

hjshin@scai.snu.ac.kr
{btzhang,ytkim}@cse.snu.ac.kr

Learning Probabilistic Graph Models for Extracting Topic Words in a Collection of Text Documents

Hyung-Joo Shin⁰ Byoung-Tak Zhang Yung Tak Kim
School of Computer Science and Engineering, Seoul National University

(dyadic) likelihood Hoc 가 EM(Expected Maximization) 가 TREC-8 Ad 가 Topic 가

1. 가

[4][5], (clustering) 가 [7], [6] mixture sparse generative model [3][8]. Aspect Model EM [3] Aspect Model (indexing) (learning topic) latent variable mixture sparse generative model [3][8]. Aspect Model EM [3] Aspect Model (statistical mixture models), EM(expectation Maximization) (fitting) (layer) latent variable, (generate) latent variable

2. Aspect model class variable $z_k \in \mathcal{Z} = \{z_1, \dots, z_k\}$ $w_m \in \mathcal{W} = \{w_1, \dots, w_M\}$ $d_n \in \mathcal{D} = \{d_1, \dots, d_N\}$ 가 [3][5]. $(d_n, w_m), n = 1, \dots, N, m = 1, \dots, M$ generative model [3].

$$P(d_n, w_m) = P(d_n)P(w_m | d_n), \tag{1}$$

$$P(w_m | d_n) = \sum_{k=1}^K P(w_k | d_n)P(z_k | d_n) \tag{2}$$

(1), (2) Bayes' rule

$$P(d_n, w_m) = \sum_{k=1}^K P(z_k)P(w_m | z_k)P(d_n | z_k) \tag{3}$$

가 (d_n, w_m) 가 iid (independently distributed) w latent variable d log-likelihood [1][2][5]. EM

$$L = \sum_{n=1}^N \sum_{m=1}^M n(d_n, w_m) \log P(d_n, w_m) \tag{4}$$

(4) E-step (local maximum) EM

$$P(z_k | d_n, w_m) = \frac{P(z_k)P(d_n | z_k)P(w_m | z_k)}{\sum_{k=1}^K P(z_k)P(d_n | z_k)P(w_m | z_k)} \quad (5)$$

M-step

$$P(w_m | z_k) = \frac{\sum_{n=1}^N n(d_n, w_m)P(z_k | d_n, w_m)}{\sum_{m=1}^M \sum_{n=1}^N n(d_n, w_m)P(z_k | d_n, w_m)} \quad (6)$$

$$P(d_n | z_k) = \frac{\sum_{m=1}^M n(d_n, w_m)P(z_k | d_n, w_m)}{\sum_{m=1}^M \sum_{n=1}^N n(d_n, w_m)P(z_k | d_n, w_m)} \quad (7)$$

$$P(z_k) = \frac{1}{R} \sum_{m=1}^M \sum_{n=1}^N n(d_n, w_m)P(z_k | d_n, w_m), \quad (8)$$

$$R \equiv \sum_{m=1}^M \sum_{n=1}^N n(d_n, w_m)$$

Aspect model

1. $P(w_m | z_k)$ latent variable z_k 가 $P(d_n | z_k)$ d_n 가 z_k 가 $P(w_m | z_k)$ 가 $K=L$, latent variable $z_k (k=1, \dots, K)$ $P(w_m | z_k)$ 가 w topic $c_l (l=1, \dots, L)$ $I(d_n \in c_l) = 1$, c_l , c_l

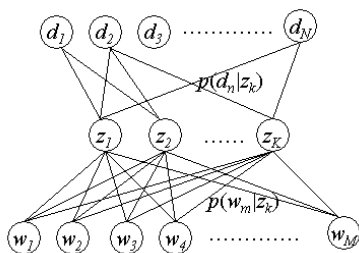


그림 1. Aspect Model

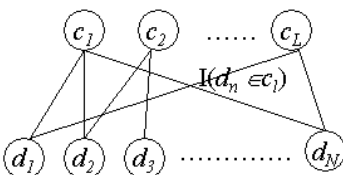


그림 2. 문서의 주어진 주제별 분류

topic relevant 가 ()
 relevant 가 434(347), topic
 401(300) relevant FBIS, FT, LATIMES
 . (가 6KB ,
 16KB 169 .) Porter
 stemming 524 stopwords 7826
 . 7826 frequency 5
 1 7000
 5610 2216 2216 . 169
 2216 .
 TREC-8 adhoc topic 434, 401 3 .

Topic 434
 Title: *Estonia, economy*
 Description: What is the state of the economy of Estonia?
 Relevant Document: **Documents that give concrete economic information such as economic statistics, entering economic unions and treaties, or monetary performance are relevant, as are discussions of economic issues such as transportation or pollution.**

Topic 401
 Title: *foreign minorities, Germany*
 Description: What language and cultural differences impede the integration of foreign minorities in Germany?
 Relevant Document: A relevant document will focus on the causes of the lack of integration in a significant way; that is, the mere mention of immigration difficulties is not relevant. Documents that discuss immigration problems unrelated to Germany are also not relevant.

그림 3. Topic 434, 401의 Definition

(밑줄은 stopwords를 제거하고 stemming한 후 뽑아 낸 단어임. 밑줄 친 단어 중 이탤릭 체로 된 것은 Topic의 특징을 정확히 표현하는 단어임)

Latent variable topic 2 ($K=L=2$) EM
 100 iteration (6), (7), (8)
 1 $k=0,1$ $P(w_m | z_k)$ 가 10 $P(w_m | z_k)$
 (3 topic)
 $P(w_m | z_k)$ 3 topic
 $P(w_m | z_k)$
 1 3 가 Topic Definition
 2, 3 Topic
 Definition Topic
 4 (a) $k=0,1$ $P(d_n | z_k)$ 가 80 , 89
 Topic 401, Topic434
 clustering , N
 K' clustering latent variable $K=K'$
 $k' (k'=1, \dots, K')$
 $P(d_n | z_{k'})$ cluster k' 4
 (b) Topic401 80 $P(d_n | z_1) > P(d_n | z_0)$
 Topic434 89 $P(d_n | z_0) >$
 $P(d_n | z_1)$ classification
 5 topic relevant ,
 Topic relevant $P(d_n | z_k)$,
 Topic relevant relevant

3.

TREC-8 ¹⁾ adhoc task
 DTDS, FR94, FT, FBIS, LATIMES
 50 topics 401-450가

¹⁾ Text Retrieval Conference 8, <http://trec.nist.gov>
²⁾ McCallum Bowlibrary
<http://www.cs.cmu.edu/~mccallum/bow>

k=0	k=1
estonia	germani
percent	immigr
state	integ
russian	minor
estonian	union
bank	cultur
russia	foreign
baltic	asian
econom	language
invest	unity

표 1. k=0,1에 대해 $P(w_m|z_k)$ 가 큰 단어 10개를 $P(w_m|z_k)$ 의 내림차순으로 정렬한 것

	k=0	k=1
minor	0	0.0010
germani	0.0010	0.0170
language	0.0010	0.0010
cultur	0	0.0010
differ	0.0010	0.0010
integ	0.0010	0.0010
immigr	0	0.0050
estonia	0.0160	0
economi	0.0050	0
state	0.0120	0.0050
econom	0.0050	0
statist	0.0010	0
union	0.0030	0.0030
treati	0	0.0010
monetary	0.0020	0
transport	0.0010	0
pollution	0.0010	0

표 2. 각 topic에 중요하다고 여겨지는 단어에 대해 $P(w_n|z_k)$ (이탤릭체는 Topic401의 topic definition에 해당하는 단어이다.)

	Topic 434	Topic 401
최대	0.0160/0.0160	0.0050/0.0050
평균	0.0029/0.0004	0.0021/0.0003
최소	0/0	0/0

표 3. 각 topic에 주어진 단어의 집합에 대해, 모든 단어의 집합에 대해 $P(w_m|z_k)$ 의 최대, 평균, 최소값

	Topic434	Topic 401t
k=0	89 개	0 개
k=1	4 개	76 개

(a)

	$P(d_n z_0) > P(d_n z_1)$	$P(d_n z_1) > P(d_n z_0)$
Topic434	87 개	2 개
Topic 401	3 개	77 개

(b)

표 4. (a) k=0,1에 대해 $P(d_n|z_k)$ 가 큰 문서 상위 80개, 89개 중 Topic 401, Topic434에 속한 개수로 이 방법이 문서의 clustering에 응용될 수 있음을 보인다. (b) Topic401에 속한 문서 80개 중 $P(d_n|z_1) > P(d_n|z_0)$ 인 문서의 개수와 Topic434에 속한 문서 89개 중 $P(d_n|z_0) > P(d_n|z_1)$ 인 문서의 개수로 이 방법이 문서의 classification에 응용될 수 있음을 보인다.

	Topic434 relevant	Topic 434 not relevant	Topic 434에 속한 모든 문서
최대	0.0240	0.0190	0.0240
평균	0.0107	0.0031	0.0085
최소	0.0030	0.0000	0.0000

	Topic401 relevant	Topic 401 not relevant	Topic 401에 속한 모든 문서
최대	0.0260	0.0120	0.0260
평균	0.0085	0.0012	0.0059
최소	0.0000	0.0000	0.0000

표 5. 각 topic에 relevant한 문서의 집합에 대해, 그 외 모든 문서에 대해, 그리고 모든 문서에 대해 $P(d_n|z_k)$ 의 최대, 최소, 평균값

($P(w_m|z_k)$) ($P(d_n|z_k)$) 가
 Latent variable (classification) class,
 clustering prototype
 가
 latent variable , latent
 variable $P(w_m|z_k)$ 가 가 $P(d_n|z_k)$
 가
 clustering,
 Topic 2
 가
 (98-199)

[1] Dempster, A.P., N.M.Laird, and D.B. Rubin, "Maximum likelihood from incomplete data via the EM algorithm (with discussion)", *J. Roy. Stat. Soc.*, B39, 1-38, 1977.
 [2] Jeff A. Bilmes, "A Gentle Tutorial of EM Algorithm and its Application to Parameter Estimation for Gaussian Mixture and Hidden Markov Models", International Computer Science Institute, 1998.
 [3] Thomas Hofmann, Jan Puzicha, "Unsupervised Learning from Dyadic Data", in *Advances in Neural Information Processing System*, 1998
 [4] Deerwester, S., Dumais, S. T., Furnas, G. W., Landauer, T. K., and Harshman, R., "Indexing by latent semantic analysis", *Journal of the American Society for Information Science*, 1990.
 [5] Thomas Hoffmann, "Probabilistic Latent Semantic Indexing", in *Proc. of the 22nd International Conference on Research and Development in Information Retrieval (SIGIR '99)*, 1999.
 [6] Yiming, Yang., "Noise Reduction in a Statistical Approach to Text Categorization", in *Proceedings of the 18th Annual international ACM SIGIR Conference on Research and Development in Information Retrieval (SIGIR '95)*, pp. 256-263, 1995
 [7] Thomas Hoffmann, "Learning and Representing Topic. A Hierarchical Mixture Model for Word Occurrences in Document Databases", in *Conferences for Automated Learning and Discovery*, 1998.
 [8] Hinton, G. E. and Ghahramani, Z., "Generative Models for Discovering Sparse Distributed Representations. *Phil. Trans. Roy. Soc. London B*, 352:1177-1190, 1997.