

CS341 Journal

Jo, Dilli and Howon

April 14, 2014

Jo: 4-2-2014

This is a lil journal system I like using. Basically, you put a new file whose name is the date (this one is April 2nd, that is, `4-2.tex`) with your journal entry. Then a composite pdf is generated from the python script `GenerateJournal.py`. Its simple and offline, so we still have a collaborative thing but without the annoyances of sharing a google doc. The pdf file could be generated during the pre-commit hook, if you like that sort of thing.

Howon: 4-6-2014

More resources, I found. More thoughts about the data, I had.

<https://dev.twitter.com/docs/platform-objects> is probably most important for not having to go through the json fields themselves. We will have the problem of nulled fields, empty sets, and this field guide (get it? birds?) has the defaults.

<https://github.com/kevinweil/elephant-bird/> is a good resource for dealing with json in hadoop

The dataset is about 400 gb. It seems like a 1% sample would be good, as it allows us to fit it on main memory even on a crappy box like mine. This I approximated a little bit by taking a 1-file sample, which is actually just about a gigabyte: ideally, we would use something like `TABLESAMPLE` in Hive or `SAMPLE` in Pig Latin, because then it would be a good pseudorandom sample.

One troubling thing about taking from the *firehose*, at least at first glance, is that the tweets are taken at the moment of conception. So it's much more annoying to get the favorites data on the tweets.

Histograms about the facts about the users I found are also forthcoming.

Howon: 4-8-2014

Todos

Howon

1. Email the guy about getting the graph data, so we can actually sample the users and we need that graph for other stuff
2. Maybe do an actual sample to think about distributions myself to unbias the tweets (maybe we're ok with just getting people who tweet more?)

Dilli

1. Keywords for one file
2. Think about reconstructing the graph

Jo

Work on implementation of the TwitterRank and the topic distributions.

Brainstorming.

Some ideas may be worse than others.

1. We could have more users
2. We could construe each tweet as a separate document
3. Think of a 3-layer bayes net, where user has distribution over topics and the tweet has distribution over topics where prior is over user's distribution over topics.
4. Find a user that does poorly in twitterrank even though they do well in real life or in our algo (someone strong in 2 different fields)
5. What if we pick our own topics, and do actual labelling?
6. Syntax? Stop words? Use NLTK
7. Can we reconstruct the retweeting graph and do everything in the proposal first?

Jo: 4-9-2014

Topic Distillation

The module `Topic.py` has a function `topic_model` which takes a number of topics and a list of word count histograms for some documents and performs LDA on it to assign each document to a topic and to topic a distribution of words. The TwitterRank paper recommends running this where each document is all words of all a user's tweets. We can try this with users picked in various ways and words chosen in various ways. Does reasonably on the generative model in `TopicTest.py`.

Everything runs in-memory and not distributed. We should be able to fix this so that the histograms actually live on disk and are streamed when necessary, but I don't know if we'll need to yet. It depends on the size of the Twittersphere we wish to rank and the number of words considered.

This analysis is required to compute the edge weights for the GraphRank algorithm.

GraphRank

Given as input is a JSON graph where each node has its weighted outgoing edges and teleportation probability. Outputs the weighted graph rank. The TwitterRank paper recommends running this once per topic where the edges are weighted by topic similarity and number of tweets and the teleportation probabilities are the topic ratings for each user. We should try to use the explosive retweetiness for edge probabilities.

TODO

Link these together (even on a fake graph), and make this happen on a fraction of our data. Probably will require a map reduce run to get histograms.

Dilli: 4-9-2014

Update on Preprocessing

I have added a sample retweet and scripts to compute wordcounts and follower graphs.

Howon: 4-9-2014

Notes from Jeff

We should seriously consider taking only a few of the topics, for the same reason the TwitterRank people (for parameter reduction). But you might also think about the combination of what people tweet about and the person themselves.

That is, there might be some relation between what a person's followers are and what they tweet about. Eg., the second part of the TwitterRank thing, with the PageRank and weighting.

Jo: 4-14-2014

0.1 Multi-Topic TwitterRank and Lots of Parameters

We've mentioned parameter reduction a few times and been afraid of introducing lots of parameters in our model, but in this context I don't think it's important.

There is an obvious model for multi-topic twitter rank. Here is how it works. To generate user's corpus of tweets, we first sample k topics from the overall topic distribution, possibly with weights. Then, for each tweet, the user samples a topic t from their k topics and for each word in that tweet they sample a word from t 's distribution.

Let T be the number of topics, W the number of words, U be the number of users and N be the number of tweets. Then this model has $T(1 + W) + kU + N$ parameters as opposed to the $T(1 + W) + U$ parameters of the original model. Here, N is HUGE. However, I don't think it matters.

There is another Bayesian process with only $T(1 + W) + kU$ parameters to learn, but the posterior of these parameters are exactly the same as the posterior of the network described above. We get this process simply by marginalizing out the topic-per-user parameters, meaning that the conditional distribution of the words per tweet is more complicated, but this doesn't matter. The point is that the posterior estimate for the topic-per-user, which is what we care about, is no better than if we had all those extra parameters, so (unless it becomes a computational problem) parameter reduction is actually not an issue here, so long as we don't actually use the topic-per-tweet estimates. Thus coming up with a better model is actually super easy, and some variant of the one described above will probably work.

0.2 Priors

We shouldn't use flat priors for the topic distributions since because symmetry. We expect some topics to be more popular than others (even if we don't know which), so it might be wise if the priors reflect that. With a flat prior, the true posterior is symmetric about switching any two topic labels, which means there will be ties in the MAP estimates. This means our estimate of the MAP might be problematic, better to just avoid this ambiguity.