

(make copies)
 2 handouts:
 outline (the 'red parts' @ top of each sheet)
 some helpfulness cues from other papers
 Set up Santa tabs for figs, talk slides for other figures.

fake, joke, sarcasm!

Review "quality/helpfulness":

reviews considered independently

~~reviews are not~~ potentially also considering context

(diff. connotations)

[Intro] (1) prediction ~~in and of itself~~

features, techniques for ~~labels~~ labels. (features helpful for projects)

[Focus] (2) ~~provide~~ a lens for studying social influence ~~interesting techniques~~

(naturally, will talk about techniques)

We talked about review helpfulness last time (remember we ~~and thought~~ wrote down some features of helpfulness and evaluated)

Helpfulness ^{are} ratings being used as a navigation tool by sites like Amazon, since so many reviews

When these evaluations are produced by users, you have social navigation. (is this from the OK reading?)

You could ask, why do this "socially" if it ~~can~~ be done automatically, and needless to say, there's been quite a bit of ~~energy~~ research energy devoted here.

w.r.t. considering this to be a prediction problem, there are a # of options:

- the "given" label is "x out of y" found this helpful.
 - regression problem?
 - binary classification (threshold on ratio x/y)
- what features?

Ottobacher (have handout; also project using desktop 1)

→ First example: ^{Johna} Ottobacher '09 'Helpfulness' ~~as a measure of online~~, Table 3 (pg 958)

(use reference tab, so people can see pub info (screen width/resolution permitting))

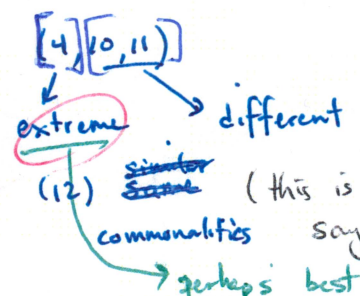
numerous better.

I wanted to show this table b/c it divides the description nicely into the "general concept" of each feature vs. how it was implemented.
 Also, the thing is super-convenient. (so, in terms of presentation, this is nice)
 • These are what she tried, not all of which turned out to be significant
 And, ~~not all~~ she was not the first to introduce many of these
 Link to survey of previous work on the course homepage.

Side notes:
 You might want to discuss a few for your plot slide
 in the context of review mining.

4th column (Explanation/Justification is perhaps easiest to go through?)
Observations:

intrinsic quality: (13, 14, 16, 17) = readability
(1-3): "fact" commonly used for coverage/objectivity of a review
Note: citation to [6] is probably wrong.



g: perplexity?
q: should one expect 'extreme' to be ↑ or ↓ correlated? (can do regression/classification w/ hypothesis)
let p_1, p_2 two distributions identified with two samples s_1, s_2 . (We'll blur the distinction). $H(p_1, p_2) = -\sum p_1(x_i) \log(p_2(x_i))$. p_1 is considered "true", p_2 an "alternate code". Perplexity is 2^H , and converts to bits as unit.

(12) similar same (this is like the Gilbert, Karahalios found some of their interviewees saying: "A completely unique review wouldn't serve any purpose")
perhaps best to consider extreme w/ other review ratings
[5-9]: ~~prop~~ author properties (some authorial, some may be for 'smoothing from past')
contextual rev: "as we saw of our example last time, sometimes helpfulness depends on other reviews, if you

seems to be a nod to the idea that relationship to other reviews can matter. But, the actual features in this category don't seem that linked to that concept.

~~note that Amazon didn't use to sort reviews by helpfulness, we think.~~

Looking @ the 2011 version of Ghose, Ipeirotis '11, table 1.
as another example of features that have been tried.

(note: some of these features were used to predict sales rank, so things like ~~is the reviewer a star~~, helpfulness rating not used in helpfulness prediction)

more reviewer characteristics
more ~~for~~ readability features (5th row)

interest in "SMOG" the "Simple Measure of Gobbledygook" Explained constant as result of fitting to data.

* subjectivity measures, using same idea that the product description can be treated as ~~is~~ definitely objective material (as it happens, to train the subj. detector)
~~comment about desprpb: add to get confidence for subj. prob, or <stg smart I can read~~

And other studies look @ similar features, as well.

desprpb: relate to the VBC slashdot/thread visualizer we talked about before.

~~So, these were the kinds of things that were in the air~~

Any comments/questions, esp. re: people's potential pilot studies for AI? recommendations

one q: finding a very specific feature that few people have talked about, like the fault in the review example we saw last time?

Discussion re: fake reviews, sarcasm, etc: note that machine classification results are surprisingly high.
q: are fake reviews 'product-dependent'?

How well do such features (esp text-only features) do?

J. Liu et al, EMNLP-CoNLL

Re-annotate reviews according to own written standard b/c of biases:

Mismatch between
- bias toward writing as 'helpful' (half of 23k sample has >80% helpful)
- some reviews don't get eval'd (early bird bias) rich-get-richer
(temporal: earlier get more.)

Are there social factors behind mismatch?

(2) Study Using quality/helpfulness as a lens on social influence

→ focus of this course

[Sipos, Ghosh, Joachims WWW '14]: (mis)-~~ranking~~ ranking by community
(to some degree - also proprietary Amazon factors: tau = .84 for
revs w/ ≥ 10 votes)

topic

stress the daily-snapshot
idea

- "true" quality by "final" ~~future~~ ranking
4 mos in future
Attempts to avoid self-fulfilling prophecy labels
- effect on helpfulness vote, and on whether they vote.

for helpfulness ratio vs. actual Amazon ranking

technique for biased labels

→ correction of mis-ranking
color = prob of next eval coming in as 'helpful'
point out 'high' rank = 4.

- g: rc: how to measure participation decision:
- paper assumes "constant # of pageviews (for each position) in each period between snapshots".
 - Doing relative # of pageviews then should be ok. (take it or leave it).
 - also a measure of 'participation' to handle smoothing of sparse data.
 - motivated to participate to correct mis-ranking (of paper)

Note a point of the message: this is an 'opposite' effect to how clicks on rankings of web-search results accumulate to the already-highly ranked.

- [non-ranking sitting] w/d.
 - overall, seeing other \oplus ~~was~~ increased prob. of \oplus eval. \ominus : not much effect.
 - ~~identity~~ ~~quality~~ ~~two~~ ~~quality~~
 Muchnik, Aral, Taylor, Science 2013 - manipulation study (~~identity~~ ~~quality~~ ~~two~~ ~~quality~~)
 comment rating, so ranking is not a factor (this is not Amazon).
 opt-in

the Web site performed "best" by the experimenters, therefore

Artificially seeded the comments' initial vote w/ a \oplus or \ominus @ evaluation, \oplus or \ominus comments altogether
 4019 \oplus 1942 \ominus
 reflect "natural" prior
 lots of "no initial eval" as control.

most demonstrative
 < Fig. 1.c > (don't show for time just make ref. on web?)
 explain bias on confidence interval

data analysis

so, perhaps not so relevant in terms of students being able to use

- note asymmetric effect: herding for \oplus .
- use the "no artificial 1st eval" as baseline.

[Daneson-Niculescu-Mizel, Kossinets, Kleinberg, Lee WWW '09] : effect of conformance to group opinion
cultural diff in

: natural experiment: "plagiarized" reviews to sidestep true quality

clipped slides.

explain what grey bars are (and expect, since for 5-star avg no one can be +4 stars away)

not just "more (+) = more helpful", but "for same degree of deviation, better to be (+)

Note that wrt one ~~review~~ version of a plagiarized review being more appropriate, the example plagiarized pair show the latter to be more helpful, and our studies showed no sig diff of helpfulness for @ the pair.

[Cheng, Danescu-Niculescu-Mizel, Leskovec ICWSM '14] = effect of evaluations on the author
: propensity matching; ~~to~~ 'natural' experiment, but requires measuring text quality to pair 'similar' reviews.

post quality: evals drop significantly after a neg. eval.

rise is not significant after pos. eval.

Other controls in pairing: # words, # posts written, general ^{quality} ~~helpfulness~~ evaluations previously).

fig 4 in site .

* = sig diff.

g: maybe the people who got \ominus eval were trolls, and so were motivated to write worse posts b/c that's what ~~they~~ trolls do?

- interesting that people who received "no" feedback most likely to drop out