

Lecture 10

Latent Dirichlet Allocation (LDA)

Instructor: Shibo Li

shiboli@cs.fsu.edu



Department of Computer Science
Florida State University

- Latent Dirichlet Allocation
- Variational inference

- Latent Dirichlet Allocation
- Variational inference

- A classical text mining model that extract topics from the text corpus
- Broadly used in all kinds of text mining and related tasks: information retrieval, text classification, advertisement keywords extraction, sentimental analysis,
- <https://medium.com/@fatmafatma/industrial-applications-of-topic-model-100e48a15ce4>
- A very good example to study Bayesian learning

“Arts”	“Budgets”	“Children”	“Education”
NEW	MILLION	CHILDREN	SCHOOL
FILM	TAX	WOMEN	STUDENTS
SHOW	PROGRAM	PEOPLE	SCHOOLS
MUSIC	BUDGET	CHILD	EDUCATION
MOVIE	BILLION	YEARS	TEACHERS
PLAY	FEDERAL	FAMILIES	HIGH
MUSICAL	YEAR	WORK	PUBLIC
BEST	SPENDING	PARENTS	TEACHER
ACTOR	NEW	SAYS	BENNETT
FIRST	STATE	FAMILY	MANIGAT
YORK	PLAN	WELFARE	NAMPHY
OPERA	MONEY	MEN	STATE
THEATER	PROGRAMS	PERCENT	PRESIDENT
ACTRESS	GOVERNMENT	CARE	ELEMENTARY
LOVE	CONGRESS	LIFE	HAITI

The William Randolph Hearst Foundation will give \$1.25 million to Lincoln Center, Metropolitan Opera Co., New York Philharmonic and Juilliard School. “Our board felt that we had a real opportunity to make a mark on the future of the performing arts with these grants an act every bit as important as our traditional areas of support in health, medical research, education and the social services,” Hearst Foundation President Randolph A. Hearst said Monday in announcing the grants. Lincoln Center’s share will be \$200,000 for its new building, which will house young artists and provide new public facilities. The Metropolitan Opera Co. and New York Philharmonic will receive \$400,000 each. The Juilliard School, where music and the performing arts are taught, will get \$250,000. The Hearst Foundation, a leading supporter of the Lincoln Center Consolidated Corporate Fund, will make its usual annual \$100,000 donation, too.

Latent Dirichlet Allocation

David M. Blei

*Computer Science Division
University of California
Berkeley, CA 94720, USA*

BLEI@CS.BERKELEY.EDU

Andrew Y. Ng

*Computer Science Department
Stanford University
Stanford, CA 94305, USA*

ANG@CS.STANFORD.EDU

Michael I. Jordan

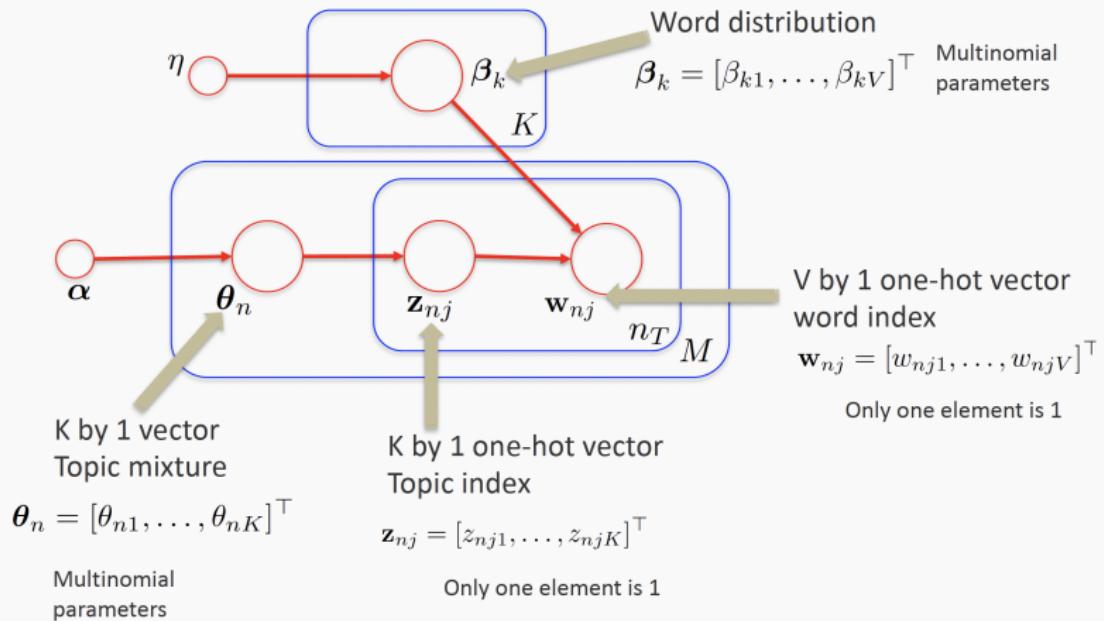
*Computer Science Division and Department of Statistics
University of California
Berkeley, CA 94720, USA*

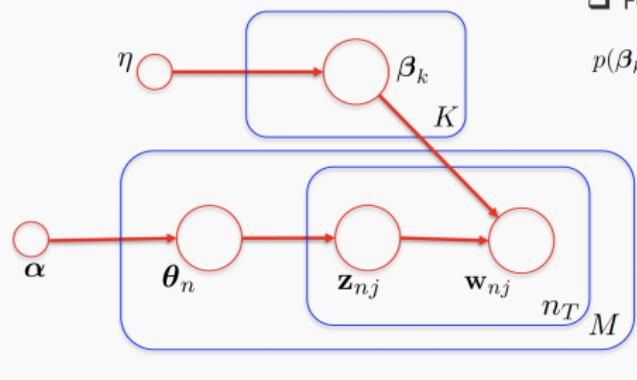
JORDAN@CS.BERKELEY.EDU

- ▶ <https://www.jmlr.org/papers/volume3/blei03a/blei03a.pdf>
- ▶ <https://arxiv.org/pdf/1601.00670>

- Given K topics
- First sample K topics (word distributions)
- For each document in the corpus
 - Sample topic mixture distribution
 - For each word in the document
 - Sample the topic index, according to which to sample the word

Suppose we have V words, M documents, document n has n_T words





□ For each topic

$$p(\beta_k|\eta) = \text{Dir}(\beta_k|\eta) = \frac{\Gamma(V\eta)}{\prod_{v=1}^V \Gamma(\eta)} \prod_{v=1}^V \beta_{kv}^{\eta-1}$$

□ For each document

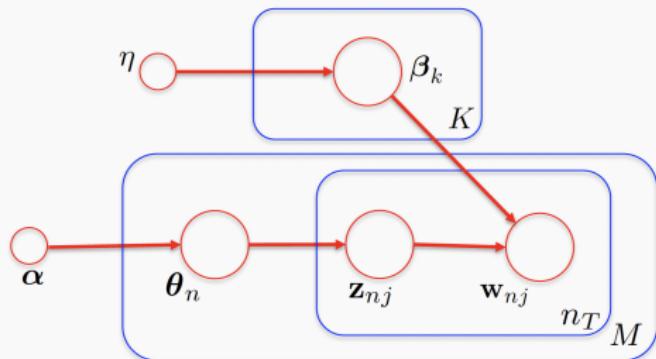
$$p(\theta_n|\alpha) = \text{Dir}(\theta_n|\alpha) = \frac{\Gamma(\sum_{k=1}^K \alpha_k)}{\prod_{k=1}^K \Gamma(\alpha_k)} \prod_{k=1}^K \theta_{nk}^{\alpha_k - 1}$$

□ For each word

$$p(z_{nj}|\theta_n) = \text{Mul}(z_{nj}|\theta_n) = \prod_{k=1}^K \theta_{nk}^{z_{nj}k}$$

$$p(w_{nj}|z_{nj}=1, \beta) = \text{Mul}(w_{nj}|\beta) = \prod_{v=1}^V \beta_{kv}^{w_{nj}v}$$

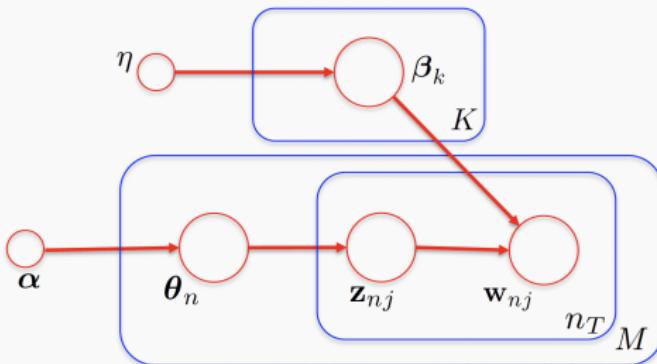
$$p(w_{nj}|z_{nj}, \beta) = \prod_{k=1}^K \prod_{v=1}^V \beta_{kv}^{z_{nj}k w_{nj}v}$$



$$p(\beta, \theta, Z, W | \eta, \alpha)$$

$$= \prod_{k=1}^K \text{Dir}(\beta_k | \eta) \prod_{n=1}^M \text{Dir}(\theta_n | \alpha) \left\{ \prod_{j=1}^{n_T} \text{Mul}(z_{nj} | \theta_n) \left[\prod_{k=1}^K \prod_{v=1}^V \beta_{kv}^{z_{njk}} w_{njv} \right] \right\}$$

- Latent Dirichlet Allocation
- Variational inference



How can we compute the posterior of

$\beta = \{\beta_1, \dots, \beta_K\}$ Topic words: critical for numerous tasks

$\theta = \{\theta_1, \dots, \theta_M\}$ Topic mixture: low-rank representation of docs

True posterior is intractable to compute

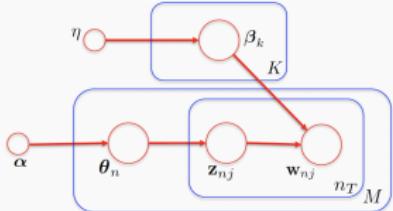
- We use variational EM algorithm (empirical Bayes)

- ❑ E step: mean-field update

$$q(\beta, \theta, \mathbf{Z}) = \prod_{k=1}^K q(\beta_k) \prod_{n=1}^M \left[q(\theta_n) \prod_{j=1}^{n_T} q(\mathbf{z}_{nj}) \right]$$

- ❑ M step

Maximize variational lower bound of the model evidence w.r.t α, η



$$p(\beta, \theta, Z, W | \eta, \alpha)$$

$$= \prod_{k=1}^K \text{Dir}(\beta_k | \eta) \prod_{n=1}^M \text{Dir}(\theta_n | \alpha) \left\{ \prod_{j=1}^{n_T} \text{Mul}(z_{nj} | \theta_n) \left[\prod_{k=1}^K \prod_{v=1}^V \beta_{kv}^{z_{njk}} w_{njv} \right] \right\}$$

Update each $q(z_{nj})$

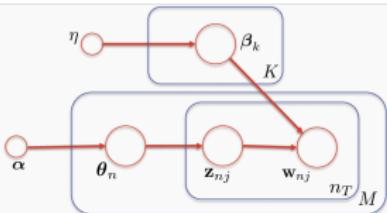
$$q(z_{nj}) \propto \exp \left(\underbrace{\mathbb{E}_q \log \left[\text{Mul}(z_{nj} | \theta_n) \prod_{k=1}^K \prod_{v=1}^V \beta_{kv}^{z_{njk}} w_{njv} \right]}_{\exp \left(\sum_{k=1}^K z_{njk} \left[\mathbb{E}_q[\log \theta_{nk}] + \sum_{v=1}^V w_{njv} \mathbb{E}_q[\log \beta_{kv}] \right] \right)} \right)$$

$$q(z_{nj}) = \text{Mul}(z_{nj} | \phi_{nj})$$

$$\phi_{njk} \propto \exp \left(\mathbb{E}_q \log[\theta_{nk}] + \sum_{v=1}^V w_{njv} \mathbb{E}_q \log \beta_{kv} \right)$$

$$p(\beta, \theta, Z, W | \eta, \alpha)$$

$$= \prod_{k=1}^K \text{Dir}(\beta_k | \eta) \prod_{n=1}^M \text{Dir}(\theta_n | \alpha) \left\{ \prod_{j=1}^{n_T} \text{Mul}(z_{nj} | \theta_n) \left[\prod_{k=1}^K \prod_{v=1}^V \beta_{kv}^{z_{njk}} w_{nv} \right] \right\}$$



Update each $q(\theta_n)$

$$q(\theta_n) \propto \exp \left(\mathbb{E}_q \left[\log[\text{Dir}(\theta_n | \alpha)] + \sum_{j=1}^{n_T} \log[\text{Mul}(z_{nj} | \theta_n)] \right] \right)$$

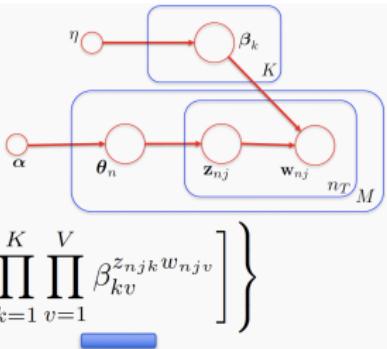
$$\sum_{k=1}^K (\alpha_k + \sum_{j=1}^{n_T} \mathbb{E}_q[z_{njk}] - 1) \log \theta_{nk}$$

$$q(\theta_n) = \text{Dir}(\theta_n | \gamma_n)$$

$$\gamma_{nk} = \alpha_k + \sum_{j=1}^{n_T} \mathbb{E}_q[z_{njk}]$$

$$p(\beta, \theta, Z, W | \eta, \alpha)$$

$$= \prod_{k=1}^K \text{Dir}(\beta_k | \eta) \prod_{n=1}^M \text{Dir}(\theta_n | \alpha) \left\{ \prod_{j=1}^{n_T} \text{Mul}(z_{nj} | \theta_n) \left[\prod_{k=1}^K \prod_{v=1}^V \beta_{kv}^{z_{njk}} w_{njv} \right] \right\}$$



Update each $q(\beta_k)$

$$q(\beta_k) \propto \exp \left(\mathbf{E}_q \left[\underbrace{\log[\text{Dir}(\beta_k | \eta) + \sum_{n=1}^M \sum_{j=1}^{n_T} \sum_{v=1}^V z_{njk} w_{njv} \beta_{kv}]}_{\sum_{v=1}^V \log[\beta_{kv}] (\eta + \sum_{n=1}^M \sum_{j=1}^{n_T} \mathbb{E}_q[z_{njk}] w_{njv} - 1)} \right] \right)$$

$$q(\beta_k) = \text{Dir}(\beta_k | \psi_k)$$

$$\psi_{kv} = \eta + \sum_{n=1}^M \sum_{j=1}^{n_T} \mathbb{E}_q[z_{njk}] w_{njv}$$

- How to compute the required moments in the update?

$$q(\mathbf{z}_{nj}) = \text{Mul}(\mathbf{z}_{nj} | \boldsymbol{\phi}_{nj})$$

$$\mathbb{E}_q[z_{njk}]$$

$$q(\boldsymbol{\theta}_n) = \text{Dir}(\boldsymbol{\theta}_n | \boldsymbol{\gamma}_n)$$

$$\mathbb{E}_q \log[\theta_{nk}]$$

$$q(\boldsymbol{\beta}_k) = \text{Dir}(\boldsymbol{\beta}_k | \boldsymbol{\psi}_k)$$

$$\mathbb{E}_q \log \beta_{kv}$$

Leave it as your review!

$$p(\beta, \theta, \mathbf{Z}, \mathbf{W} | \eta, \alpha)$$

$$= \prod_{k=1}^K \text{Dir}(\beta_k | \eta) \prod_{n=1}^M \text{Dir}(\theta_n | \alpha) \left\{ \prod_{j=1}^{n_T} \text{Mul}(\mathbf{z}_{nj} | \theta_n) \left[\prod_{k=1}^K \prod_{v=1}^V \beta_{kv}^{z_{njk}} w_{njv} \right] \right\}$$


 $\frac{\Gamma(\sum_{k=1}^K \alpha_k)}{\prod_k \Gamma(\alpha_k)} \prod_{k=1}^K \theta_{nk}^{\alpha_k - 1}$

Variational lower bound

$$\begin{aligned} \mathcal{L}(\alpha) &= \sum_{n=1}^M \log \Gamma\left(\sum_{k=1}^K \alpha_k\right) - \sum_{n=1}^M \sum_{k=1}^K \log \Gamma(\alpha_k) \\ &\quad + \sum_{n=1}^M \sum_{k=1}^K (\alpha_k - 1) \mathbb{E}_q[\log \theta_{nk}] + \text{const} \end{aligned}$$

Variational lower bound

$$\begin{aligned}\mathcal{L}(\boldsymbol{\alpha}) &= \sum_{n=1}^M \log \Gamma\left(\sum_{k=1}^K \alpha_k\right) - \sum_{n=1}^M \sum_{k=1}^K \log \Gamma(\alpha_k) \\ &\quad + \sum_{n=1}^M \sum_{k=1}^K (\alpha_k - 1) \mathbb{E}_q[\log \theta_{nk}] + \text{const}\end{aligned}$$



$$\frac{\partial \mathcal{L}}{\partial \alpha_k} = M \Psi\left(\sum_{k=1}^K \alpha_k\right) - M \Psi(\alpha_k) + \sum_{n=1}^M \mathbb{E}_q[\log \theta_{nk}]$$

Digamma function

- Derive $\frac{\partial \mathcal{L}}{\partial \eta}$ Note that η is a scalar

Leave it as your exercise

Use any gradient based algorithm with constraints

$\alpha > 0, \eta > 0$ e.g., LBFGS-B

- Write down LDA sampling procedure and joint probability
- Derive the variational E updates and gradients for M step
- Implement an algorithm for LDA inference and test it on real-world data (see homework assignments).