



Modeling Online Discourse with Coupled Distributed Topics



Nikita Srivatsan, Zachary Wojtowicz, Taylor Berg-Kirkpatrick

nsrivats@cmu.edu, zdw@andrew.cmu.edu, tberg@cs.cmu.edu

ArXiv

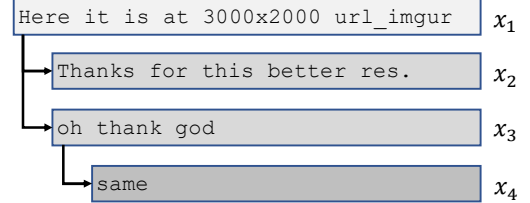
<https://arxiv.org/abs/1809.07282>



Summary

We propose a deep, globally normalized topic model that captures discursive interactions along observed reply links, typical of online data, in addition to traditional topic information. Our model incorporates latent distributed representations arranged in a deep architecture, which enables efficient GPU-based variational inference. We apply our model to a new dataset of 13M comments mined from Reddit.

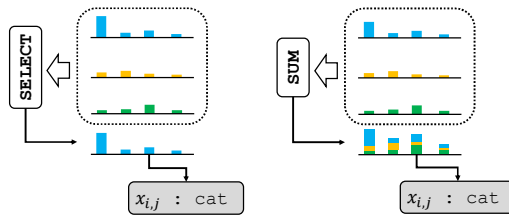
Online Discourse



Visualization of a branching Reddit thread with observed reply links

Distributed Topics

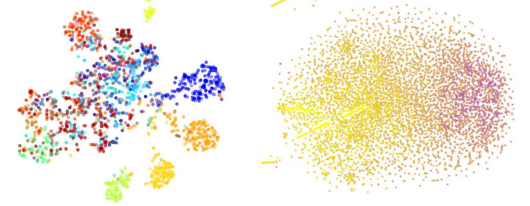
Discrete (LDA) vs distributed (RS) topic modeling approaches for explaining some particular word



Analysis

t-SNE plot of thread-level representations colored by subreddit of origin

t-SNE plot of comment-level representations colored by comment length



Strongest word emissions for particular thread-level topic bits (single bit active)

Bit 1
maduro
venezuelan
ballot
puerto
catalonia
rican
quak
skateboard
venezuela
quebec

Bit 2
comey
pede
macron
pgl3
maga
globalist
ucf
committe
cuck
distributor

Bit 3
btc
gameplay
tutori
cyclist
dev
currenc
kitti
bitcoin
rpg
crypto

Strongest word emissions for particular comment-level topic bits (single bit active)

Bit 1
irl
riamverysmart
legend
omfg
riski
aboard
favr
madman
skillset
tunnel

Bit 2
faq
tldr
pms
165
til
keyword
questions
feedback
chat
pm

Bit 3
funniest
mah
tfw
teleport
fav
hoo
plz
bah
whyd
dumbest

Strongest word emissions based on inferred latent comment-level topic representations (multiple bits active)

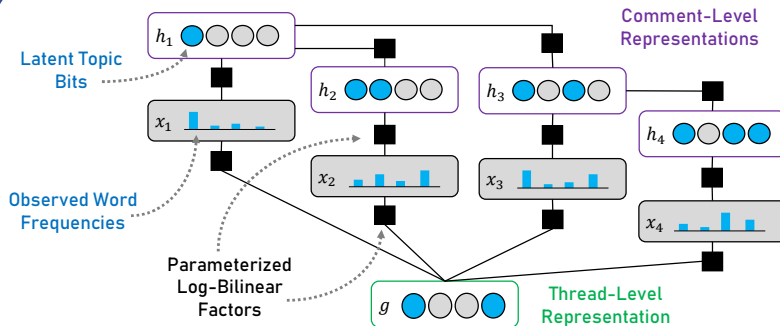
Embedding 1
reev
keanu
christoph
murphy
walken
vincent
chris
til
wick
roger

Embedding 2
reddit
shill
question
background
user
subreddit
answer
relev
discord
guild

Embedding 3
moron
douchebag
stupid
dipshit
snitch
jackass
dickhead
idioci
hypocrit
riddanc

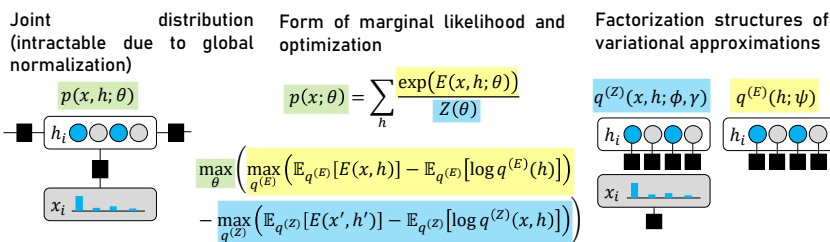
Qualitative analysis of strongest word emissions for selected topic bits

Model



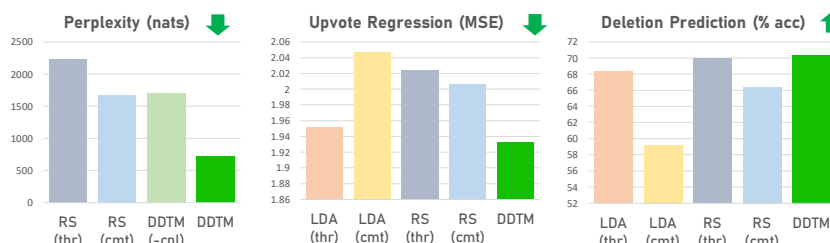
Factor graph visualization of Discursive Distributed Topic Model (DDTM)

Learning



Likelihood approximation breaks down into two variational lower bounds

Results



Performance of our model (DDTM), with and without (-cpt) coupling potentials, vs baselines at both thread (thr) and comment (cmt) resolutions