Announcements:

- Please fill HW Survey
- Weekend Office Hours starting this weekend (Hangout only)
- Proposal: Can use 1 late period

Probabilistic Contagion and Models of Influence

CS224W: Analysis of Networks Jure Leskovec, Stanford University http://cs224w.stanford.edu



Models of Cascading Behavior

- So far:
 Decision Based Models
 - Utility based
 - Deterministic



- "Node" centric: A node observes decisions of its neighbors and makes its own decision
- Require us to know too much about the data
- Next: Probabilistic Models
 - Lets you do things by observing data
 - We lose "why people do things"

Epidemic Model Based on Trees

Simple probabilistic model of cascades where we will learn about the **reproductive number**

Probabilistic Spreading Models

Epidemic Model based on Random Trees

- (a variant of a branching processes)
- A patient meets *d* other people
- With probability *q > 0* she infects each of them
- Q: For which values of d and q does the epidemic run forever?
 - **Run forever:** $\lim_{h \to \infty} P \begin{bmatrix} \text{At least 1 infected} \\ \text{node at depth h} \end{bmatrix}$



|>0

= 0

Die out:

Probabilistic Spreading Models

- *p_h* = prob. there is an infected node at depth *h* We need: lim_{h→∞} *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

$$p_{h} = 1 - \underbrace{(1 - q \cdot p_{h-1})^{d}}_{\text{No infected node}}$$



- $\lim_{h \to \infty} p_h$ = 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: $f(x) = 1 - (1 - qx)^d$



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

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



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, $0 \le x, q \le 1$, 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]!

10/18/17

Fixed Point: When is this zero?



Important Points

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

Only R₀ matters:

- $R_0 \ge 1$: epidemic never dies and the number of infected people increases exponentially
- $R_0 < 1$: Epidemic dies out exponentially quickly

Models of Disease Spreading

We will learn about the epidemic threshold

Spreading Models of Viruses

Virus Propagation: 2 Parameters:

- (Virus) Birth rate β:
 - probability than an infected neighbor attacks
- (Virus) Death rate δ:
 - 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

Susceptible $\xrightarrow{\beta}$ Infected $\xrightarrow{\delta}$ Recovered

dt

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 is: $\frac{dS}{dI} = -\beta SI \qquad \frac{dR}{dI} = \delta I$





Jure Leskovec, Stanford CS224W: Analysis of Networks, http://cs224w.stanford.edu

SIS Model

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



SIS Model



Models flu:

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

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

dI $-\delta I$

Question: Epidemic threshold *t*

SIS Model:

- **Epidemic threshold of an arbitrary**
- graph G is τ, 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:



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

Experiments (AS graph)



Experiments

Does it matter how many people are initially infected?



Example: Ebola



[Gomes et al., Assessing the International Spreading Risk Associated with the 2014 West African Ebola Outbreak, *PLOS Current Outbreaks*, 2014]

10/18/17

Jure Leskovec, Stanford CS224W: Analysis of Networks, http://cs224w.stanford.edu

[Gomes et al., 2014]

Example: Ebola, R_o=1.5-2.0



Read an article about how to estimate R₀ of ebola.

10/19/17

Jure Leskovec, Stanford CS224W: Analysis of Networks, http://cs224w.stanford.edu

[Gomes et al., 2014]

Example: Ebola



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 \boldsymbol{v} with prob. $\boldsymbol{p}_{\boldsymbol{uv}}$

Activations spread through the network!

Independent Cascade Modal

- 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



[KDD '12]

Exposure Curves

Exposure curve:

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





Exposure Curves

From exposures to adoptions

- Exposure: Node's neighbor exposes the node to information
- Adoption: The node acts on the information
 Adoption curve:



Example Application

- Marketing agency would like you to adopt/buy product X
- They estimate the adoption curve
- Should they expose you to X three times?
 Or, is it better to expose you X,
- then Y and then X again?



[Leskovec et al., TWEB '07]

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



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?



[Backstrom et al., KDD '06]

Exposure Curve: LiveJournal

LiveJournal group membership



Exposure Curve: Information



- 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

۵

- 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)

Modeling Interactions Between Contagions

Information Diffusion

So far we considered pieces of information as **independently** propagating. **Do pieces of information interact?**



Jure Leskovec, Stanford CS224W: Analysis of Networks, http://cs224w.stan

Modeling Interactions

- Goal: Model interaction between many pieces of information
 - Some pieces of information may help each other in adoption
 - Other may compete for attention



You are reading posts on Twitter:

- You examine posts one by one
- Currently you are examining X
- How does your probability of reposting X

depend on what you have seen in the past? Contagions adopted by neighbors

(X is exposed by them in order):



We assume K most recent exposures effect a user's adoption:

• $P(adopt X=c_0 | exposed Y_1=c_1, Y_2=c_2, ..., Y_K=c_k)$ Contagion the user is Contagions the user

viewing now.

previously viewed.

Contagions adopted by neighbors:



We assume K most recent exposures effect a user's adoption:

P(adopt X=c₀ | exposed Y₁=c₁, Y₂=c₂, ..., Y_K=c_k)
 Contagion the user is viewing now.
 Contagions the user previously viewed.

Contagions adopted by neighbors:



The Model: Problem

- We want to estimate: P(X | Y₁, ... Y₅)
- What's the problem?
 - What's the size of probability table P(X | Y₁, ... Y₅)?
 = (Num. Contagions)⁵ ≈ 1.9x10²¹
- Simplification: Assume Y_i is independent of Y_j $P(X|Y_1, ..., Y_K) = \frac{1}{P(X)^{K-1}} \prod_{k=1}^K P_k(X|Y_k)$

We apply Bayes theorem twice and use the independence assumption

Goal: Model P(adopt X | Y₁,..., Y_K)

Assume:

$$P(X = u_j | Y_k = u_i) \approx \underbrace{P(X = u_j)}_{\text{Prior infection}} + \underbrace{\Delta_{cont.}^{(k)}(u_i, u_j)}_{\text{Interaction term}}$$

- Next, assume "topics":
 - $\Delta_{cont}^{(k)}(u_i, u_j)$ models the change in infection prob. of u_j given that exposure *k*-steps ago was u_j
 - We estimate P(X) and $\Delta^{(k)}(u_i, u_j)$ by simply counting
 - P(X) ... fraction of people exposed to X that got infected by X
 - Δ^(k)(u_i, u_j) ... P(X) fraction of people first exposed to u_i and then to u_j and then got infected by u_j.

Dataset: Twitter

Data from Twitter

- Complete data from Jan 2011: 3 billion tweets
- All URLs tweeted by at least 50 users: 191k
- Task:

Predict whether a user will post URL X

What do we learn from the model?

How do Tweets Interact?

How P(post u₂ | exp. u₁) changes if ...

- u₂ and u₁ are similar/different in the content?
 - LCS (low content similarity), HCS (high content similarity)
- u₁ is highly viral? Prob. of infection P(u):



Final Remarks

Modeling contagion interactions

• 71% of the adoption probability comes from the topic interactions!

Conclusion

R₀: Epidemics die out if R₀<1</p>

R₀: reproductive number

- Epidemic Threshold: Virus "strength" $s = \beta / \delta$ < τ the epidemic can not happen (it eventually dies out)
- Shape of the adoption curve:

 Modeling interactions between contagions

