

Modeling Social Annotation Data with Content Relevance using a Topic Model

T. Iwata, T. Yamada and N. Ueda
NIPS 2009

presented by J. Silva, Duke University

Introduction

- Context: social annotations, aka collaborative tagging or folksonomy
- Users freely annotate objects such as webpages, photos, blog posts, videos, music and scientific papers
- Examples: Delicious, Flickr, Technorati, YouTube, Last.fm, CiteULike

Introduction

- Problem: users often write *noisy*, or content-irrelevant annotations (e.g. “great”, “to read”)
- This paper proposes a generative model for topics and annotations that takes into account relevance/irrelevance of the annotations
- It is an extension of Blei and Jordan’s Correspondence Latent Dirichlet Allocation (Corr-LDA), which assumes the annotations are always relevant

Proposed Method

- We have a set of D documents, and each consists of a pair of words and annotations $(\mathbf{w}_d, \mathbf{t}_d)$, where $\mathbf{w}_d = \{w_{dn}\}_{n=1}^{N_d}$ and $\mathbf{t}_d = \{t_{dm}\}_{m=1}^{M_d}$

Table 1: Notation

| Symbol | Description |
|----------|--|
| D | number of documents |
| W | number of unique words |
| T | number of unique annotations |
| K | number of topics |
| N_d | number of words in the d th document |
| M_d | number of annotations in the d th document |
| w_{dn} | n th word in the d th document, $w_{dn} \in \{1, \dots, W\}$ |
| z_{dn} | topic of the n th word in the d th document, $z_{dn} \in \{1, \dots, K\}$ |
| t_{dm} | m th annotation in the d th document, $t_{dm} \in \{1, \dots, T\}$ |
| c_{dm} | topic of the m th annotation in the d th document, $c_{dm} \in \{1, \dots, K\}$ |
| r_{dm} | relevance to the content of the m th annotation of the d th document, $r_{dm} = 1$ if relevant, $r_{dm} = 0$ otherwise |

Note: all tables and figures taken from the original paper

Proposed Method

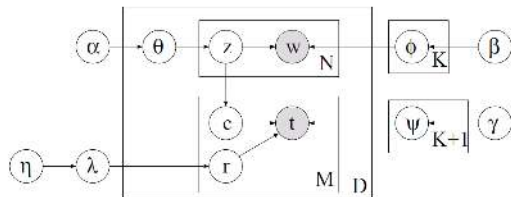


Figure 1: Graphical model representation of the proposed topic model with content relevance

- There are $K + 1$ annotation distributions (Ψ), since Ψ_0 is a topic-unrelated distribution that applies to irrelevant annotations

Proposed Method: Generative Model

1. Draw relevance probability $\lambda \sim \text{Beta}(\eta)$
2. Draw content-unrelated annotation probability $\psi_0 \sim \text{Dirichlet}(\gamma)$
3. For each topic $k = 1, \dots, K$:
 - (a) Draw word probability $\phi_k \sim \text{Dirichlet}(\beta)$
 - (b) Draw annotation probability $\psi_k \sim \text{Dirichlet}(\gamma)$
4. For each document $d = 1, \dots, D$:
 - (a) Draw topic proportions $\theta_d \sim \text{Dirichlet}(\alpha)$
 - (b) For each word $n = 1, \dots, N_d$:
 - i. Draw topic $z_{dn} \sim \text{Multinomial}(\theta_d)$
 - ii. Draw word $w_{dn} \sim \text{Multinomial}(\phi_{z_{dn}})$
 - (c) For each annotation $m = 1, \dots, M_d$:
 - i. Draw topic $c_{dm} \sim \text{Multinomial}(\{\frac{N_{kd}}{N_d}\}_{k=1}^K)$
 - ii. Draw relevance $r_{dm} \sim \text{Bernoulli}(\lambda)$
 - iii. Draw annotation $l_{dm} \sim \begin{cases} \text{Multinomial}(\psi_0) & \text{if } r_{dm} = 0 \\ \text{Multinomial}(\psi_{c_{dm}}) & \text{otherwise} \end{cases}$

Proposed Method: Inference

- The joint distribution is $P(\mathbf{W}, \mathbf{T}, \mathbf{Z}, \mathbf{C}, \mathbf{R} | \alpha, \beta, \gamma, \eta) = P(\mathbf{Z} | \alpha) P(\mathbf{W} | \mathbf{Z}, \beta) P(\mathbf{T} | \mathbf{C}, \mathbf{R}, \gamma) P(\mathbf{R} | \eta) P(\mathbf{C} | \mathbf{Z})$

- In the expression above, we have:

$$\mathbf{W} = \{\mathbf{w}_d\}_{d=1}^D, \mathbf{T} = \{\mathbf{t}_d\}_{d=1}^D, \mathbf{Z} = \{\mathbf{z}_d\}_{d=1}^D, \mathbf{C} = \{\mathbf{c}_d\}_{d=1}^D, \\ \mathbf{c}_d = \{c_{dm}\}_{m=1}^{M_d}, \mathbf{R} = \{\mathbf{r}_d\}_{d=1}^D, \mathbf{r}_d = \{r_{dm}\}_{m=1}^{M_d}$$

- θ, Φ, Ψ and λ are marginalized out

Proposed Method: Inference

- $P(\mathbf{Z}|\alpha) = \prod_{d=1}^D \int P(\mathbf{Z}|\boldsymbol{\theta}_d)P(\boldsymbol{\theta}_d|\alpha)d\boldsymbol{\theta}_d = \left(\frac{\Gamma(\alpha K)}{\Gamma(\alpha)^K}\right)^D \prod_d \frac{\prod_k \Gamma(N_{kd}+\alpha)}{\Gamma(N_{kd}+\alpha K)}$
- $P(\mathbf{W}|\mathbf{Z}, \beta) = \left(\frac{\Gamma(\beta W)}{\Gamma(\beta)^W}\right)^K \prod_k \frac{\prod_w \Gamma(N_{kw}+\beta)}{\Gamma(N_{kw}+\beta W)}$
- $P(\mathbf{T}|\mathbf{C}, \mathbf{R}, \gamma) = \left(\frac{\Gamma(\gamma T)}{\Gamma(\gamma)^T}\right)^{K+1} \prod_{k'} \frac{\prod_t \Gamma(N_{k't}+\gamma)}{\Gamma(N_{k't}+\gamma T)}$, where $k' \in \{0, \dots, K\}$
- $P(\mathbf{R}|\eta) = \frac{\Gamma(2\eta)}{\Gamma(\eta)^2} \frac{\Gamma(M_0+\eta)\Gamma(M-M_0+\eta)}{\Gamma(M+2\eta)}$
- $P(\mathbf{C}|\mathbf{Z}) = \prod_d \prod_k \left(\frac{N_{kd}}{N_d}\right)^{M'_{kd}}$
- Inference of the latent $\mathbf{Z}|\mathbf{W}, \mathbf{T}$ is done using collapsed Gibbs sampling
- The hyperparameters are estimated by maximizing the joint distribution, using a fixed-point iteration method

Proposed Method: Inference

- We have the following expressions, where $j = (d, n)$, $i = (d, m)$ and $\setminus j$ denotes the count excluding the n -th word in the d -th document

$$P(z_j = k | \mathbf{W}, \mathbf{T}, \mathbf{Z}_{\setminus j}, \mathbf{C}, \mathbf{R}) \propto \frac{N_{kd\setminus j} + \alpha}{N_{d\setminus j} + \alpha K} \frac{N_{kw_j\setminus j} + \beta}{N_{k\setminus j} + \beta W} \left(\frac{N_{kd\setminus j} + 1}{N_{kd\setminus j}} \frac{N_d - 1}{N_d} \right)^{M'_{kd}}$$

$$P(r_i = 0 | \mathbf{W}, \mathbf{T}, \mathbf{Z}, \mathbf{C}, \mathbf{R}_{\setminus i}) \propto \frac{M_{0\setminus i} + \eta}{M_{\setminus i} + 2\eta} \frac{M_{0t_i\setminus i} + \gamma}{M_{0\setminus i} + \gamma T},$$

$$P(r_i = 1 | \mathbf{W}, \mathbf{T}, \mathbf{Z}, \mathbf{C}, \mathbf{R}_{\setminus i}) \propto \frac{M_{\setminus i} - M_{0\setminus i} + \eta}{M_{\setminus i} + 2\eta} \frac{M_{c_i t_i\setminus i} + \gamma}{M_{c_i\setminus i} + \gamma T}.$$

$$P(c_i = k | r_i = 0, \mathbf{W}, \mathbf{T}, \mathbf{Z}, \mathbf{C}_{\setminus i}, \mathbf{R}_{\setminus i}) \propto \frac{N_{kd}}{N_d},$$

$$P(c_i = k | r_i = 1, \mathbf{W}, \mathbf{T}, \mathbf{Z}, \mathbf{C}_{\setminus i}, \mathbf{R}_{\setminus i}) \propto \frac{M_{kt_i\setminus i} + \gamma}{M_{k\setminus i} + \gamma T} \frac{N_{kd}}{N_d}.$$

Experiments

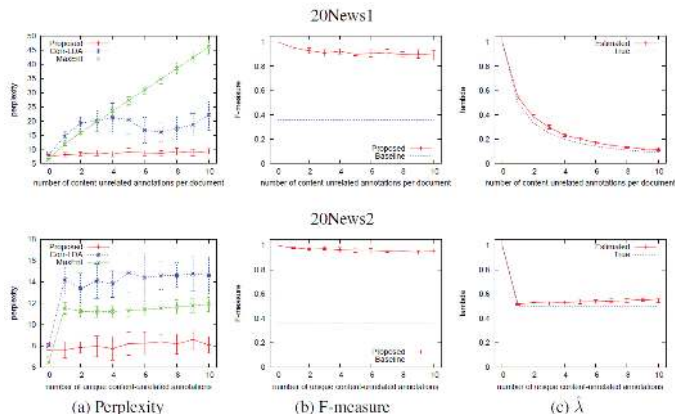
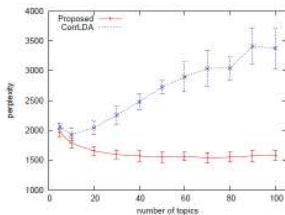


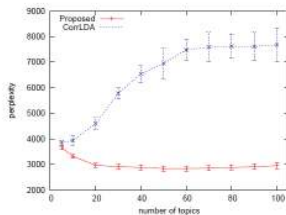
Figure 2: (a) Perplexities of the held-out content-related annotations, (b) F-measures of content relevance, and (c) Estimated content-related annotation ratios in 20News data.

Figure 2 (c) shows the content-related annotation ratios as estimated by the following equation. $\hat{\lambda} = \frac{M - M_0 + n}{M + 2n}$, with the proposed method. The estimated ratios are about the same as the true ratios.

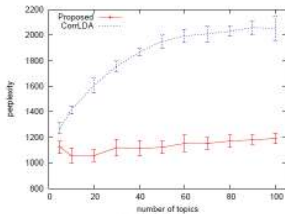
Experiments



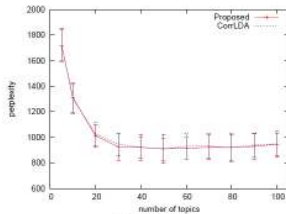
(a) Hatena



(b) Delicious



(c) Flickr



(d) Patent

Figure 3: Perplexities of held-out annotations with different numbers of topics in social annotation data (a)(b)(c), and in data without content unrelated annotations (d).

Experiments

Table 2: The ten most probable content-unrelated annotations (leftmost column), and the ten most probable annotations for some topics (other columns), estimated with the proposed method using 5C topics. Each column represents one topic. The lower half in (a) and (b) shows probable words in the content.

(a) Hatena

| unrelated | Topic1 | Topic2 | Topic3 | Topic4 | Topic5 | Topic6 | Topic7 | Topic8 | Topic9 |
|-----------|-------------|-------------|-------------|---------------|----------|------------|---------------|-------------|--------------|
| toread | programming | game | economics | science | food | linux | politics | pc | medical |
| web | development | animation | finance | research | cooking | tips | international | apple | health |
| later | dev | movie | society | biology | gourmet | windows | oversea | iphone | lie |
| great | webdev | Nintendo | business | study | recipe | security | society | hardware | government |
| document | php | movie | economy | psychology | cook | server | history | gadget | agriculture |
| troll | java | event | reading | mathematics | life | network | china | mac | food |
| " | software | xbox360 | investment | pseudoscience | fooditem | unix | world | cupidity | mentalhealth |
| ? | ruby | DS | japan | knowledge | foods | mysql | international | technology | mental |
| summary | opensource | PS3 | money | education | alcohol | mail | usa | ipod | environment |
| memo | softwaredev | animation | company | math | foodie | Apache | news | electronics | science |
| | development | game | year | science | eat | in | japan | yen | rice |
| | web | animation | article | researcher | use | setting | country | product | banana |
| | series | movie | finance | answer | omission | file | usa | digital | medical |
| | hp | story | economics | spirit | water | server | china | pc | diet |
| | technology | work | investment | question | decision | case | politics | support | hospital |
| | management | create | company | human | broil | mail | aso | in | poison |
| | source | PG | day | ehara | face | address | mr | note | eat |
| | usage | nr | management | proof | input | connection | korea | price | incident |
| | project | interesting | information | mind | miss | access | human | equipment | korea |
| | system | world | nikkei | brain | food | security | people | model | jelly |

Experiments

(b) Delicious

| | | | | | | | | | |
|---|--|---|--|--|--|---|---|---|--|
| reference web imported design internet online cool toread tools blog | money finance economics business economy Finance financial investing bailout finances | video music videos fun entertainment funny movies media Video film | opensource software programming development linux tools rails ruby webdev rubyonrails | food recipes recipe cooking Food Recipes baking health vegetarian diy | windows linux sysadmin Windows security computer microsoft network Linux ubuntu | art photo photography photos Photography Art inspiration music foto fotografia | shopping shop Shopping home wishlist buy store fashion gifts house | iphone mobile hardware games iphone apple tech gaming mac game | education learning books book language library school teaching Education research |
| | money financial credit market economic october economy banks government bank | music video link tv movie itunes film amazon play interview | project code server ruby rails source file version files development | recipe food recipes make wine made add love eat good | windows system microsoft linux software file server user files ubuntu | art photography photos camera vol digital images 2008 photo tracks | buy online price cheap product order free products rating card | iphone apple ipod mobile game games pc phone mac touch | book legal theory books law university students learning education language |

(c) Flickr

| | | | | | | | |
|---|---|--|---|---|--|--|---|
| 2008 nikon canon white yellow red photo italy california color | dance bar dc digital concert bands music washingtondc dancing work | sea sunset sky clouds mountains ocean panorama south ireland oregon | autumn trees tree mountain fall garden bortescristian geotagged mud natura | rock house party park inn coach creature halloween mallory night | beach travel vacation camping landscape texas lake cameraphone md sun | family portrait cute baby boy kids brown closeup 08 galveston | island asia landscape rock blue tour plant tourguidesoma koh samui |
|---|---|--|---|---|--|--|---|