

1. Gibbs Sampling

$$\begin{aligned}
 P^*(z) T(z, z') &= P^*(z_i | u_i) P^*(z'_{\neq i} | u_i) \quad (P^*(z_i | u_i) \text{ is the} \\
 &= P^*(z_i | u_i) P^*(u_i) P^*(z'_i | u_i) \text{ probability that random} \\
 &= P^*(z_i | u_i) P^*(z'_i, u_i) \quad \text{variable } \cancel{z_j} \text{ has value } z'_i \text{ given that} \\
 &= P^*(z'_i) T(z'_i, z) \quad z_j = z_j \text{ for all } j \text{ not equal} \\
 &\quad \text{to } i)
 \end{aligned}$$

2. Let's say that there are two documents given below focusing on pets.

d₁: The Welsh corgi is a small type of herding dog that originated in Wales.

d₂: A husky and corgi make your heart melt with what we like to call the "Horge".

query term: husky corgi

mixed with $\lambda = 1/2$

$$P(q|d_1) = \frac{1}{2} \left(\frac{0}{14} + \frac{1}{30} \right) \times \frac{1}{2} \left(\frac{1}{14} + \frac{2}{30} \right) = 0.00115$$

$$P(q|d_2) = \frac{1}{2} \left(\cancel{\frac{1}{16}} + \frac{1}{30} \right) \times \frac{1}{2} \left(\frac{1}{16} + \frac{2}{30} \right) = 0.00309$$

From the result, we can get that d₂ is better than d₁ given the query. But d₁ is still better than other document focused on computer science (say one document d₃ focused on Java. $P(q|d_3)=0$).

If we don't use mixture model then $P(q|d_1)=0$ because "husky" doesn't appear in d₁, however, d₁ is certainly related to the query.

This mixture can smooth probabilities in document language model.

- to discount non-zero probabilities and to give some probability mass to unseen words.

3. PageRank and HITS

First, I'll discuss the pros and cons about PageRank and HITS and then I'll compare these two algorithms.

PageRank - Advantages

- ① PageRank is global measure and query independent
- ② robust (somewhat) against spam because it's hard to add links to a spam page from other important pages

PageRank - Disadvantages

- ① PageRank favors older pages because a new page will not have many links
- ② Google PageRank algorithm is patented.

HITS - Advantages

- ① has the ability to rank pages according to the query topic.
- ② can find multi-typical results
- ③ has no computational burden on small graphs

HITS - Disadvantages

- ① query dependent, so the result must be computed on the fly.
- ② the real Web is not hubs and authorities
- ③ very easy to spam (by out links)

Difference between PageRank and HITS

- ① PageRank is query independent
HITS is query dependent
- ② HITS is more prone to spams. ~~one page~~ can have lots of out links to authorities to pretend to be a good hub
- ③ PageRank operates on the whole Web graph
HITS operates on a small subgraph (the seed) from the Web graph
- ④ the result of HITS is less relevant than PageRank
(from "Comparative Analysis of PageRank and HITS Algorithms")

CSE 6240

HW4

Ke Wang
9030 58889

4. pLSI

a) Given $P(d)$, $P(z|d)$, $P(w|z)$, we can generate a word from a document.

First, a document is given by $P(d)$

Then, we generate a topic from the document by $P(z|d)$

At last, we generate a word from the topic by $P(w|z)$

b) EM algorithm is used in PLSA/PLSI to learn the (latent) parameters because explicitly finding the parameters is hard

Let the likelihood be $\ell(\theta)$ where θ is parameter. Maximizing $\ell(\theta)$ explicitly might be difficult, so the strategy is to repeatedly construct a lower-bound on ℓ (E-step), and then optimizes that lower bound (M-step)

Specifically, in the Expectation step, calculate the expected value of log likelihood function; in the Maximization step, find the parameter that maximizes the quantity.

c) pLSI can provide $P(w|z)$ and $P(z|d)$

from the training data and EM algorithm, we can finally get the parameters to calculate $P(w,d)$

For two documents d_1 and d_2 , we can calculate $P(d_1|z)$ and $P(d_2|z)$ (these can be calculated by the results of EM algorithm)

Then we use $\sum_z P(d_1|z) \cdot P(d_2|z)$ to calculate the distance

between the two documents.

5. LDA

a) w_i : words word-level Z_j : topics corpus-level ϕ_d : documents da-level

α : parameters ~~over~~ the document distribution over topics
 to assign probabilities to new documents (somewhat doc-level)

β : parameters over the topic distribution of words

to assign probabilities to new words (somewhat word-level)

b) E-step: for each document, running the following iterative algorithm
 to find the optimizing values of the variational parameters:

1. initialize ϕ and β

2. while not converge

3. for $n = 1 \rightarrow N_d$ 4. update ϕ 5. ~~update~~5. normalize ϕ 6. update β

M-step: maximize the resulting lower bound on the log likelihood
 with respect to the model parameters α and β ?