SCRATCH | 9/4/14
Plan: did 'broadcast + social reaction | (1. lec 4

do we want

structure
SD
Twitter:
Youtube: Siersdorfer { length
re-entry

Outcomes:
intresting
Controversial : Gben, Siersdorfer
vote (wiki) →
speed dates
Resnik / Prabhakaran
live                          email, [Enron]
                             emal "structure"
                             single - thread was best gen
                             nomination.

OUTLINE for lec 4 on "surface" structure
of conversations.
lec notes start on next pg.

"shape statistics
  · length        [Siersdorfer]
  · for trees: depth, "bushiness"
        [Kumar]
    [ Gonzalez - Bailon, Kaltenbrunner, Banchs '10]
     · dimensions explicitly correspond to conv. aspects  (fig 1)
                                            vs. fig 2. each discussion,
       · single measure
       · note (anon. users)    → render Fig 4.
                                   show txt.
what if you a
add participats explicitly
(3) email might not even be trees
        [Prabhakaran : Rambow '14] - add: drop are actually features!
                                        recipients.
                            forwarding = fork.

① Kumar fig 4, section 4.3
     avg. # of unique authors / thread length  (log-log)
                                                 q/a.

② Backstrom et al
     distribution of # unique authors.
                         ↳ comment length.

add topics explicitly :→ quick mention of topic models.
        - show Resnik slides (for Nguyen et al MLJ '14]
why html in SD -    topics that recur
  crowd             topic shifts        (show example from Nguyen et al paper of shifts)
                      ... relate to power   ↳ "Power of confidence = polls"
        [Prob                              ↳ Enron = org-chart.

discussn can make people adopt more
Xtreme positions / Schkade et al '07
quoted in Kaltenbrunner's 'structur of
political discussions'.

extension today:
Sente tabs
WSDM poster open

: "surface-y" perspectives: what do conv. "look like"
(as opposed to discourse (intentional) structure)

Conversations: thread structure ~~post lengths~~

In our 1st lecture, we distinguished b/wn 2 sets of ~~systems which make~~ lang: social interaction manifest:

- conversation
- broadcast: social effect

We've spent the last two lectures ~~on~~ broadcast: social interaction
doing some preliminary exploration of

Now ~~Today~~, I'd like to turn to the conversation paradigm.

- - - ▶ <repeat>

[ will also survey diff. sites : corpora ]

## Thread structure

Let's assume we're dealing w/ a post ~~comment~~ reply setting:

[ "comment" is ambiguous: try to stick to "post" vs. "reply" ]

~~A basic statistic:~~
~~We took a look @~~

A basic statistic: [ length ]

~~Analyzing~~
Analyzing : mining [Siersdorfer et al 2014] : lots of good overview material
Fig 7.   Chelara
San Pedro
Attingorde
Nejel

Siersdorfer 2014
Youtube ~~comment~~ vs.
Yahoo News
But, then don't allow replies to replies in ~~youtube~~ thread structure
That's ok, it's like FB.
Predict what'll receive replies

- Comparing Yahoo! news comments to ~~youtube~~ Youtube ( ~~not badly~~ noted "cesspool".)
  (show page:) "posts" hidden by default
  ~~domains~~ "replies" hidden by default
  lots of ~~comments~~ arriving as we speak! ~~a lot of~~ ~~rare to get a~~

easier to make new post

affordances:
hackernews vs. ←easier to make new post
slashdot   ←easier to make new ~~comment~~ reply
Elsner ~~Rosé~~ ~~Rosé~~ paper on affordances:
SD vs. LiveJournal

- Histogram of lengths: as expected : not
  ~~table X value it:~~ "unreplied" vs "seeds" (≥ 1 reply)
  rare to get a convo going ( >50% posts get no comment )

note highlight : control for politics category — still Youtube threads are shorter

" what does a long thread correspond to"
Fig 9 : diff. sign of deriv: for how "good" initiator of long thread is.

~~For Yahoo, is it~~
Does Yahoo show that ~~more is~~ more length = more interest ⇒ more ⊕ ratings? ( i.e., ⊕ bias combined w/ ~~attention~~ attention ?
No. Youtube is the other way. (flame war)

likes vs # comm.

reviews:

viz of reddit.

q: ⊘ really need error bars on fig 9. esp since right-hand bars have little data.

Now what if we consider conversations as trees (explicitly, or, via thread induction, like what Elsner & Charniak '10 did)? Then we have another "dimension" besides length:
or Wang/Joshi/Rosé

2 class members were from Brown! Also [Wang...

Kumar, Mahdian, McGlohon KDD 2010, Dynamics of Conversations.

~~Twitter~~

`UseNet groups` ; Yahoo groups.   } graphs for Usenet (others are qualitatively similar)

(*Twitter)   q: "why Usenet? < which is email, btw >
            scan up to part before §3:
↓           "While Usenet is declining in pop., ... public, easy
important conv. source we          to crawl obv. thread structure ... some groups
haven't talked about yet.           still active". This is the rationale.

remark about "footnotes in reaction to reviewers": e.x.:
"Recall that X is not y ". A social interaction!

Fig 1b: empirical size vs. depth "fits a √ law (given the exponent value in the
                                                              q: about corpus of reviewer-
key).                                                          author "interaction"?
conversation
=> trees are "deeper" than a "rich-get-richer" would predict (log-depth)
^ (according to highlighted text)

(relate to arXiv:
diff. versions;
laTeX source
available)

* counter-claim: the interfaces for these settings don't
rank by branching; so an 'attention bias' towards   or etc.
the rich-get-richer ^ wouldn't happen by affordance.
→ A/B test w/ HackerNews vs. Slashdot? See Wang/Joshi/Rosé:
                                          SD vs. LJ Sim = ~~stream~~ reply "stream"
· ~~you~~ as they say, you can hide an elephant in a log-log plot.
(checked if my signing in:         revealing my cluster login).

3
Fig ~~2~~: do diff. levels have diff. or same branching?
Lines slope down b/c you expect fewer things to have large branching.
~~But~~ ~~B~~ Point is diff. btwn lines: higher levels (towards tree top) have more branches.

reminder:
... all of this: about what these trees look like.

Depth & other measures of tree "shape"

Gongalez-Bailon, Kaltenbrunner, Banchs : "structure of political ..." - on slashdot !
( for structure w/out participants explicitly labeled

abused the term "branching factor" a # of times. (cf. "branching")

max & branching                                              <scratch>

h-index = 2 measures of depth : "weight in controversy of certain comment."
                               [engagement of most active users", not relevant for lec.
                                outline ]

# of layers  X ⟹ x layers of @ least x comments in that layer

Fig 1 : again start w/ theory about (useful discussions), Fig 1.
         Vertical axis : "argumentation level", is practically "literally" depth of the tree
               better discussions have participants involved.    (nom'd : can't have negative
                                                                            depth! )

         horiz : democratic discussions involves more people
              ... ⟹ more people pile in @ a particular post, so correlate it w/ branching
              [ later, also explicit consider # of unique participants in the tree ]
                  ↑
                in pgph

         Scan thru : note highlights on "why Slash Dot" (remark: "these people really like Slashdot" )
              "had enough time to evolve & consolidate, overcoming the problems associated to
               spam or misbehaviour and proving its robustness"
                : also about handling of anon posts ("Anonymous Coward" is not a single
                  person.)

                  - don't remove - would break the tree structure.

Fig 2 : examples of 2x2 plane (presumably real slashdot trees)
         - highlight : quadrant I is the "good" one.

Fig 3 : plot political v. non-political on that plane.
         "intersection" of dotted lines : mean width & depth.
          Centroid shows the political discussions do indeed "look diff" - in Quad I -
             from non-political

         (Fig 4 : divide by category. Not in color, unfortunately.
              linux: near center
              games: very low participation (Quad III) : fewer people in less active category? )

An h-index measure, to try to summarize shape in a single #.
         - see highlighted text, under fig 4.

              also, here's translation ...
                  (note layers w/ big branching don't have to be consecutive)

Now add: participant ID to ~~there~~ conversation structure
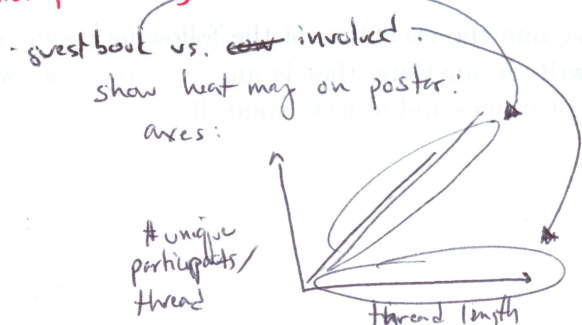
~~Kumar~~

<span style="color:red">Kumar et al again, fig 4 (and section 4.3):</span>

On a log-log, we see that the <u>avg</u> # of unique authors (red) vs. size:
a polynomial relationship.

<u>But</u>, what about the <u>distribution</u> of # of unique authors?

(.., cont ~~of "what is a long thread"~~)

<span style="color:red"><Backstrom, Kleinberg, ~~Lee~~, Panesu-Niculesu-Mizel WSDM '13></span> : <span style="color:red">(display the WSDM poster)</span>

- guestbook vs. ~~conv~~ involved.
  show heat map on poster:
  axes:



# unique
participants/
thread

thread length

- ⊙. Why's low? not
  (dyn. of conv.
  paper

SP: deals w/ avg
  comments
argue that width is good
proxy for # of ~~psych~~ involved.

→ ~~crowtin~~ says
  this my not
  be true.

business:
h-index.

→ handi of anon.

matters b/c you may want to rate posts by type of
  thread, not just length of thread
  - if it's a thread someone's likely to re-enter,
    keep them apprised.

• Also, since it came up before, note that social-connectedness
  of 1st few posters has interesting effects on the length
  of a reply thread vs. a "like" thread

Prabhakran, Enron
~~-ca add~~
Cold (enron is not for
  ~~the third big classify)~~

<span style="color:red">- in the setting of email, you may not even have trees, "because of" addition of ~~another's~~ people's
  identities</span>

Enron data set      (post ref)
     → fork
forwarding vs. replying      <need to clear at the "prior email" text or
     make sure you know ~~of~~ who wrote which parts>
dropping/adding recipients
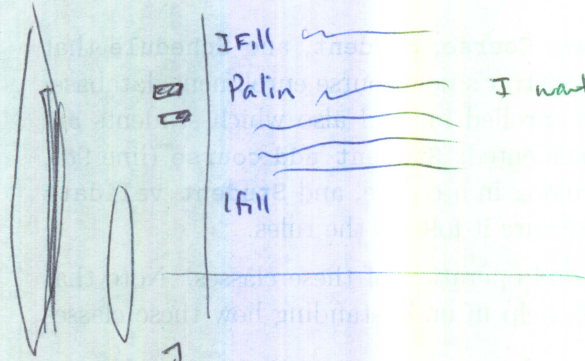power relationships known.
     → (actually used as features in Prabhakaran; Rambow
        '14)

Note also the UMass work on the North Carolinian corpus (Freedom of
  Info request?)

<span style="color:blue">Note also the W3C data (Glasgow.</span>

. Just talked about adding recipients, what about now adding context to the structure of conv?

<Nguyen, Boyd-Graber, Resnik, Cai, Midberry, Wang MLJ '18 > 95(3):381-421
14
(using Philip's slides @ the WESP on Lang & Soc Sci)

Debate
- real-time,
f2f
conversation.

IFill
Palin          I want to talk about, again, my record on energy

IFill

use topic modeling <post ref> to show which topics are active when.
Some recur throughout, some only for a short time.

when does the topic shift?
From text, Palin is abruptly shifting to a diff. topic.
IFill tries to get back to other topics.

next slide: moderator & Palin changing topic about the same amt.

next slide: criticism of IFill for not managing the conversation better

next slide: a diff (set?) of debates.

Prabhakaran & Rambow '14 - topic shifts & related to poll standings.