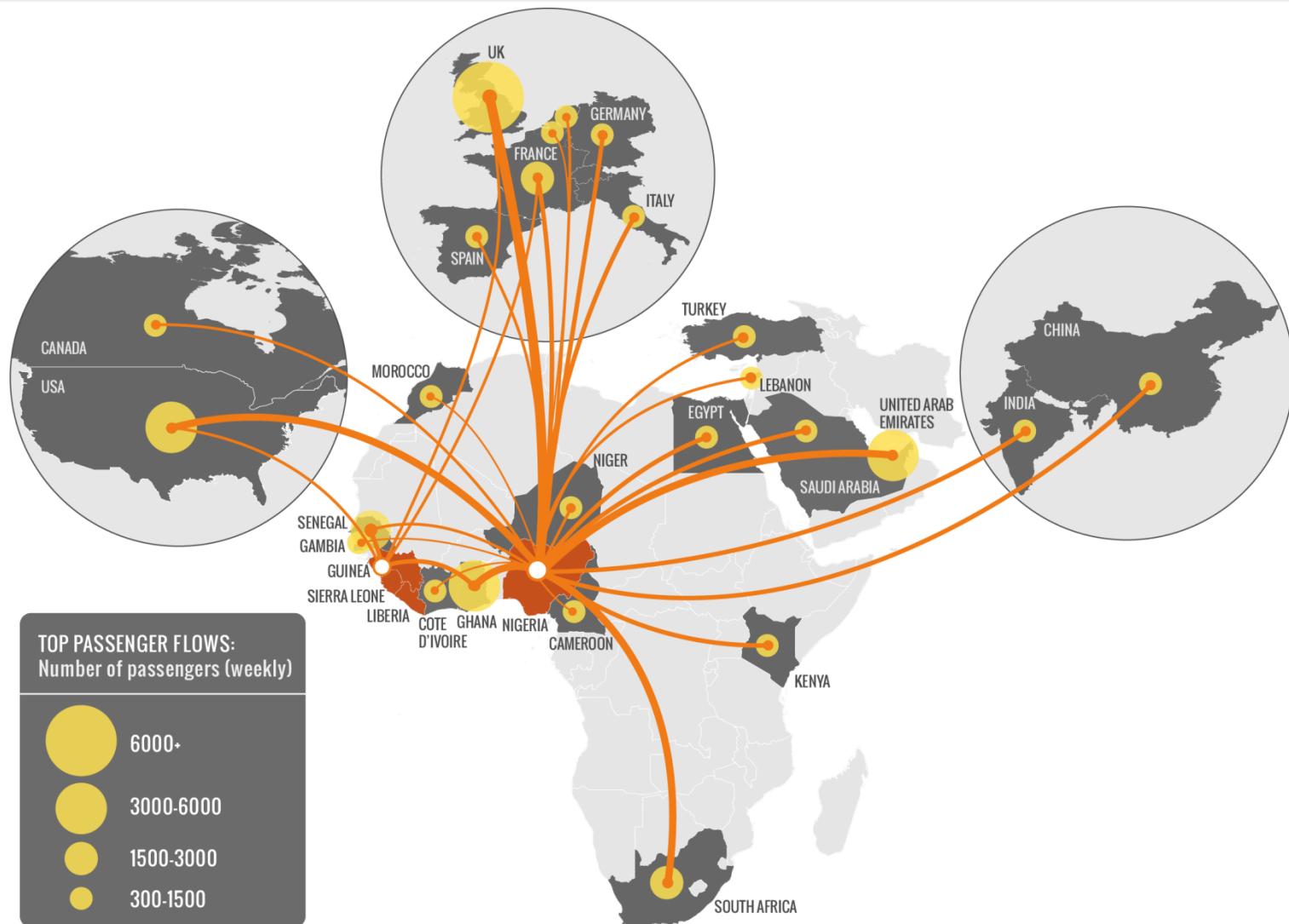


(b) At the threshold,  
 $s=1.003$

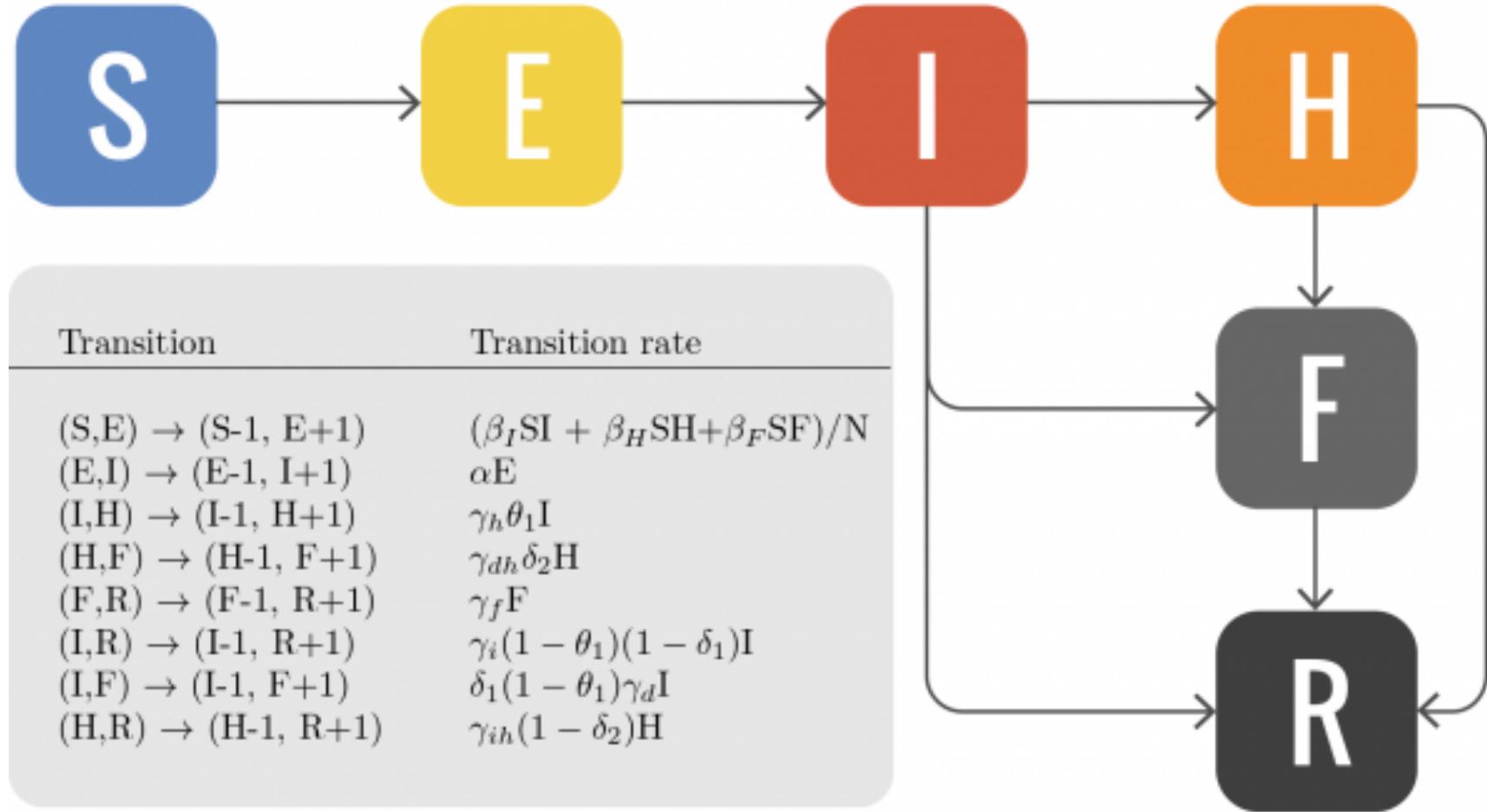


(c) Above the threshold,  
 $s=1.1$

# Modeling Ebola with SEIR

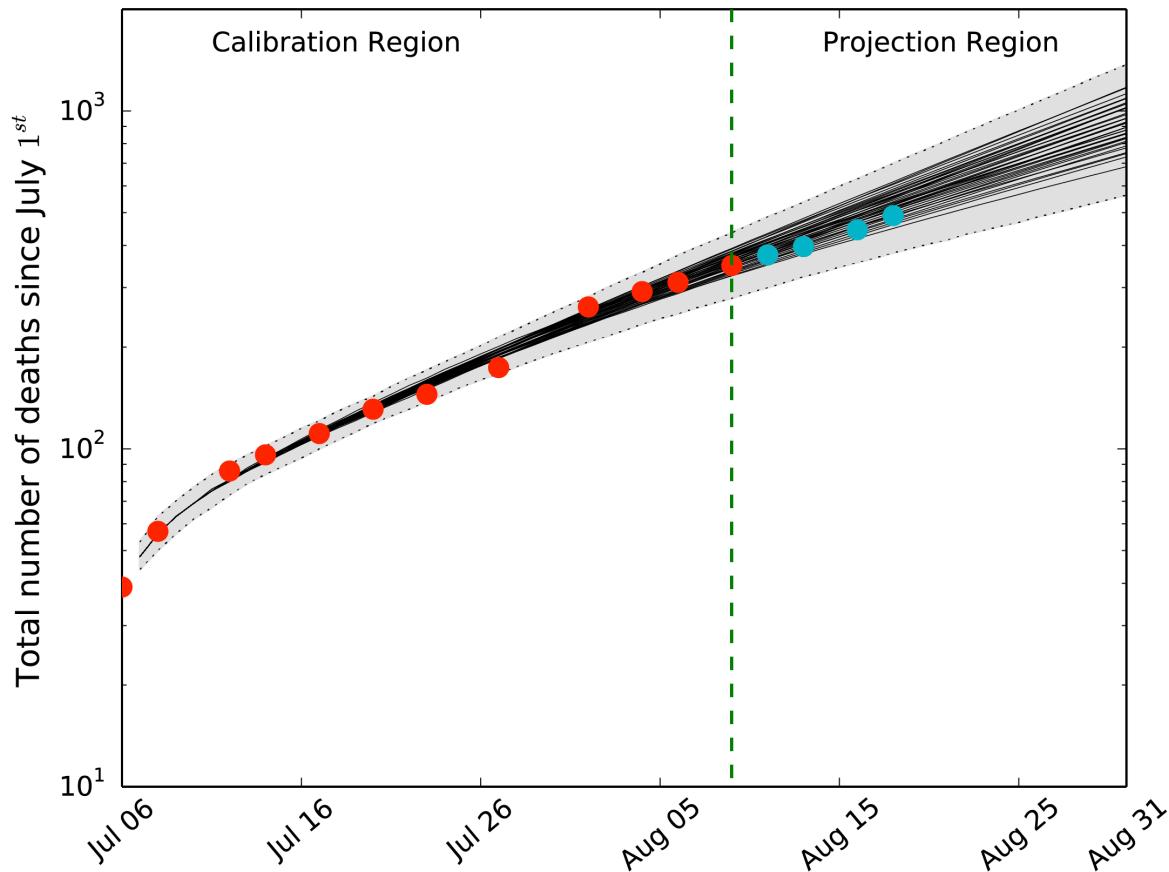


# Example: Ebola



**S:** susceptible individuals, **E:** exposed individuals, **I:** infectious cases in the community,  
**H:** hospitalized cases, **F:** dead but not yet buried, **R:** individuals no longer transmitting the disease

# Example: Ebola, $R_0=1.5-2.0$



Read an article about [how to estimate  \$R\_0\$  of ebola](#).

# **Application: Rumor spread modeling using SEIZ model**

## References:

1. Epidemiological Modeling of News and Rumors on Twitter. Jin et al. SNAKDD 2013
2. False Information on Web and Social Media: A survey. Kumar et al., arXiv :1804.08559

# SEIZ model: Extension of SIS model



Susceptible Twitter accounts

Infected Believe news / rumor, (I) post a tweet

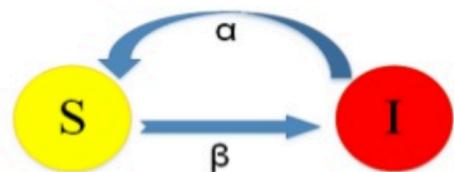
Exposed Be exposed but not yet believe

Skeptics Skeptics, do not tweet

**Disease**

**Twitter**

# Recap: SIS model



$$S = S(t), I = I(t)$$

$\beta$  = rate of contact between 2 individuals

$\alpha$  = rate of recovery

$$\frac{d[S]}{dt} = \dot{S} = -\beta SI + \alpha I$$

$$\frac{d[I]}{dt} = \dot{I} = \beta SI - \alpha I$$

## Disease Applications:

- Influenza
- Common Cold

## Twitter Application Reasoning:

- An individual either believes a rumor (I),
- or is susceptible to believing the rumor (S)

# Details of the SEIZ model

## Notation:

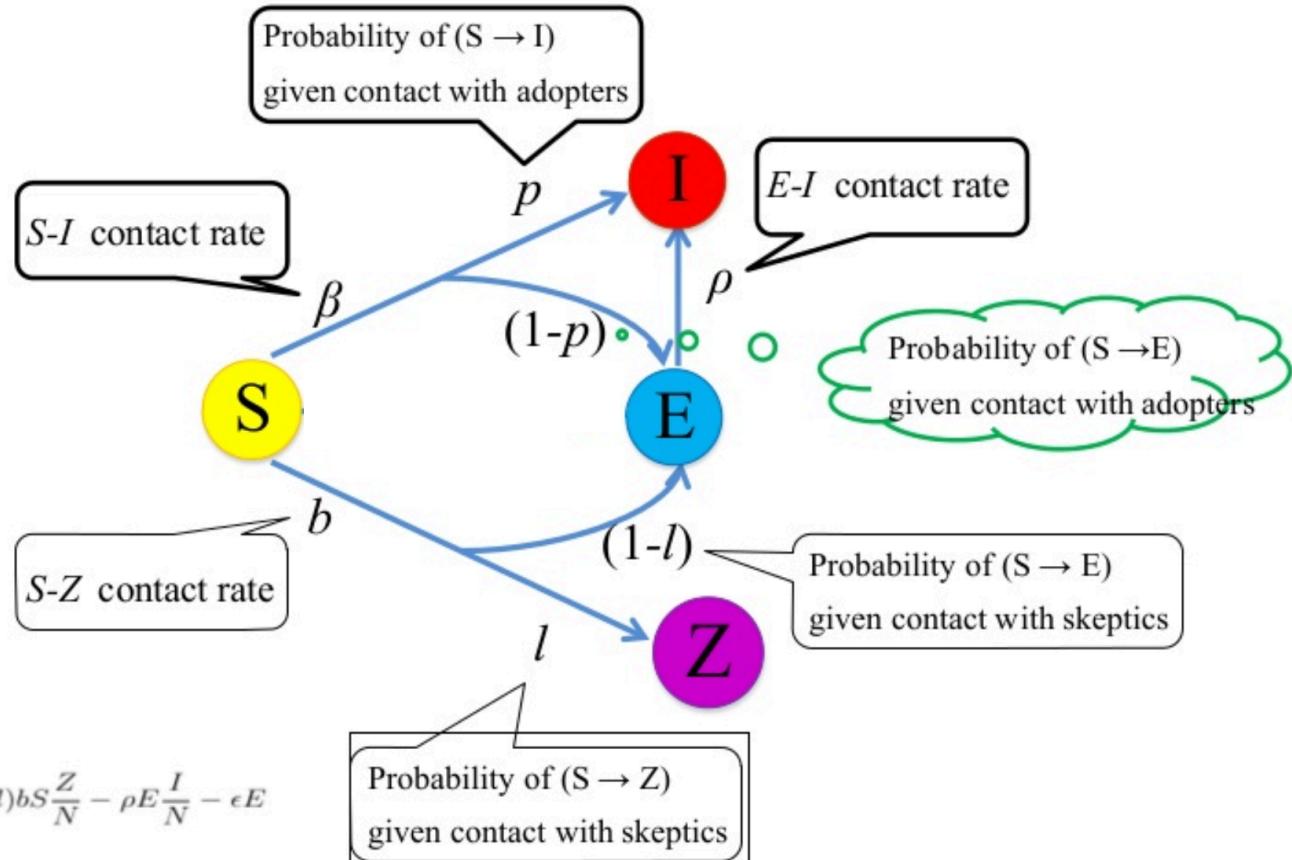
- S = Susceptible
- I = Infected
- E = Exposed
- Z = Skeptics

$$\frac{d[S]}{dt} = -\beta S \frac{I}{N} - bS \frac{Z}{N}$$

$$\frac{d[E]}{dt} = (1-p)\beta S \frac{I}{N} + (1-l)bS \frac{Z}{N} - \rho E \frac{I}{N} - eE$$

$$\frac{d[I]}{dt} = p\beta S \frac{I}{N} + \rho E \frac{I}{N} + eE$$

$$\frac{d[Z]}{dt} = lbS \frac{Z}{N}$$



# Dataset

# Tweets collected from eight stories: Four rumors and four real

---

- Boston Marathon Explosion. 04-15-2013
- Pope Resignation. 02-11-2013
- Venezuela's refinery explosion. 08-25-2012
- Michelle Obama at the 2013 Oscars. 02-24-2013

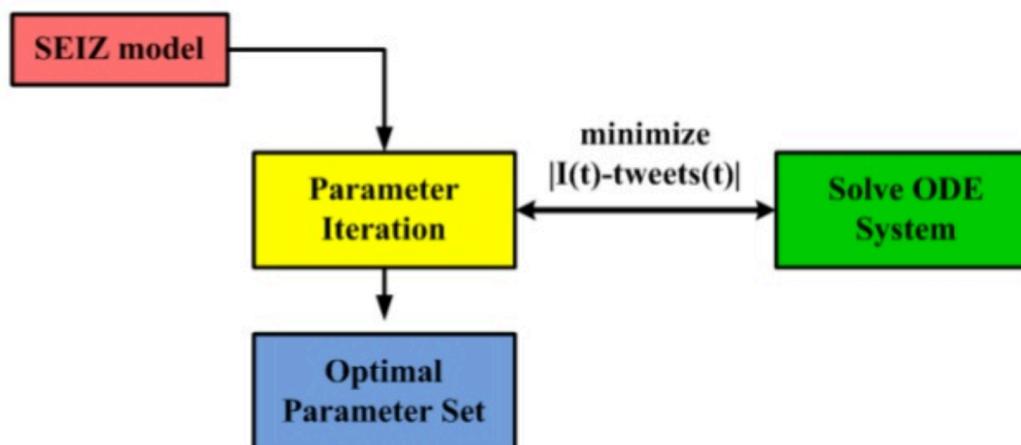
- Boston Marathon Explosion. 04-15-2013
  - Pope Resignation. 02-11-2013
  - Venezuela's refinery explosion. 08-25-2012
  - Michelle Obama at the 2013 Oscars. 02-24-2013

RUMORS  
na injured. 04-23-2013  
nsday rumor. 12-21-2012  
Castro's coming death. 10-15-2012  
and shooting in Mexico. 09-05-2012

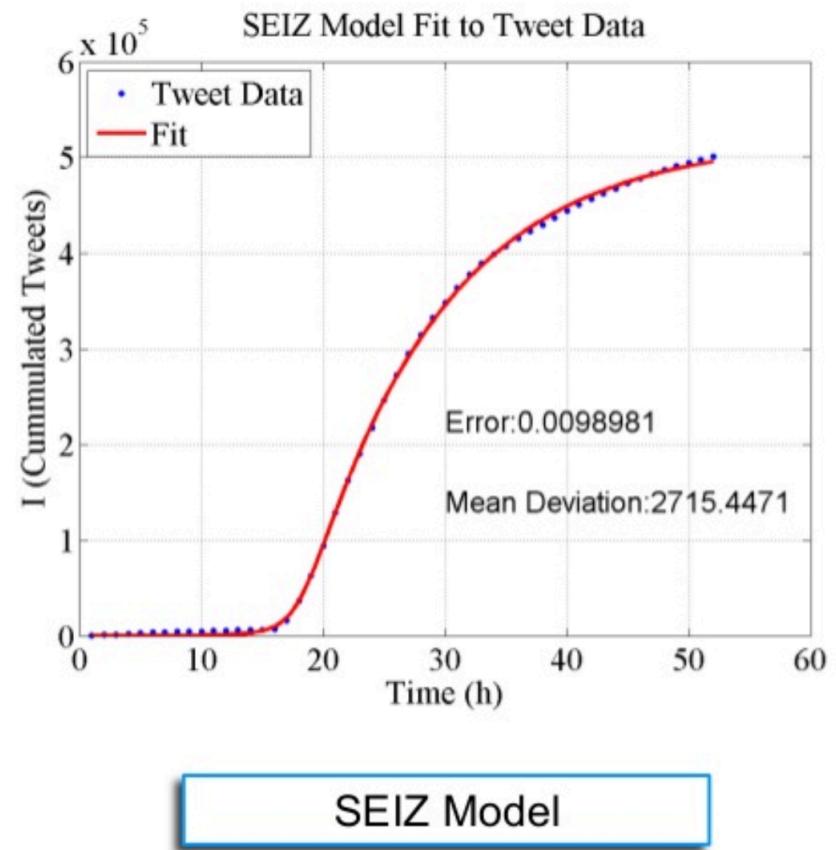
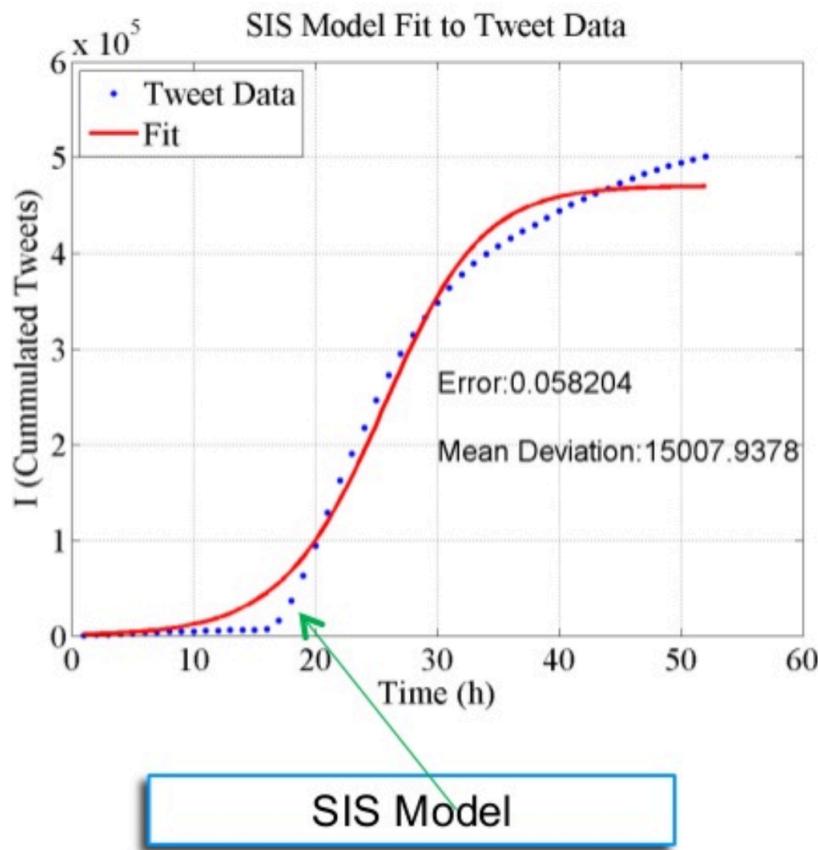


# Method: Fitting SEIZ model to data

- SEIZ model is fit to each cascade to minimize the difference  $|I(t) - \text{tweets}(t)|$ :
  - $\text{tweets}(t)$  = number of rumor tweets
  - $I(t)$  = the estimated number of rumor tweets by the model
- Use grid-search and find the parameters with minimum error



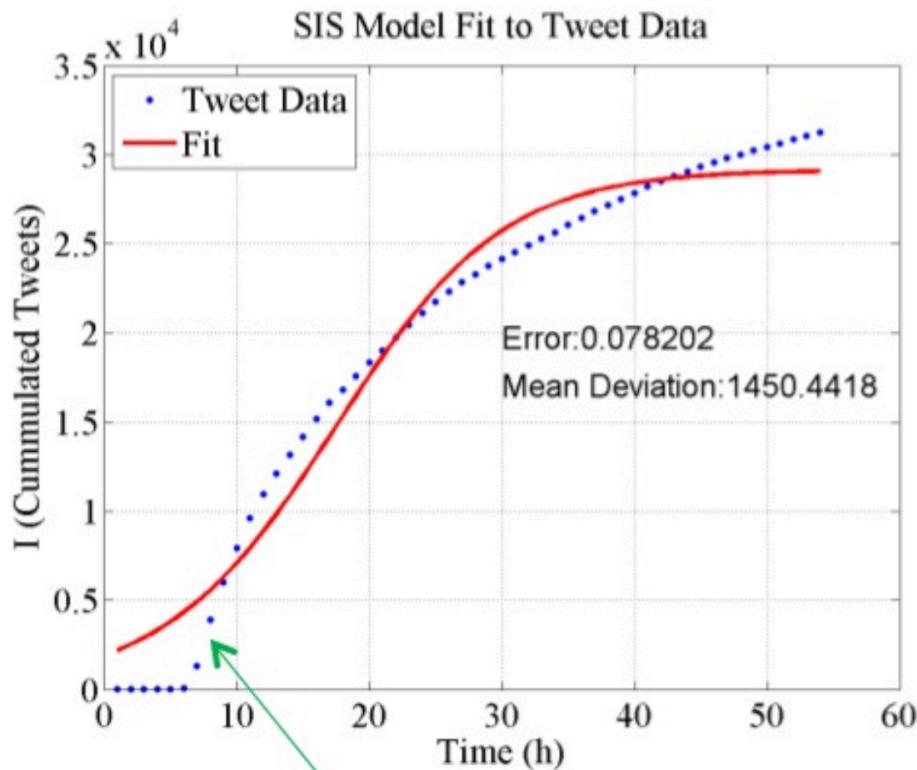
# Fitting to “Boston Marathon Bombing”



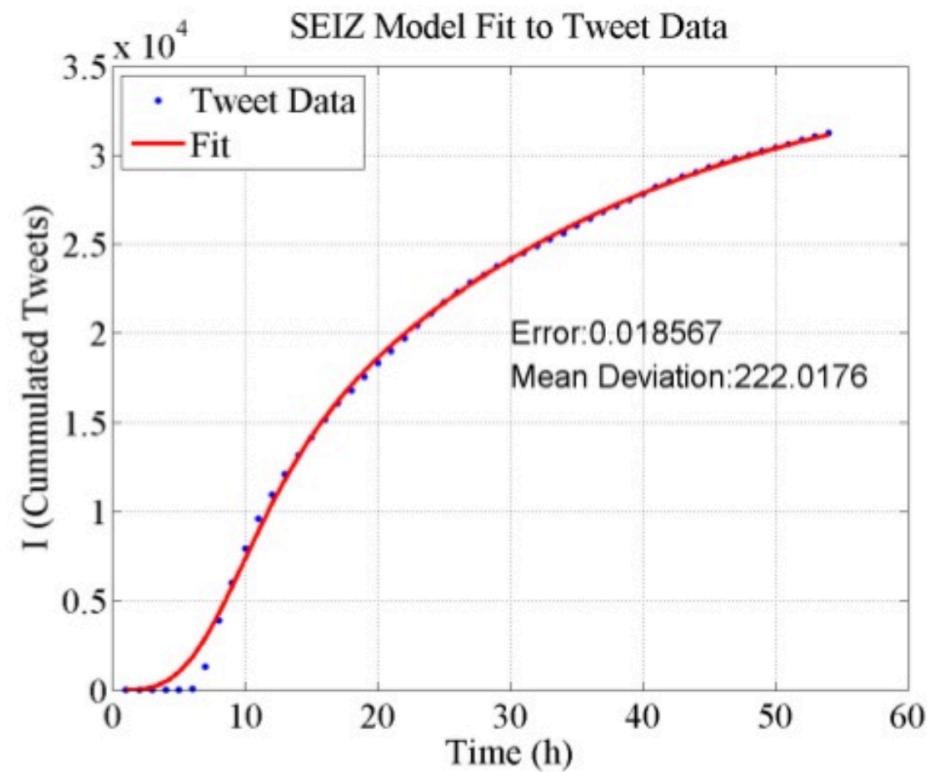
$$\text{Error} = \text{norm}(I - \text{tweets}) / \text{norm}(\text{tweets})$$

SEIZ model better models the real data, especially at initial points

# Fitting to "Pope resignation" data



SIS Model

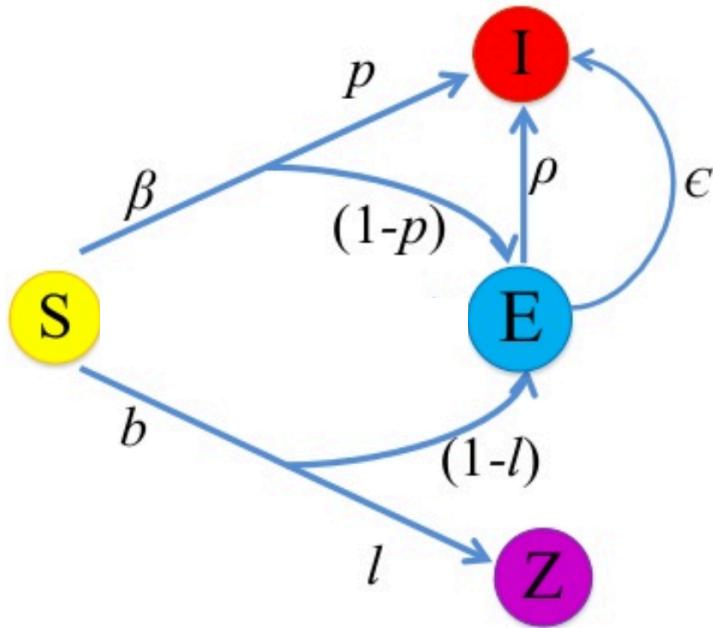


SEIZ Model

SEIZ model better models the real data, especially at initial points

# Rumor detection with SEIZ model

By SEIZ model parameters



Notation:  
S = Susceptible  
I = Infected  
E = Exposed  
Z = Skeptics

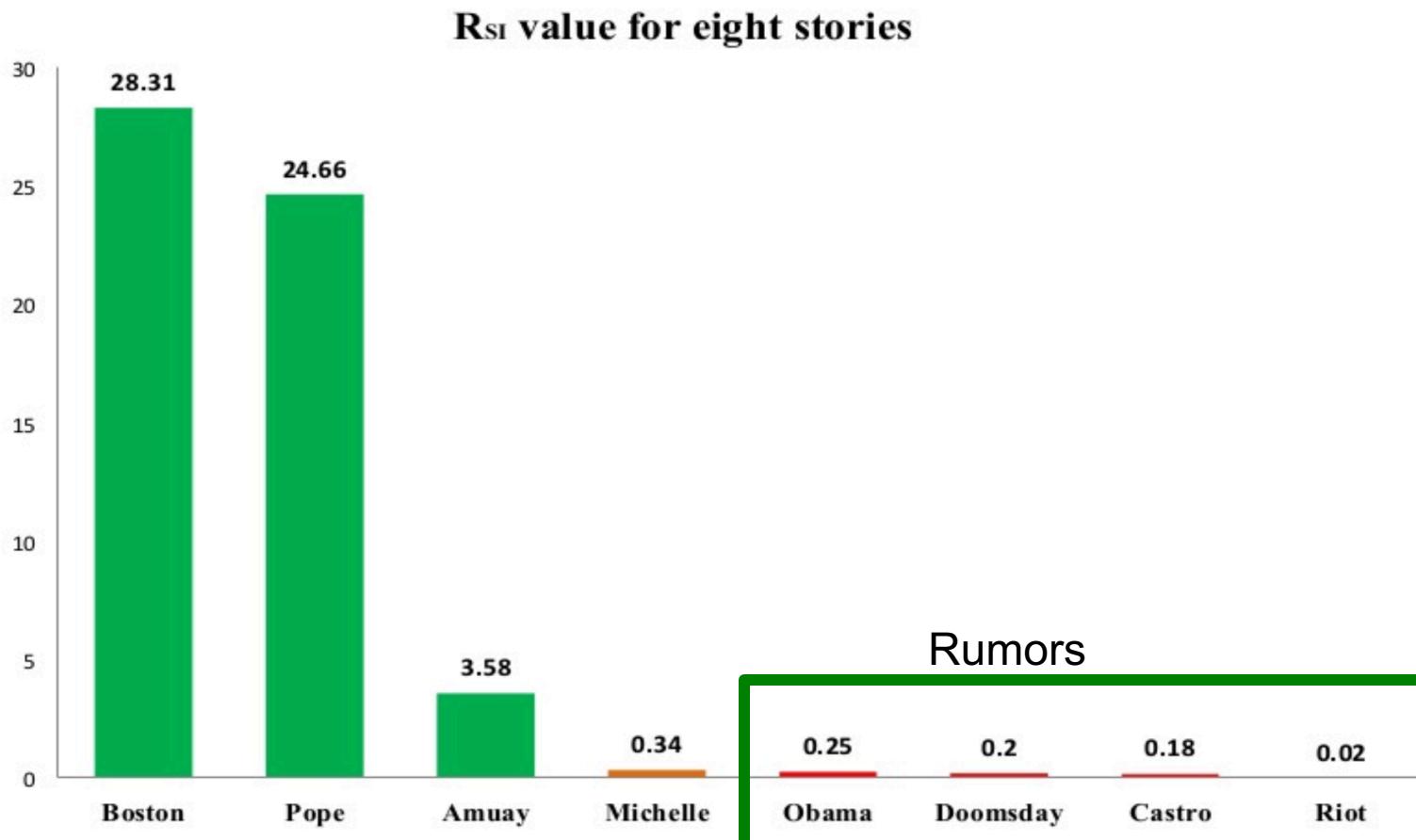
New metric:

$$R_{SI} = \frac{(1-p)\beta + (1-l)b}{\rho + \epsilon}$$

All parameters learned by model fitting to real data (from previous slides)

$R_{SI}$ , a kind of flux ratio, the ratio of effects entering E to those leaving E.

# Rumor detection by $R_{SI}$

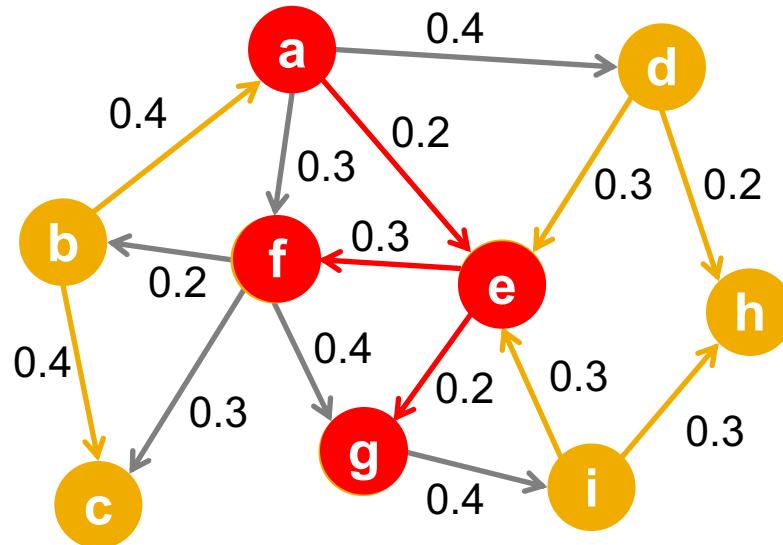


Parameters obtained by fitting SEIZ model  
efficiently identifies rumors vs. news

# Independent Cascade Model

# Independent Cascade Model

- Initially some nodes  $S$  are active
- Each edge  $(u,v)$  has probability (weight)  $p_{uv}$



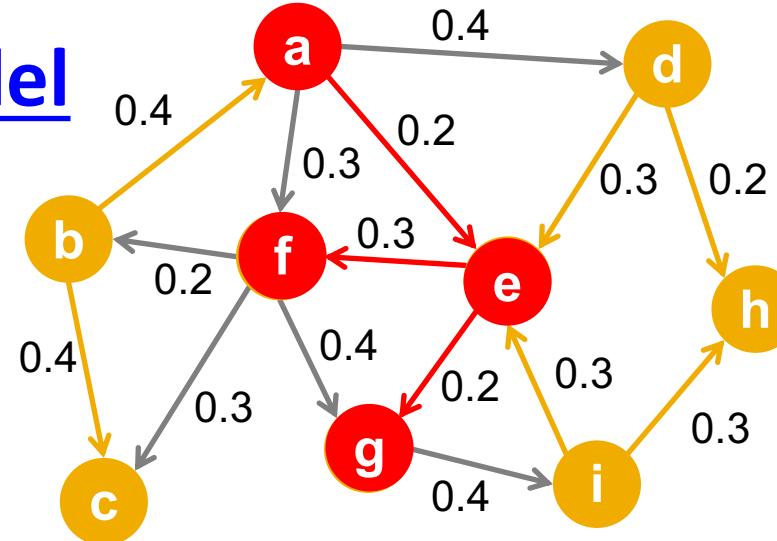
- When node  $u$  becomes active/infected:
  - It activates each out-neighbor  $v$  with prob.  $p_{uv}$
- Activations spread through the network!

# Independent Cascade Model

- Independent cascade model is simple but requires many parameters!

- Estimating them from data is very hard [Goyal et al. 2010]

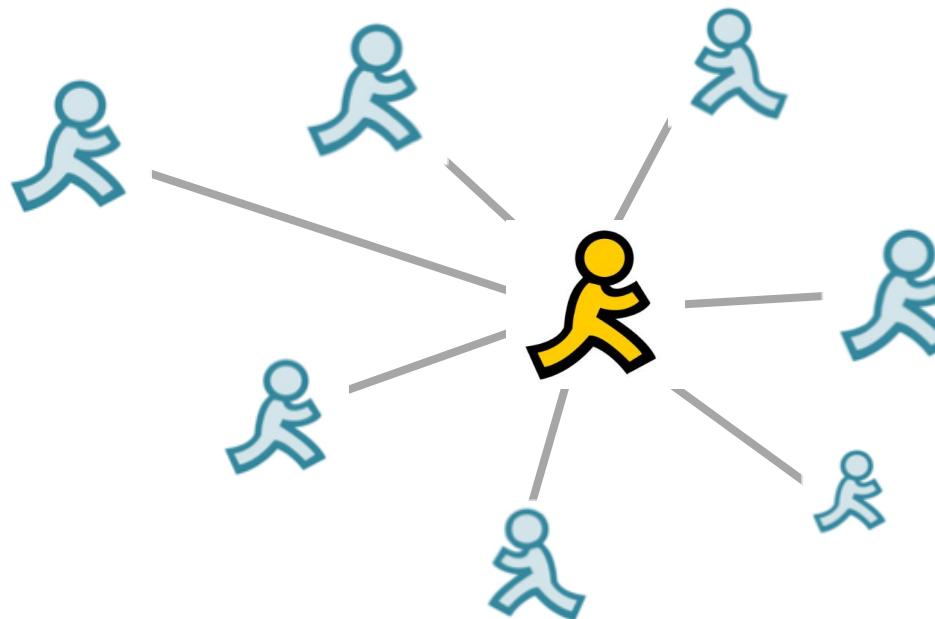
- **Solution:** Make all edges have the same weight (which brings us back to the SIR model)
  - Simple, but too simple
- **Can we do something better?**



# Exposures and Adoptions

## ■ From exposures to adoptions

- **Exposure:** Node's neighbor exposes the node to the contagion
- **Adoption:** The node acts on the contagion

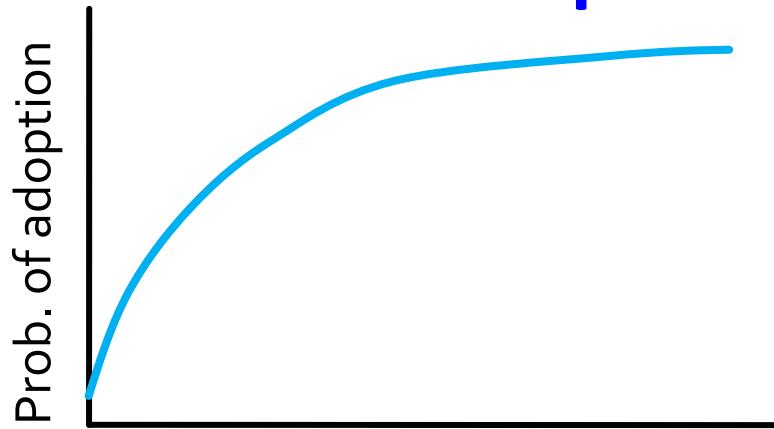


# Exposure Curves

- Exposure curve:

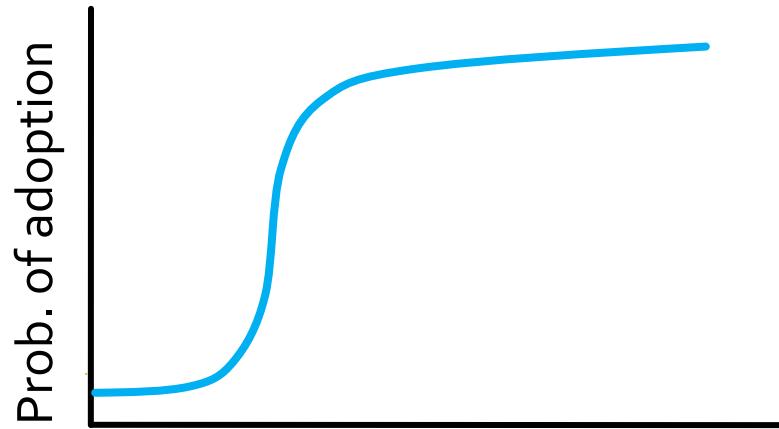
- Probability of adopting new behavior depends on the total number of friends who have already adopted

- What's the dependence?



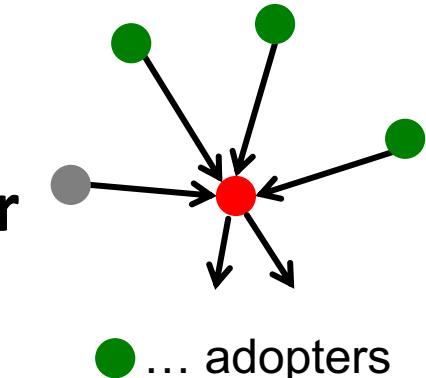
$k$  = number of friends adopting

“Probabilistic” spreading:  
Viruses, Information



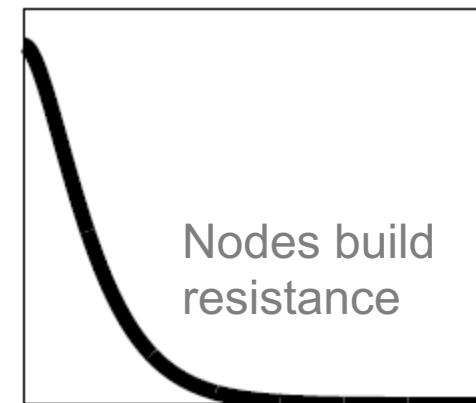
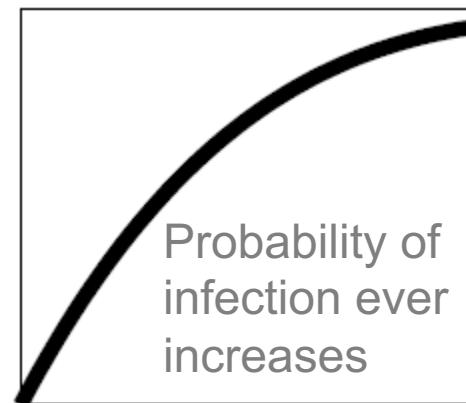
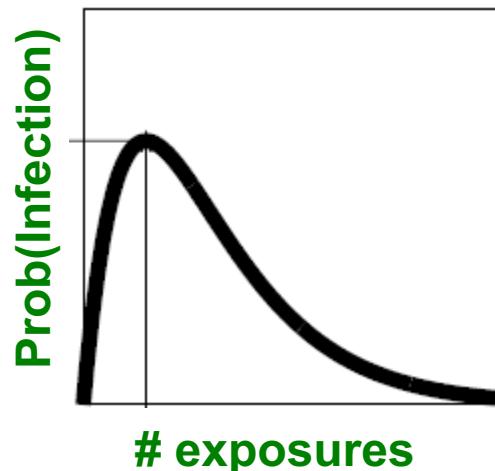
$k$  = number of friends adopting

Critical mass:  
Decision making



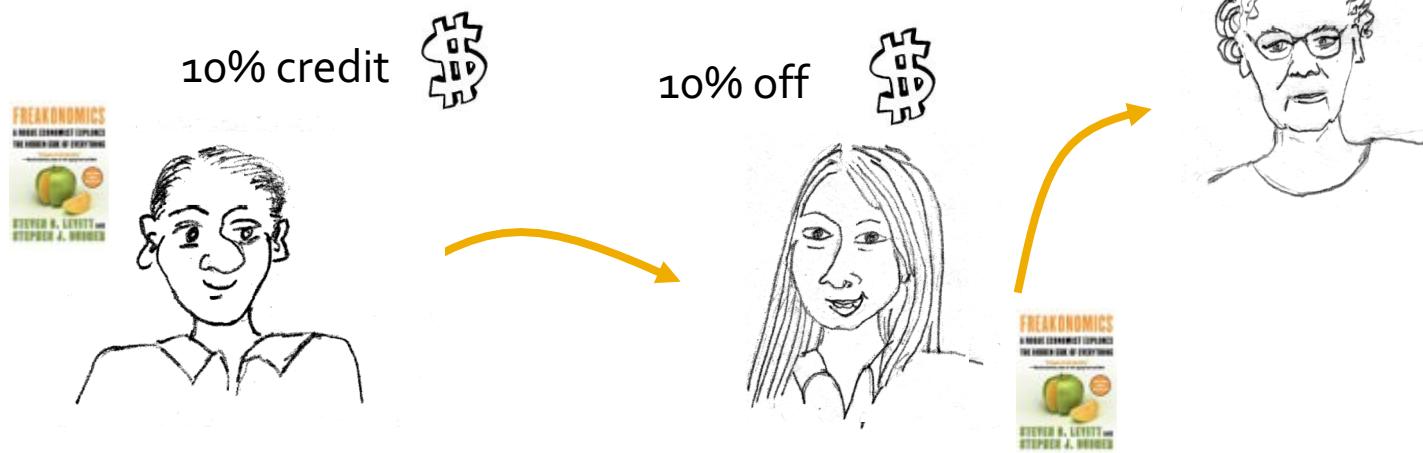
# Exposure Curves

- **From exposures to adoptions**
  - **Exposure:** Node's neighbor exposes the node to information
  - **Adoption:** The node acts on the information
- **Examples of different adoption curves:**



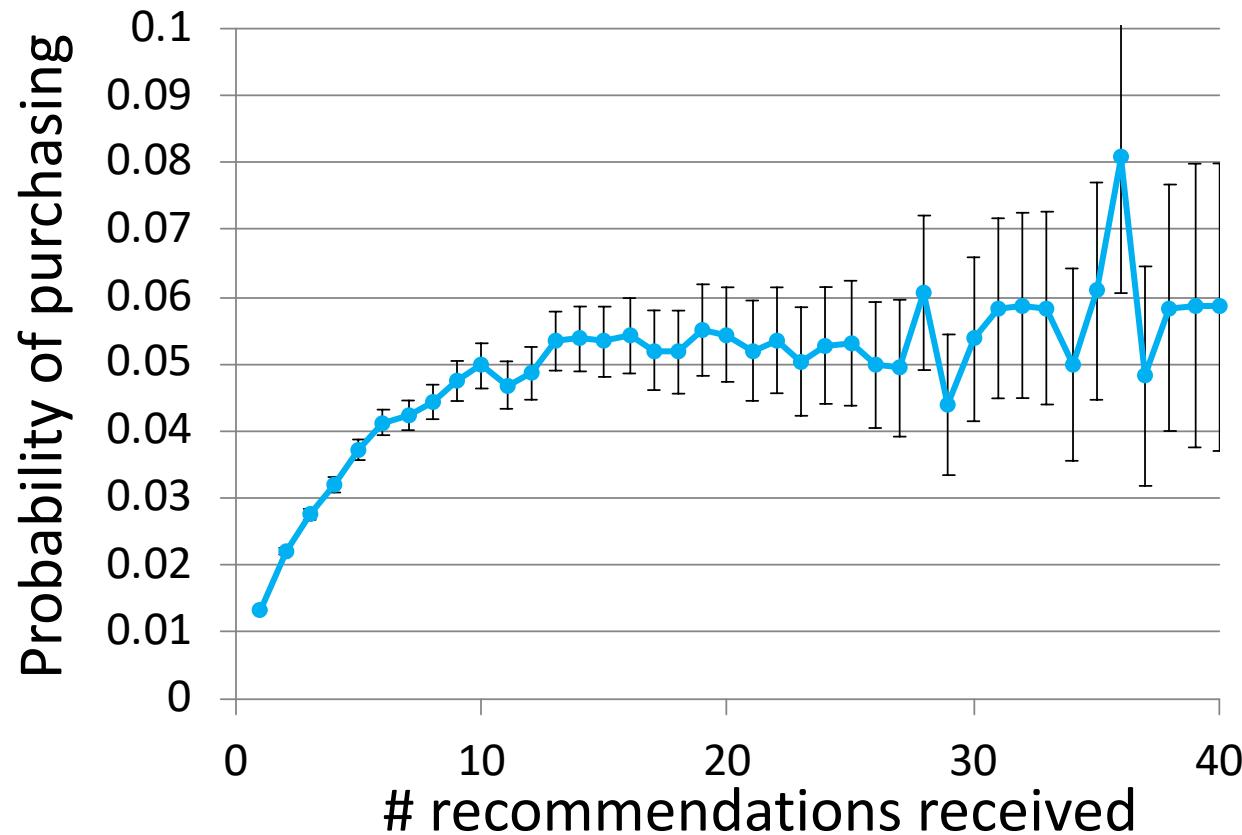
# Diffusion in Viral Marketing

- Senders and followers of recommendations receive discounts on products



- Data: Incentivized Viral Marketing program
  - 16 million recommendations
  - 4 million people, 500k products

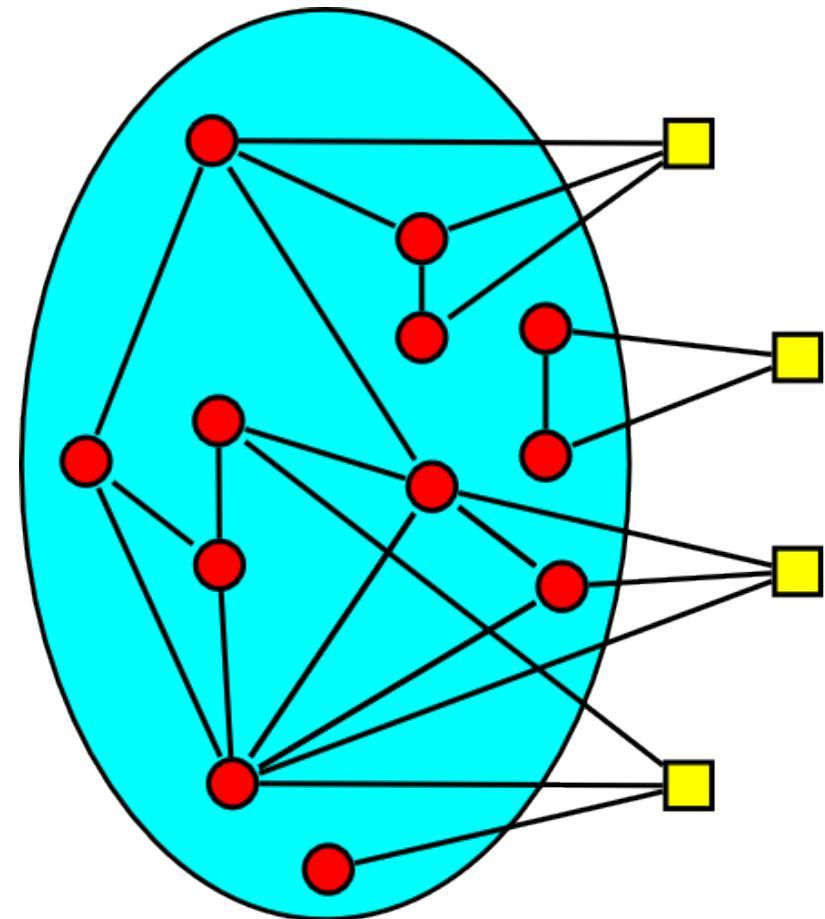
# Exposure Curve: Validation



DVD recommendations  
(8.2 million observations)

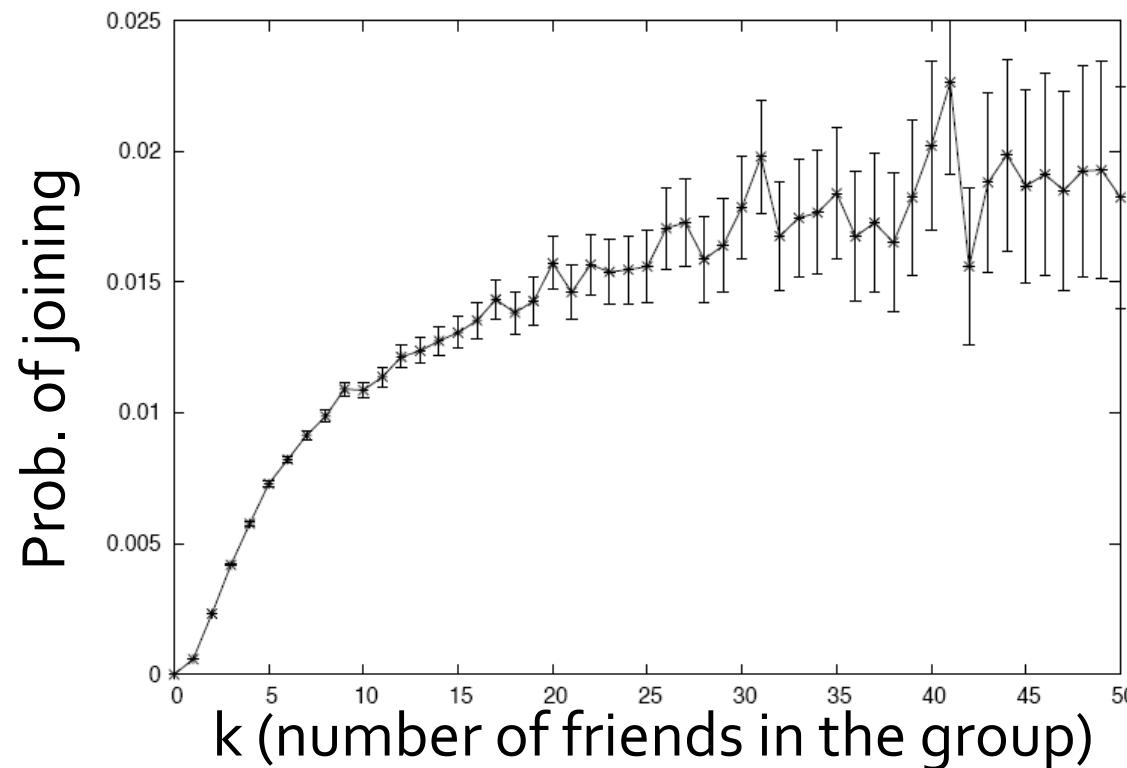
# Exposure Curve: LiveJournal

- Group memberships spread over the network:
  - Red circles represent existing group members
  - Yellow squares may join
- Question:
  - How does prob. of joining a group depend on the number of friends already in the group?



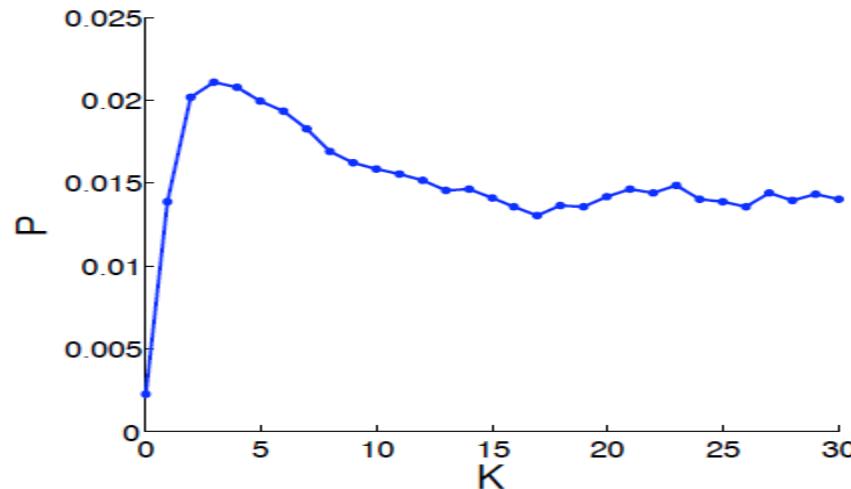
# Exposure Curve: LiveJournal

## ■ LiveJournal group membership



# Exposure Curve: Information

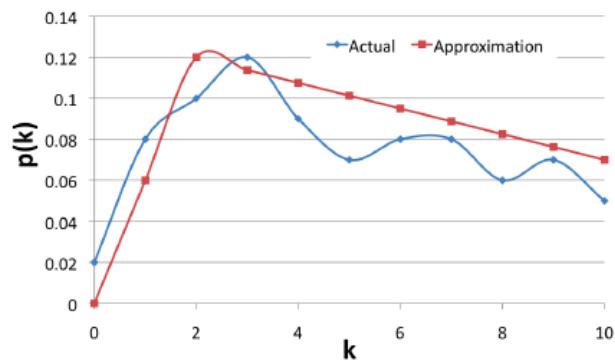
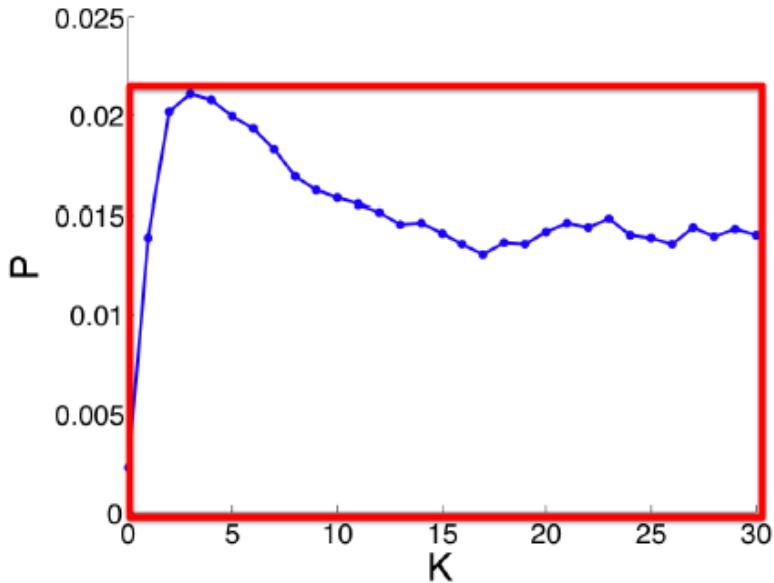
- Twitter [Romero et al. '11]
  - Aug '09 to Jan '10, 3B tweets, 60M users



- Avg. exposure curve for the top 500 hashtags
- What are the most important aspects of the shape of exposure curves?
- Curve reaches peak fast, decreases after!

# Modeling the Shape of the Curve

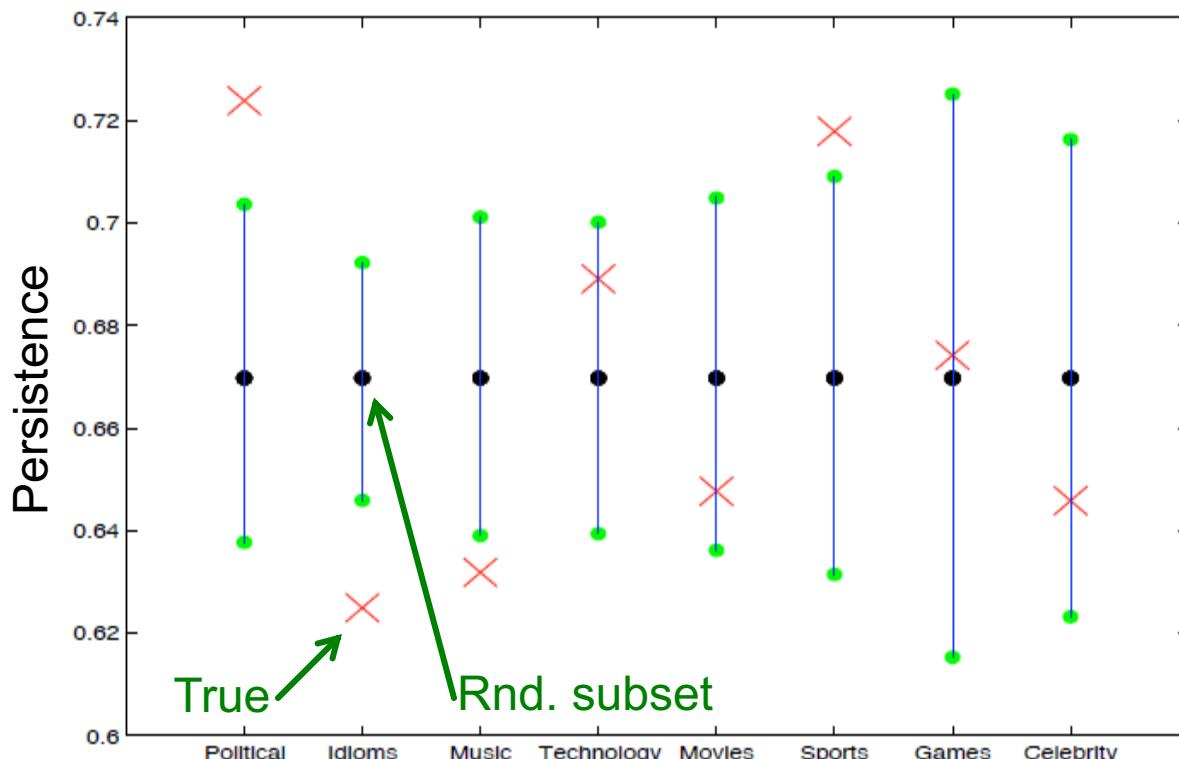
- Persistence of  $P$  is the ratio of the area under the curve  $P$  and the area of the rectangle of height  $\max(P)$ , width  $\max(D(P))$ 
  - $D(P)$  is the domain of  $P$
  - Persistence measures the decay of exposure curves
- Stickiness of  $P$  is  $\max(P)$ 
  - Stickiness is the probability of usage at the most effective exposure



# Exposure Curve: Persistence

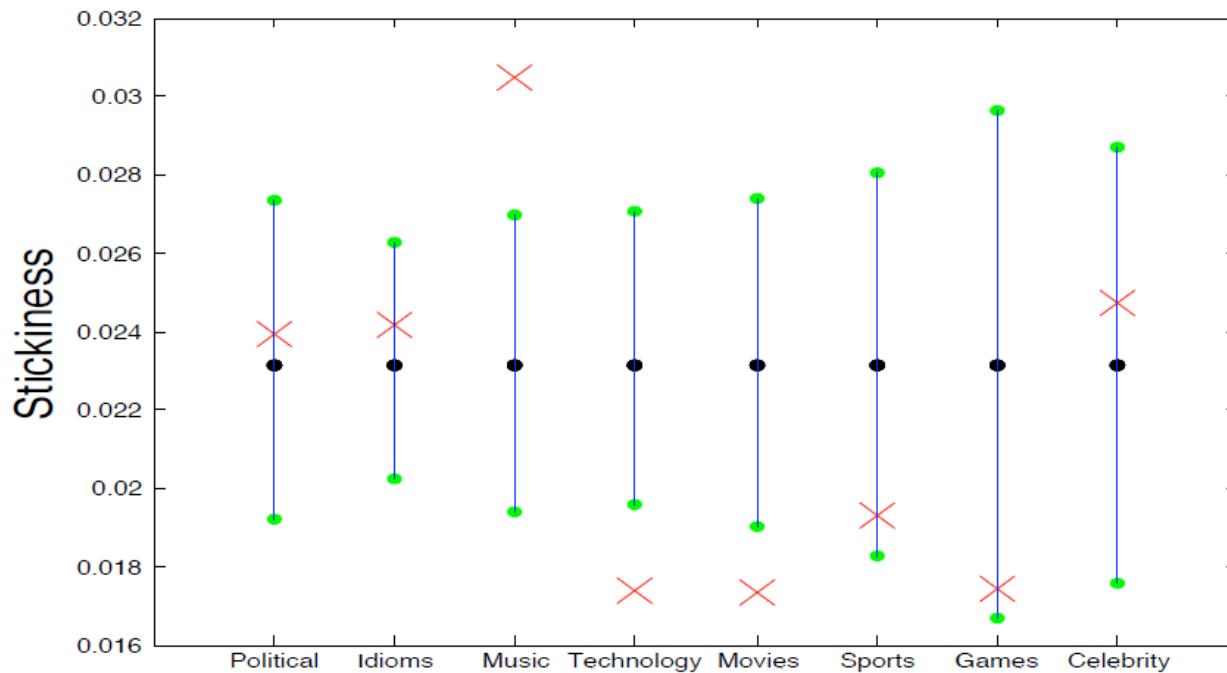
- Manually identify 8 broad categories with at least 20 HTs in each

Category	Examples
Celebrity	mj, brazilwantsjb, regis, iwantpeterfacinelli
Music	thisiswar, mj, musicmonday, pandora
Games	mafiaWars, spymaster, mw2, zyngapirates
Political	tcot, glennbeck, obama, hcr
Idiom	cantlivewithout, dontyouhate, musicmonday
Sports	golf, yankees, nhl, cricket
Movies/TV	lost, glennbeck, bones, newmoon
Technology	digg, iphone, jquery, photoshop



- Idioms and Music have lower persistence than that of a random subset of hashtags of the same size
- Politics and Sports have higher persistence than that of a random subset of hashtags of the same size

# Exposure Curve: Stickiness



- Technology and Movies have lower stickiness than that of a random subset of hashtags
- Music has higher stickiness than that of a random subset of hashtags (of the same size)

# Recap of this lecture

- Basic reproductive number  $R_0$
- **General epidemic models**
  - SIR, SIS, SEIZ
  - Independent cascade model
  - Applications to rumor spread
  - Exposure curves