

# Reasoning (III)

Mingsheng Long

Tsinghua University

# Outline

- Probabilistic Topic Model
  - Dirichlet-Multinomial Model
  - Latent Dirichlet Allocation
  - Document Generation Process
- Variational Inference
  - Mean Field Variational Inference
  - Variational Inference for LDA
  - Parameter Estimation for LDA



# Content Analysis

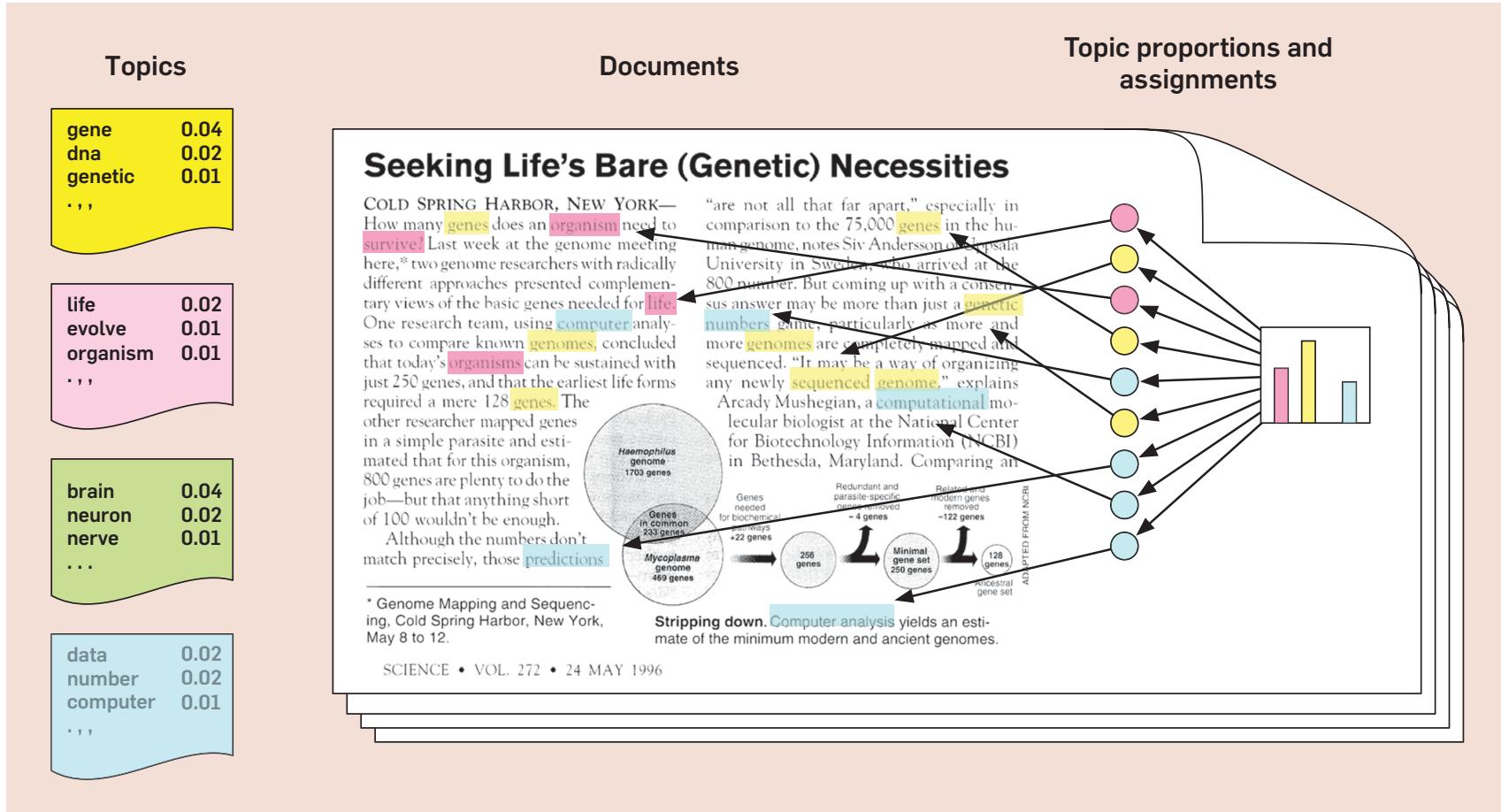
- Suppose you're given a massive corpora and asked to carry out the following tasks:
  - Organize the documents into thematic categories
  - Enable a domain expert to analyze and understand the content
  - Find relationships between the categories
  - Understand how authorship influences the content...



- Key problem: How documents are generated? Which distribution?

# Probabilistic Topic Model

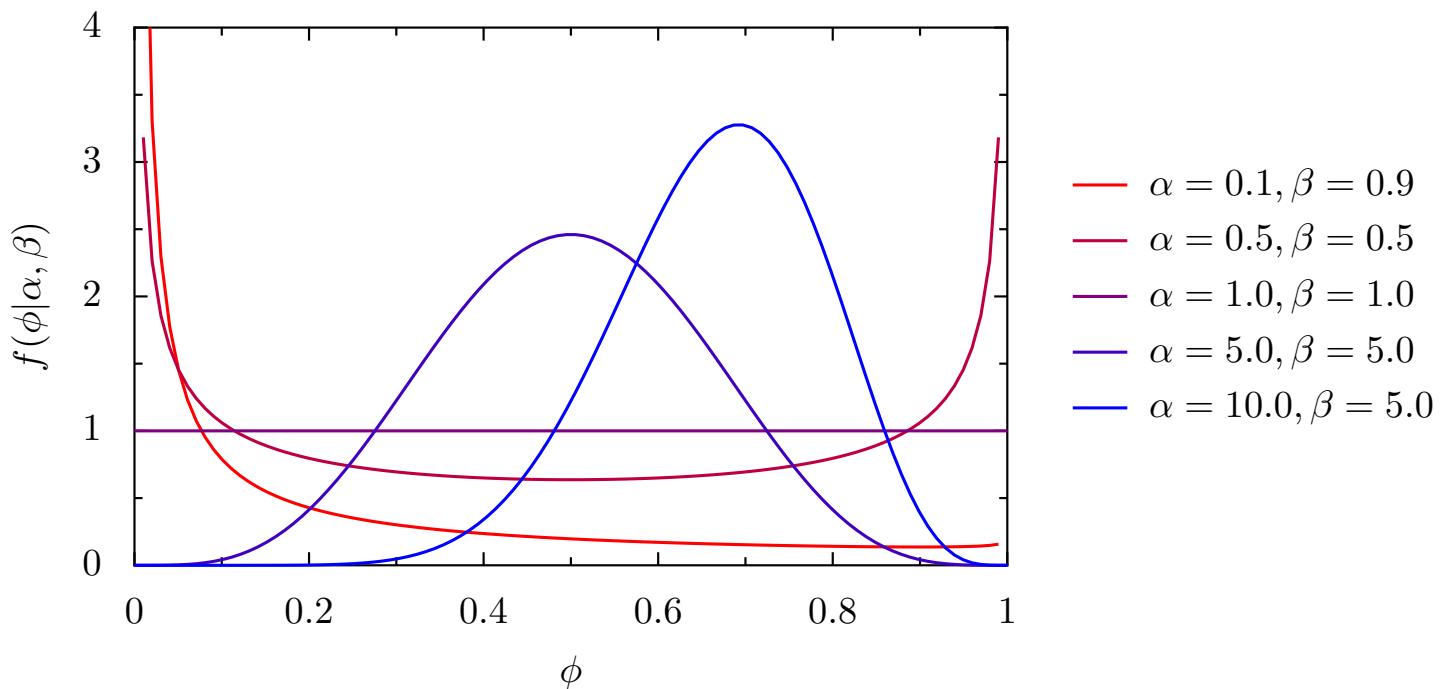
- The answer to content analysis is probabilistic topic model (PTM).



# Beta-Bernoulli Model

- Beta Distribution: conjugate prior for Bernoulli distribution

$$f(\phi|\alpha, \beta) = \frac{1}{B(\alpha, \beta)} \phi^{\alpha-1} (1-\phi)^{\beta-1}$$



# Beta-Bernoulli Model

- How are documents generated?
- Generative Process

$$\phi \sim \text{Beta}(\alpha, \beta)$$

*[draw distribution over words]*

For each word  $n \in \{1, \dots, N\}$

$$x_n \sim \text{Bernoulli}(\phi)$$

*[draw word]*

- Example corpus (heads/tails)

H	T	T	H	H	T	T	H	H	H
$x_1$	$x_2$	$x_3$	$x_4$	$x_5$	$x_6$	$x_7$	$x_8$	$x_9$	$x_{10}$

- There are only two words? We need more! Bernoulli is not enough.



# Multinomial Distribution

- Multinomial Distribution models the probability of counts of each side for rolling a  $k$ -sided dice  $n$  times.

- Each side equips with a probability  $\theta_i$ .

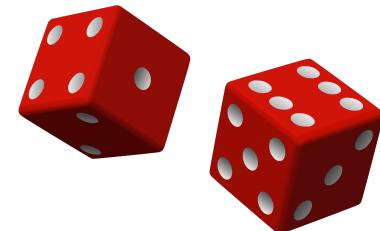
- $\sum_{i=1}^k \theta_i = 1$ .

- Write by  $\text{Mult}(n, \boldsymbol{\theta})$

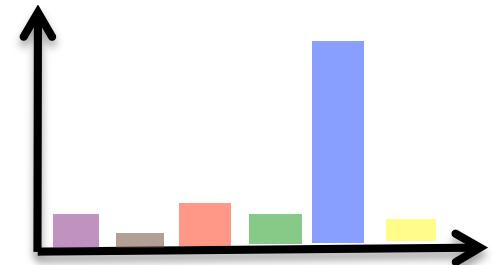
- $p(n_1, n_2, \dots, n_k) = \frac{n!}{n_1! \dots n_k!} \theta_1^{n_1} \theta_2^{n_2} \dots \theta_k^{n_k}$

- $\sum_{i=1}^k n_i = n$

- When  $n = 1$ , it is Categorical Distribution,  
 $\text{Cat}(\boldsymbol{\theta}) = \text{Mult}(1, \boldsymbol{\theta})$



Flipping dice:  
 $\sum_{l=1}^6 \theta_l = 1$

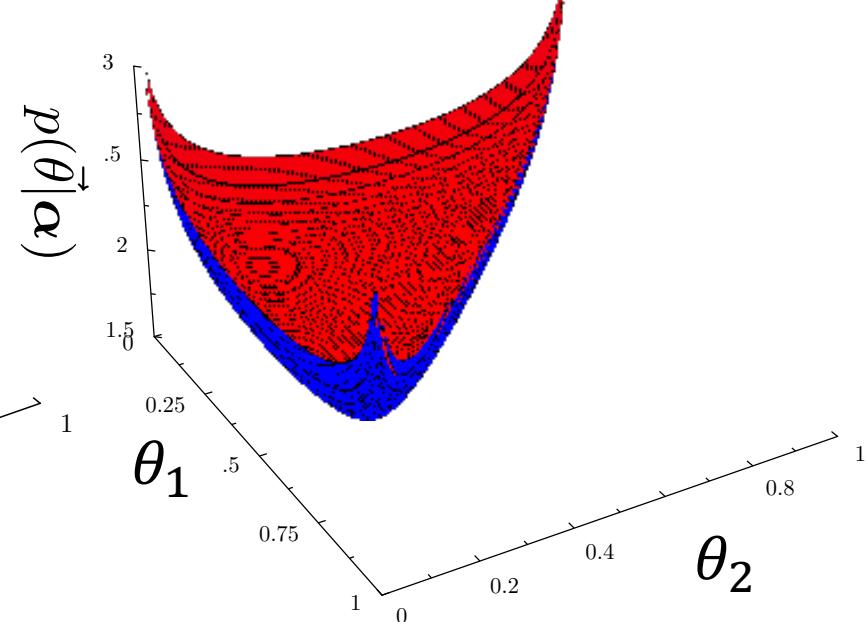
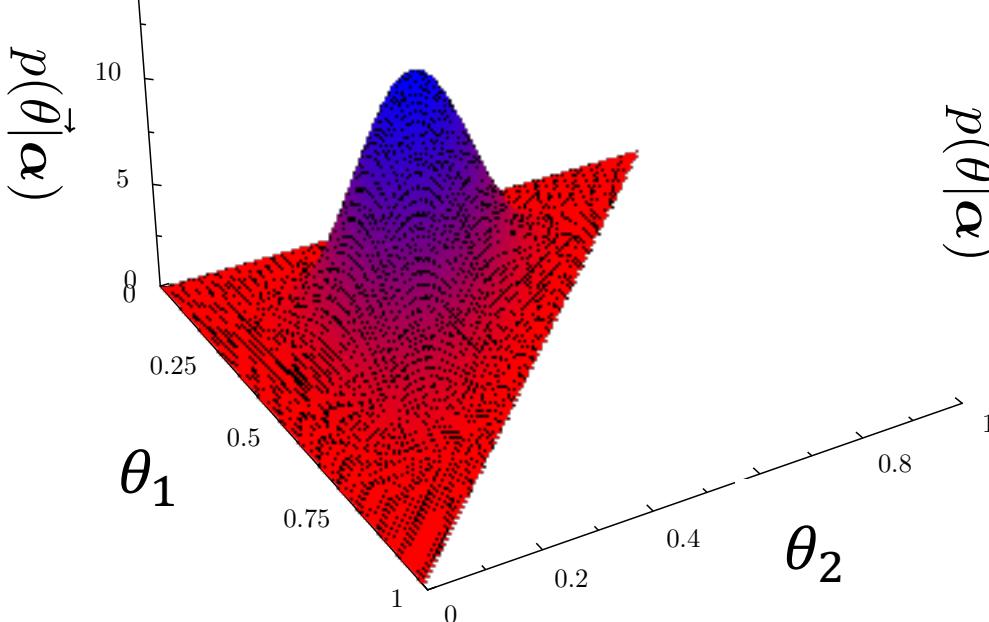


# Dirichlet-Multinomial Model

- Dirichlet Distribution: Multi-dimensional version of Beta Distribution

$$p(\vec{\theta}|\alpha) = \frac{1}{B(\alpha)} \prod_{k=1}^K \theta_k^{\alpha_k - 1}$$
 Where  $B(\alpha) = \frac{\prod_{k=1}^K \Gamma(\alpha_k)}{\Gamma(\sum_{k=1}^K \alpha_k)}$

$\sum_{i=1}^K \theta_i = 1!$  Conjugate prior for Multinomial Distribution.



# Dirichlet-Multinomial Model

- Assume that we have a vocabulary of  $V$  words.

- Generative Process

What is the dimension of  $\theta$  here?

$$\theta \sim \text{Dir}(\alpha)$$

[draw distribution over words]

$$\text{For each word } n \in \{1, \dots, N\}$$

[draw word]

$$w_n \sim \text{Mult}(1, \theta)$$

[draw word]

Categorical Distribution

- Example corpus

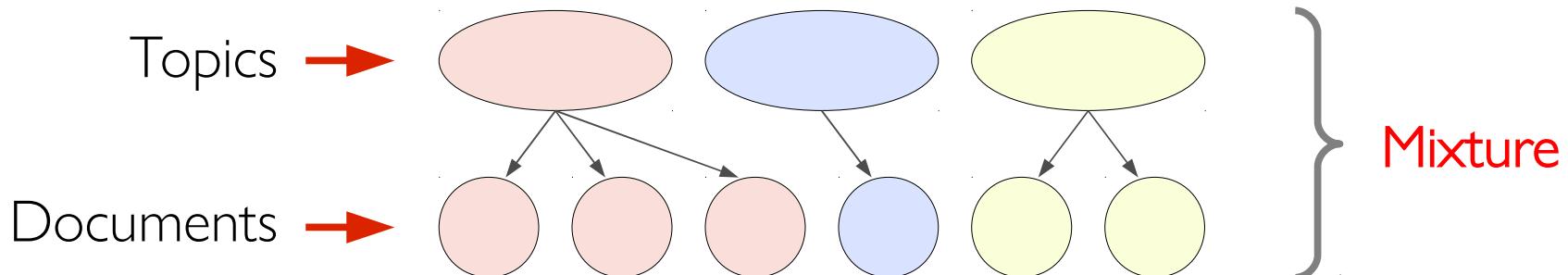
the	he	is	the	and	the	she	she	is
$w_1$	$w_2$	$w_3$	$w_4$	$w_5$	$w_6$	$w_7$	$w_8$	$w_9$

- You will see the benefit of using conjugate prior for a distribution.

# Dirichlet-Multinomial Mixture Model

- Generative Process

We have  $K$  topics.



- Example corpus

We have  $V$  words.

the	he	is
$w_{11}$	$w_{12}$	$w_{13}$

Document 1

the	and	the
$w_{21}$	$w_{22}$	$w_{23}$

Document 2

the	she	she	is
$w_{31}$	$w_{32}$	$w_{33}$	$w_{34}$

Document 3

# Dirichlet-Multinomial Mixture Model

What is the dimension of  $\boldsymbol{\theta}$ ?

- Generative Process

What is the dimension of  $\boldsymbol{\beta}_k$ ?

For each topic  $k \in \{1, \dots, K\}$  :

$$\boldsymbol{\beta}_k \sim \text{Dir}(\boldsymbol{\eta})$$

[draw distribution over words]

$$\boldsymbol{\theta} \sim \text{Dir}(\boldsymbol{\alpha})$$

[draw distribution over topics]

For each document  $d \in \{1, \dots, D\}$

$$z_d \sim \text{Mult}(1, \boldsymbol{\theta})$$

[draw topic assignment]

For each word  $n \in \{1, \dots, N_d\}$

$$w_{dn} \sim \text{Mult}(1, \boldsymbol{\beta}_{z_d})$$

[draw word]

- Example corpus

the	he	is
$w_{11}$	$w_{12}$	$w_{13}$

Document 1

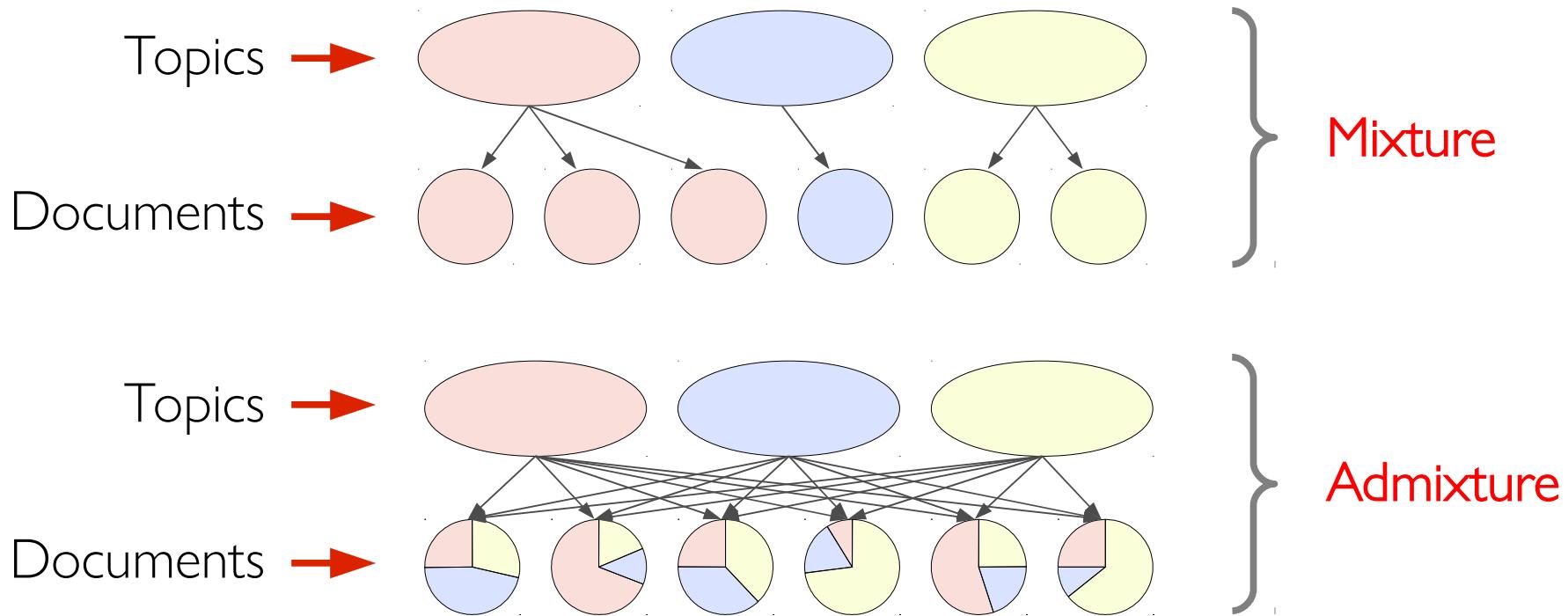
the	and	the
$w_{21}$	$w_{22}$	$w_{23}$

Document 2

the	she	she	is
$w_{31}$	$w_{32}$	$w_{33}$	$w_{34}$

Document 3

# Mixture vs. Admixture



There should be **multiple topics** in a document:  
Each document can be modeled by **a new distribution of topics!**

# Outline

- Probabilistic Topic Model
  - Dirichlet-Multinomial Model
  - Latent Dirichlet Allocation
  - Document Generation Process
- Variational Inference
  - Mean Field Variational Inference
  - Variational Inference for LDA
  - Parameter Estimation for LDA



# Latent Dirichlet Allocation (LDA)

Dirichlet-  
Multinomial  
Mixture  
Model

Latent  
Dirichlet  
Allocation

Each document has a single topic  $z$

For each topic  $k \in \{1, \dots, K\}$  :

$$\beta_k \sim \text{Dir}(\eta)$$

$$\theta \sim \text{Dir}(\alpha)$$

For each document  $d \in \{1, \dots, D\}$

$$z_d \sim \text{Mult}(1, \theta)$$

[draw distribution over words]

[draw distribution over topics]

[draw topic assignment]

For each word  $n \in \{1, \dots, N_d\}$

$$w_{dn} \sim \text{Mult}(1, \beta_{z_d})$$

[draw word]



Each document has a topic distribution  $\theta$

Each word has a topic  $z$

For each topic  $k \in \{1, \dots, K\}$  :

$$\beta_k \sim \text{Dir}(\eta)$$

[draw distribution over words]

For each document  $d \in \{1, \dots, D\}$

$$\theta_d \sim \text{Dir}(\alpha)$$

[draw distribution over topics]

For each word  $n \in \{1, \dots, N_d\}$

$$z_{dn} \sim \text{Mult}(1, \theta_d)$$

[draw topic assignment]

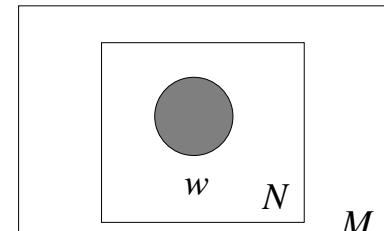
$$w_{dn} \sim \text{Mult}(1, \beta_{z_{dn}})$$

[draw word]

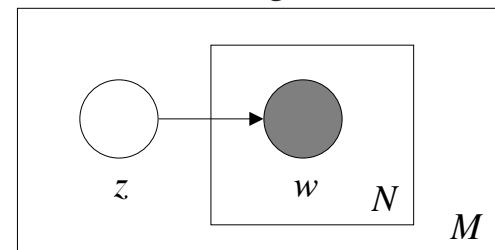


# Probabilistic Graphical Model

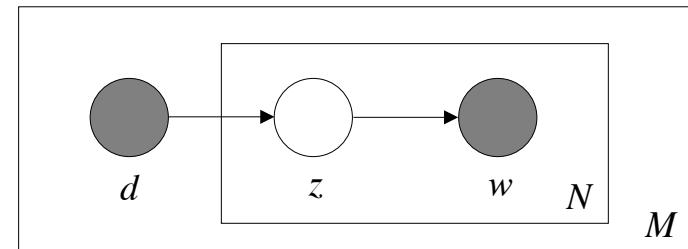
- Box with number on the bottom right:
  - Hierarchical samples and Number of samples.
- White/Gray circle: Hidden/Observed Variables.
- Arrow: Generating path (conditional prob).
- Read (b):
  - We generate  $M$  samples of  $z$ .
  - Each contains  $N$  kid samples  $w$ .
  - E.g.  $M$  clusters, each with  $N$  samples in it.
  - $w$  is observed, but  $z$  is hidden.
- Can you read (c) Probabilistic Latent Semantic Indexing (PLSI) now?



(a) unigram

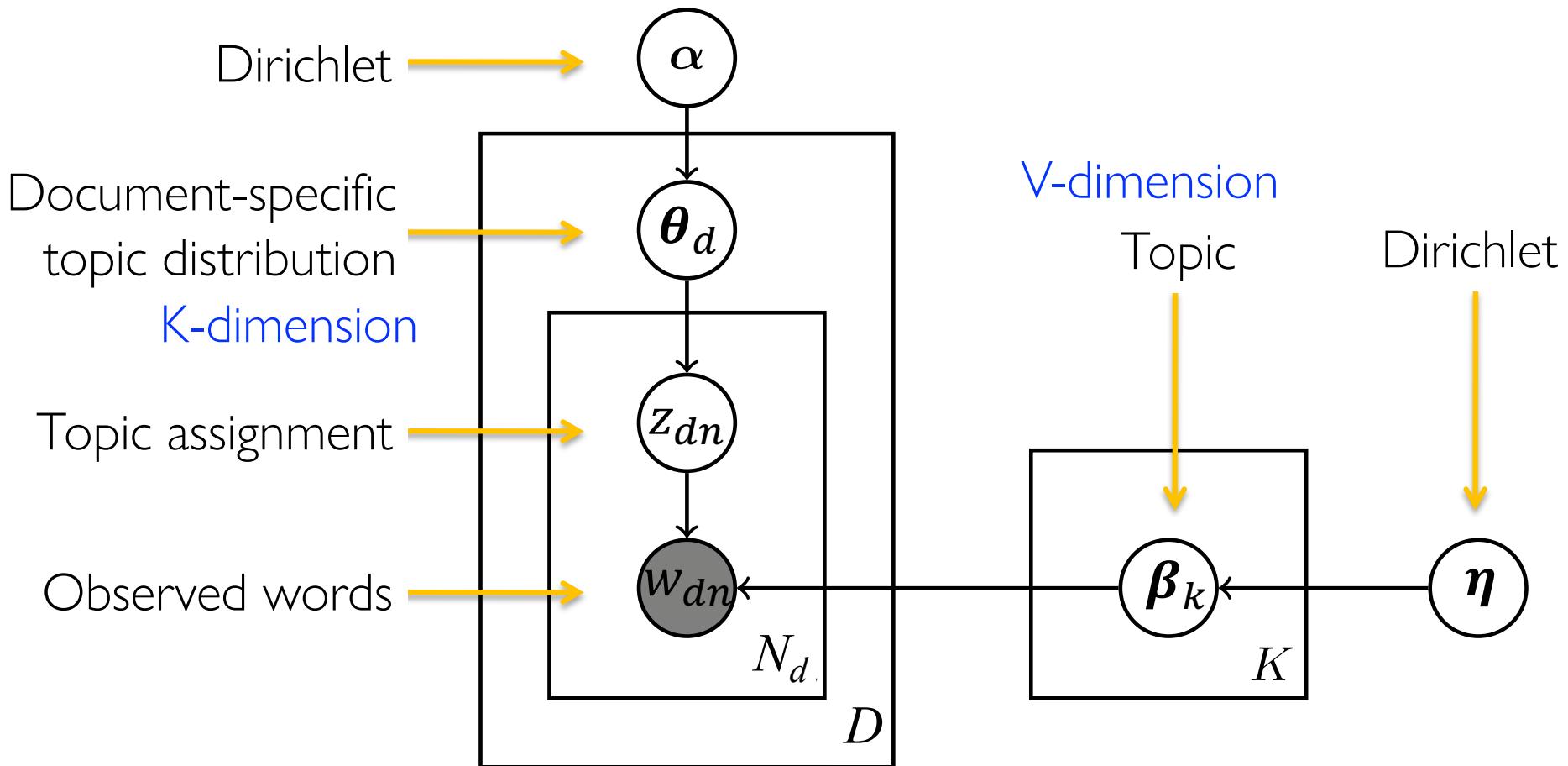


(b) mixture of unigrams



(c) pLSI/aspect model

# Latent Dirichlet Allocation (LDA)



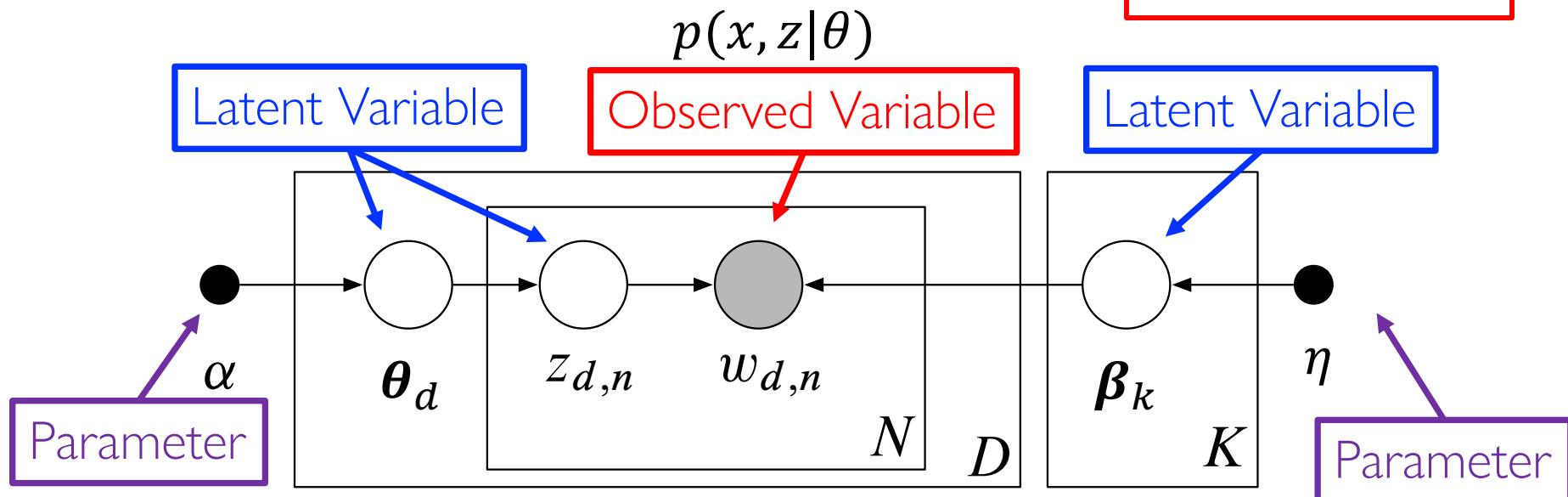
Can you explain graph of LDA now?

$$p(\beta, \theta, z, w | \alpha, \eta) = p(\beta | \eta) p(\theta | \alpha) p(z | \theta) p(w | z, \beta)$$

# LDA as Latent Variable Model

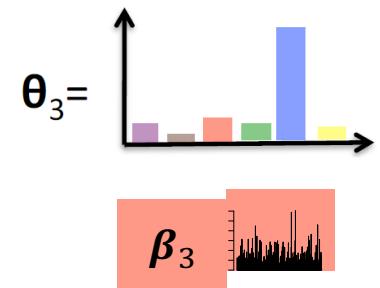
- Two sets of random variables:  $z$  and  $x$ 
  - $z$  consists of unobserved hidden variables.
  - $x$  consists of observed variables.
- Joint probability model parameterized by  $\theta \in \Theta$ :

Where is  $x, z, \theta$  in LDA?



# How to Use LDA

- After we learn the model, we can do inference of all latent variable.
  - Learning: estimate parameters  $\alpha, \eta$
  - Inference: for latent variable  $\theta_{1:D}, \beta_{1:K}, z_{1:D, 1:N_d}$
- From latent variable  $\theta_d$ :
  - We can know the main topic of this document.
- From latent variable  $\beta_k$ :
  - We will know the hot words in this topic.
  - We can infer what this topic is about.
- Then we classify the documents, with no supervision.
- An important unsupervised learning approach to content analysis.



# Latent Dirichlet Allocation (LDA)

[PDF] **Latent dirichlet allocation**

DM Blei, AY Ng, MI Jordan - Journal of machine Learning research, 2003 - jmlr.org

We describe **latent Dirichlet allocation** (LDA), a generative probabilistic model for collections of discrete data such as text corpora. LDA is a three-level hierarchical Bayesian model, in ...

☆ Save    ↗ Cite    **Cited by 46501**    Related articles    All 90 versions    ➞



David Blei



Michael I. Jordan



Andrew Ng

# Outline

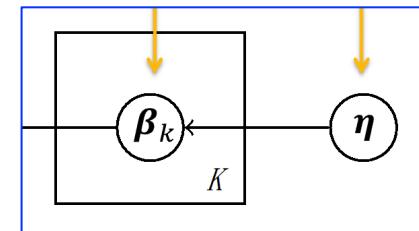
- Probabilistic Topic Model
  - Dirichlet-Multinomial Model
  - Latent Dirichlet Allocation
  - Document Generation Process
- Variational Inference
  - Mean Field Variational Inference
  - Variational Inference for LDA
  - Parameter Estimation for LDA



# LDA for Topic Modeling

How many topics?  
Hyperparameter!

Dirichlet ( $\eta$ )



$\beta_1$

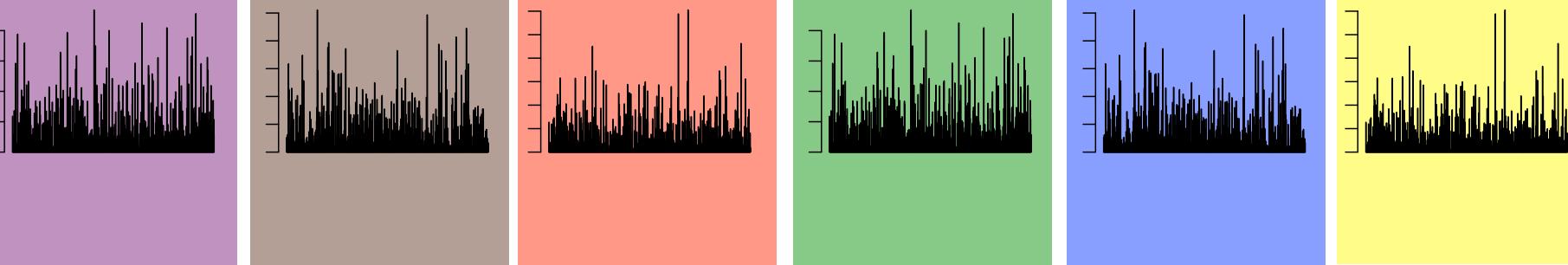
$\beta_2$

$\beta_3$

$\beta_4$

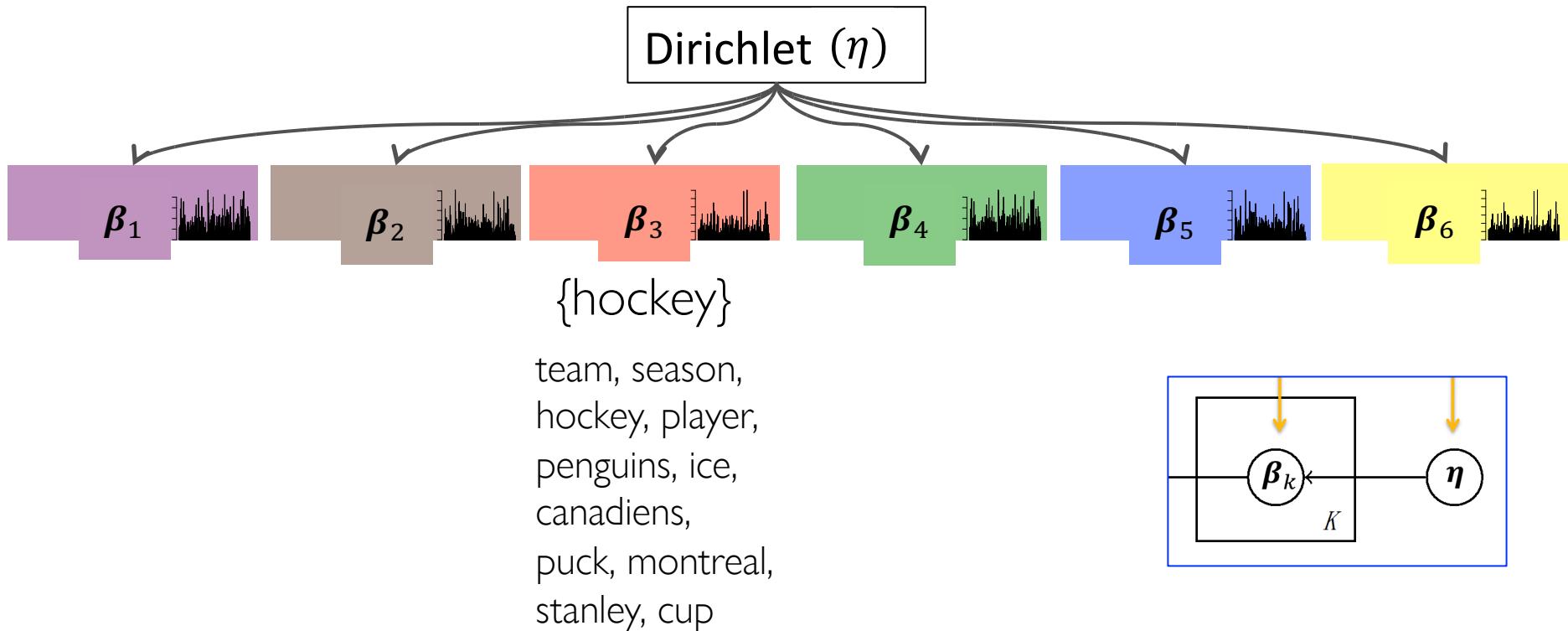
$\beta_5$

$\beta_6$



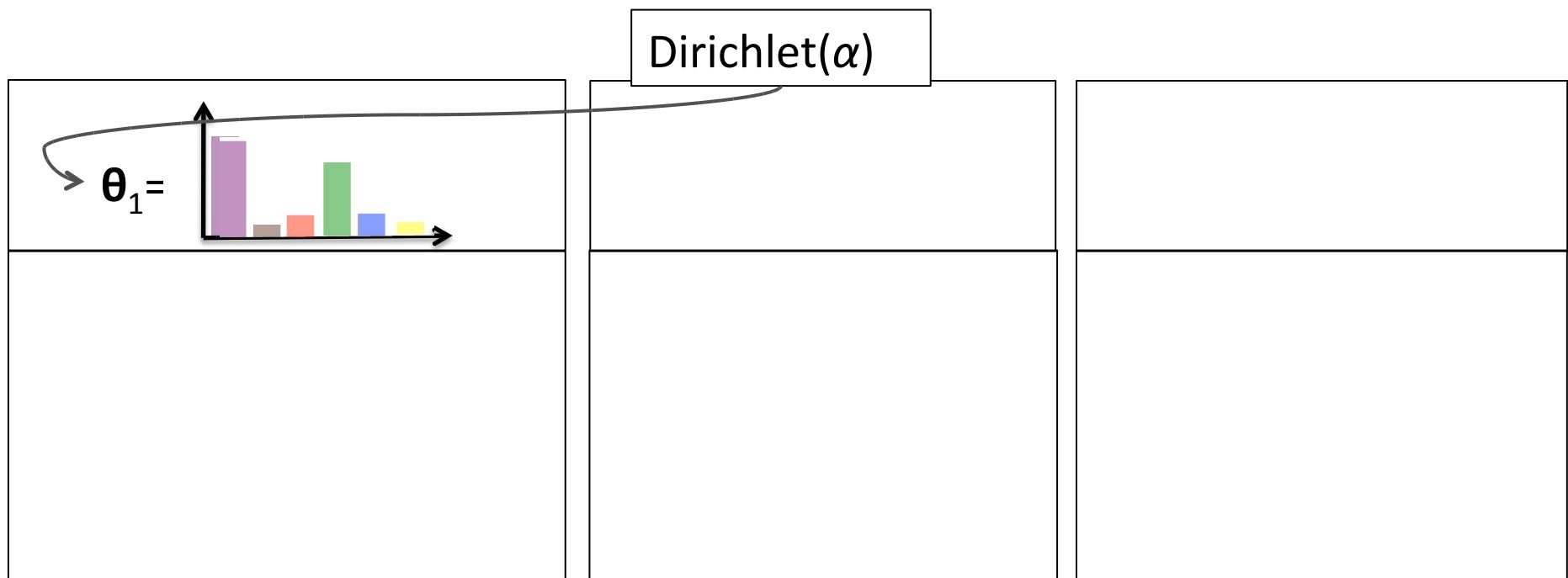
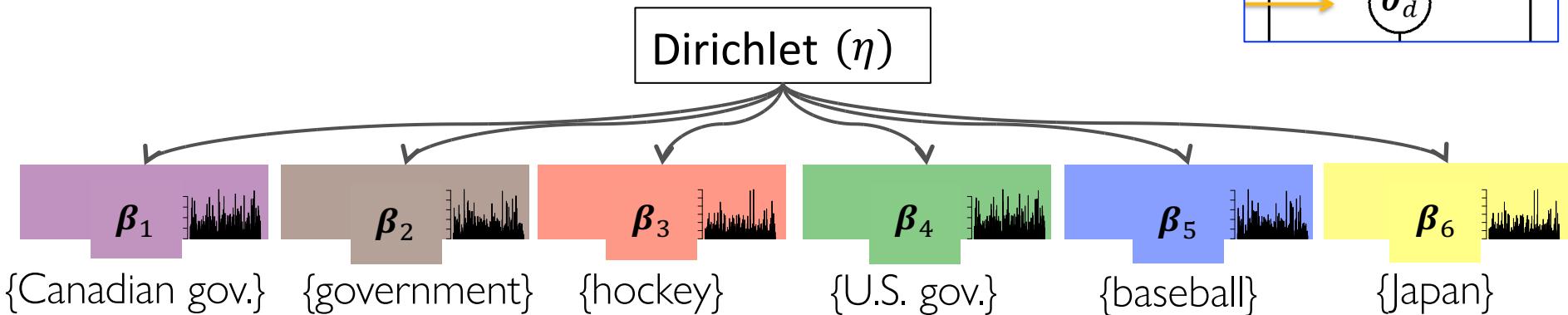
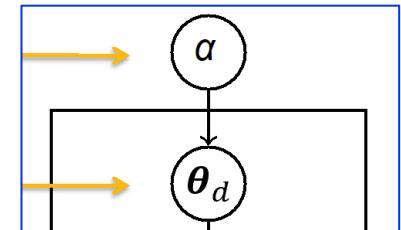
- The **generative story** begins with a Dirichlet prior over the topics.
- Each **topic** is defined as a **Multinomial distribution** over the vocabulary, parameterized by  $\beta_k$ .

# LDA for Topic Modeling

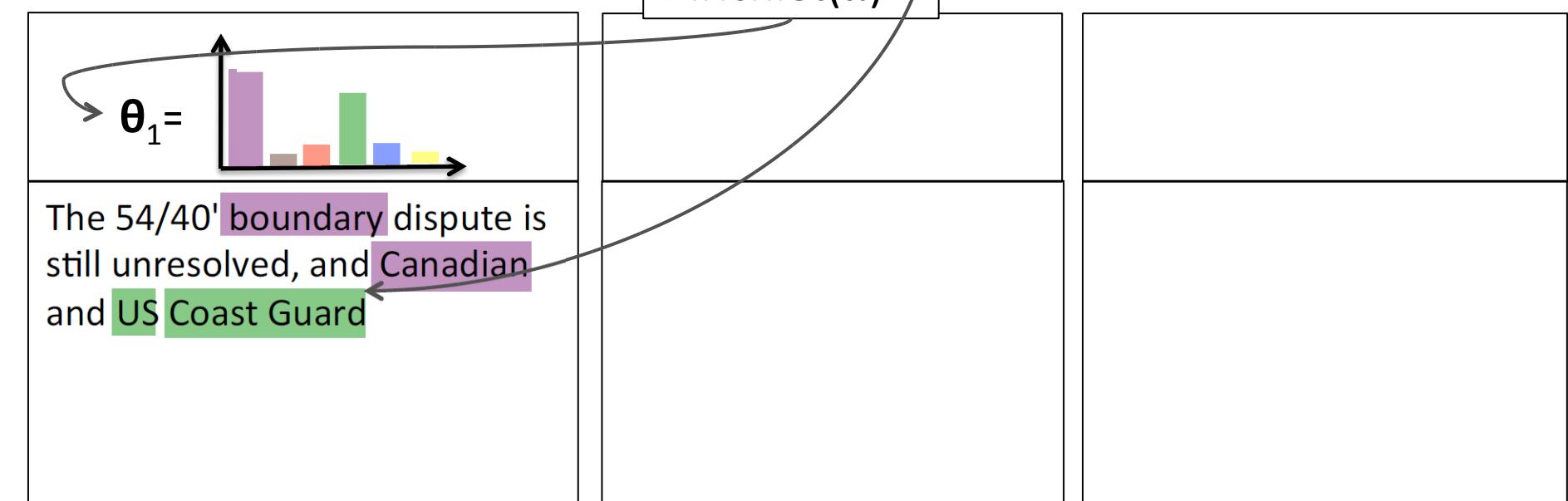
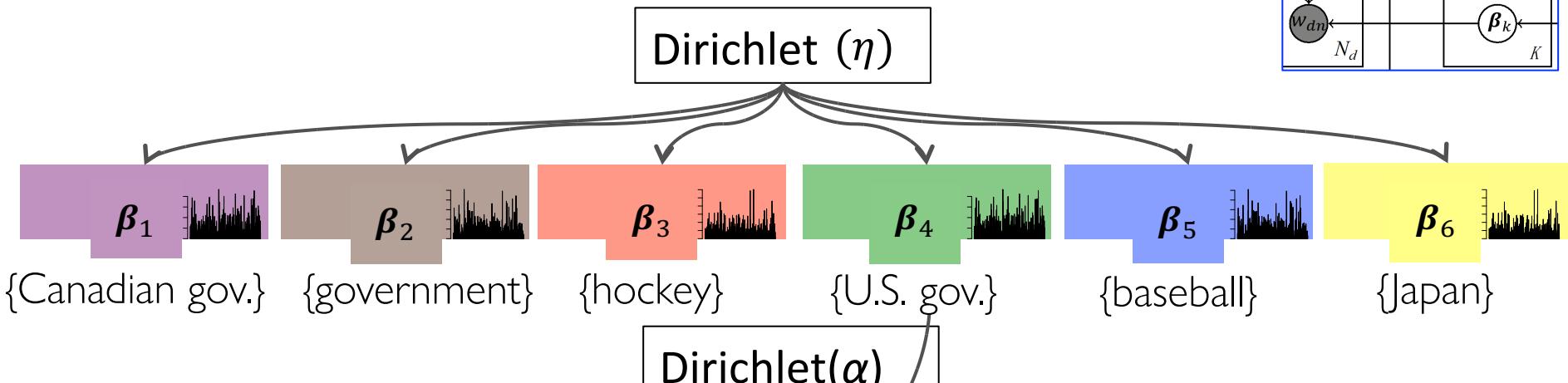


- A topic is visualized as its **high probability words**.
- A **pedagogical label** is used to identify the topic.

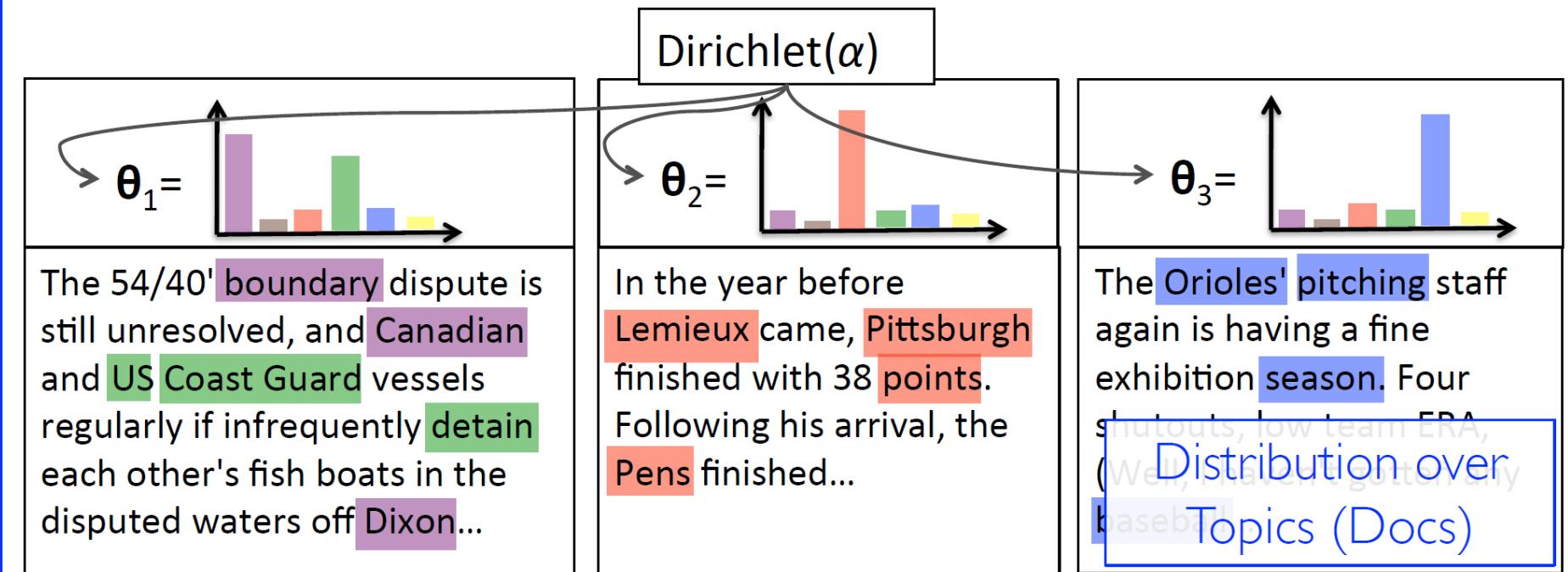
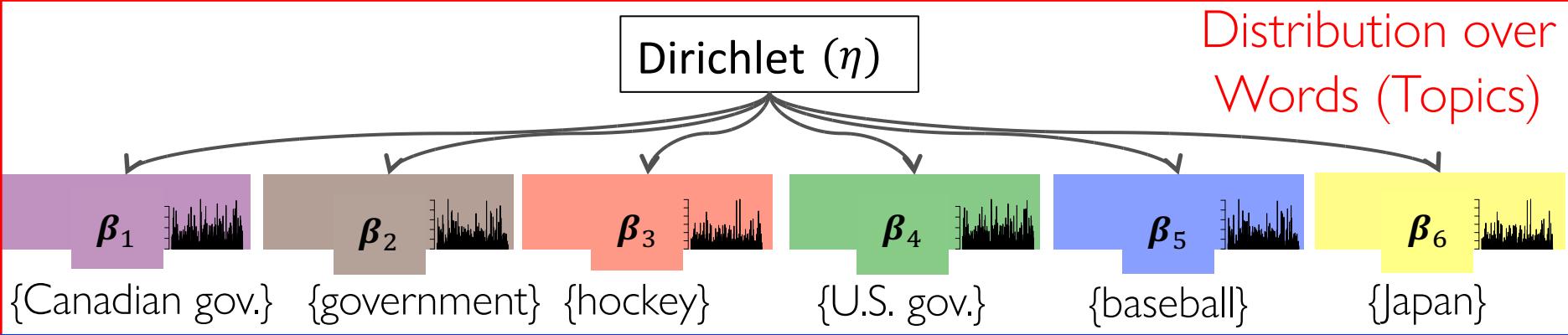
# LDA for Topic Modeling



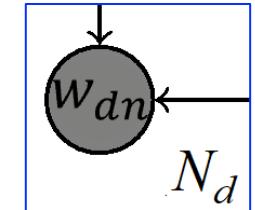
# LDA for Topic Modeling



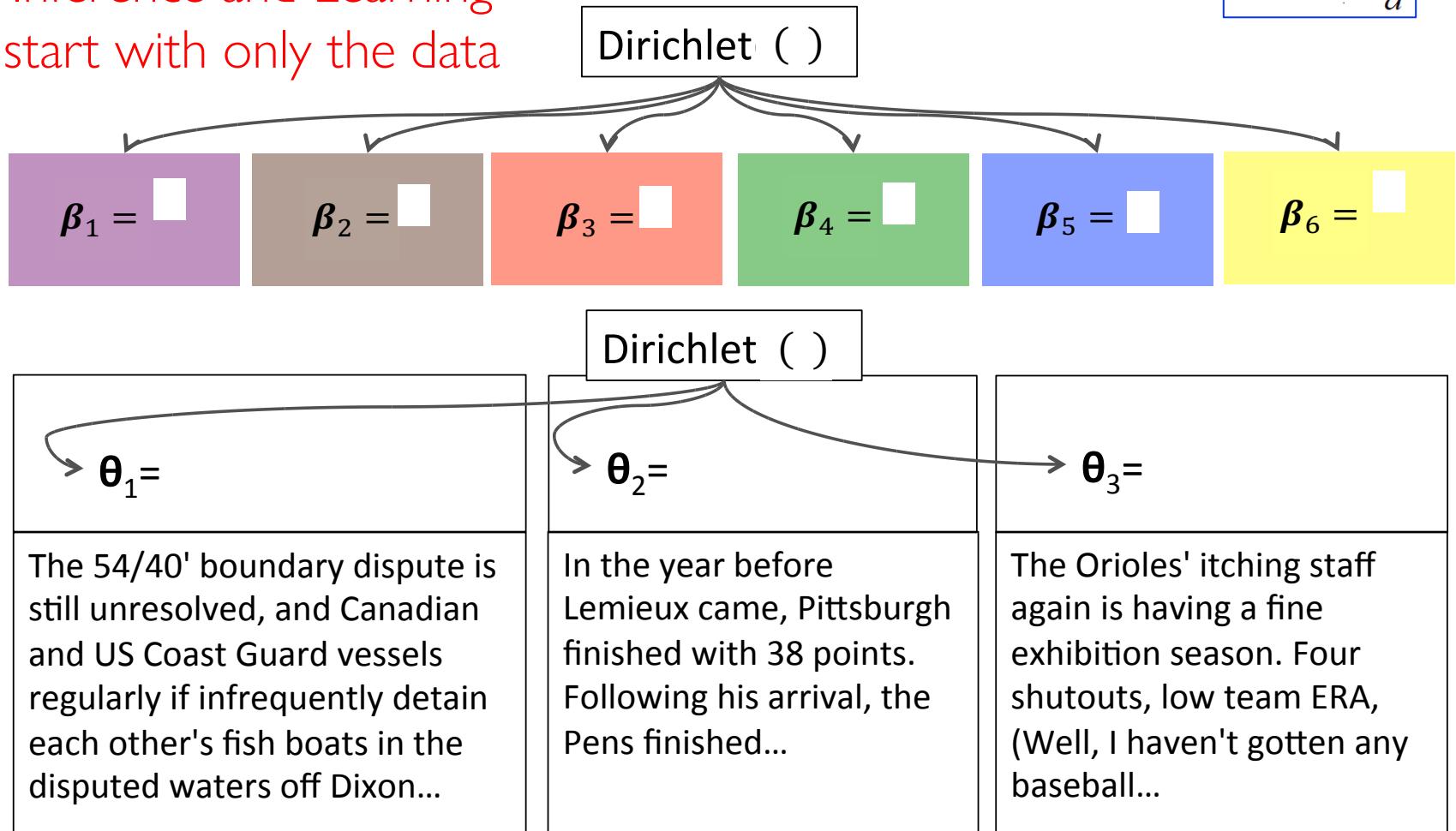
# LDA for Topic Modeling



# LDA for Topic Modeling



Inference and Learning  
start with only the data



# LDA for NIPS Papers

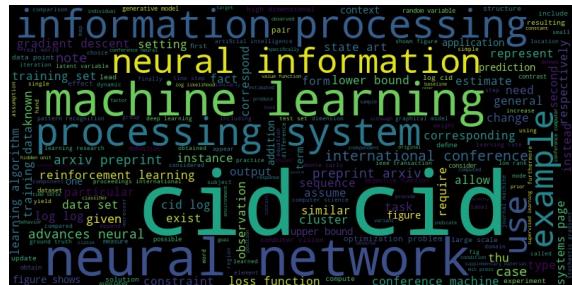
- <https://proceedings.neurips.cc/paper/2015>

NeurIPS Proceedings ➔

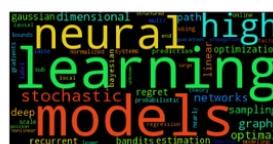
Advances in Neural Information Processing Systems 28 (NIPS 2015)

Edited by: C. Cortes and N. Lawrence and D. Lee and M. Sugiyama and R. Garnett

- Quartz: Randomized Dual Coordinate Ascent with Arbitrary Sampling *Zheng Qu, Peter Richtarik, Tong Zhang*
  - Associative Memory via a Sparse Recovery Model *Arya Mazumdar, Ankit Singh Rawat*
  - Policy Gradient for Coherent Risk Measures *Aviv Tamar, Yinlam Chow, Mohammad Ghavamzadeh, Shie Mannor*
  - A fast, universal algorithm to learn parametric nonlinear embeddings *Miguel A. Carreira-Perpinan, Max Vladymyrov*



# Topic I: Neural

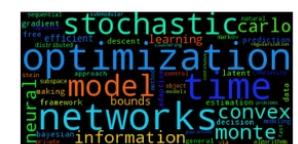


sample selection  
deep learning  
methods, bandits  
models, gaussian  
spectral, gradient, rank  
learning, bayesian  
inference, data  
adaptive, variational  
supervised, causal  
selection, random  
networks, hidden, semi  
latent, sparse, sparse

## Topic 2: Bayesian



# Topic 3: Inference



# Topic 4: Optimization

# LDA for NIPS Papers

- <https://cs.stanford.edu/people/karpathy/nips2015/>

Below every paper are TOP 100 most-occurring words in that paper and their color is based on LDA topic model with k = 7.  
(This is very hard but it looks like 0 = graphical models?, 1 = reinforcement learning?, 2 = deep learning, 3 = kernels?, 4 = theory?, 5 = optimization, 6 = matrix factorization?)

Toggle LDA topics to sort by: [TOPIC0](#) [TOPIC1](#) [TOPIC2](#) [TOPIC3](#) [TOPIC4](#) [TOPIC5](#) [TOPIC6](#)

Also note that you can filter by the day each poster appears: click on the day next to a paper to filter by that day.

## A Nonconvex Optimization Framework for Low Rank Matrix Estimation

Tuo Zhao, Zhaoran Wang, Han Liu



[pdf] [rank by tf-idf similarity to this]

[abstract]

Poster: Mon Dec 7th 07:00 - 11:59 PM @ 210 C #101

[exact, observation, key] [framework, based, existing, computationally, class, set] [arxiv, preprint] [log, computational, divergence, sample, family, parameter, introduce, analytic, point] [lemma, proof, algorithm, corollary, implies, prove, exists, complexity, number] [optimization, nonconvex, gradient, minimization, oracle, convex, projected, descent, convergence, rate, problem, global, max, satisfy, main, assumption, iteration, coordinate, assume, large, linear, updating, converge, general, step, established, optimum] [matrix, alternating, rank, establish, low, singular, theorem, completion, min, estimation, decomposition, appendix, geometric, provided, suppose, factorization, sensing, condition, error, high, orthonormal, result, broad, ieee, benjamin, entry, nsf, smaller, establishes, theoretical, including, solving]

0 =graphical models?, 1 = reinforcement learning?, 2 = deep learning, 3 = kernels?, 4 = theory?, 5 = optimization, 6 = matrix factorization?

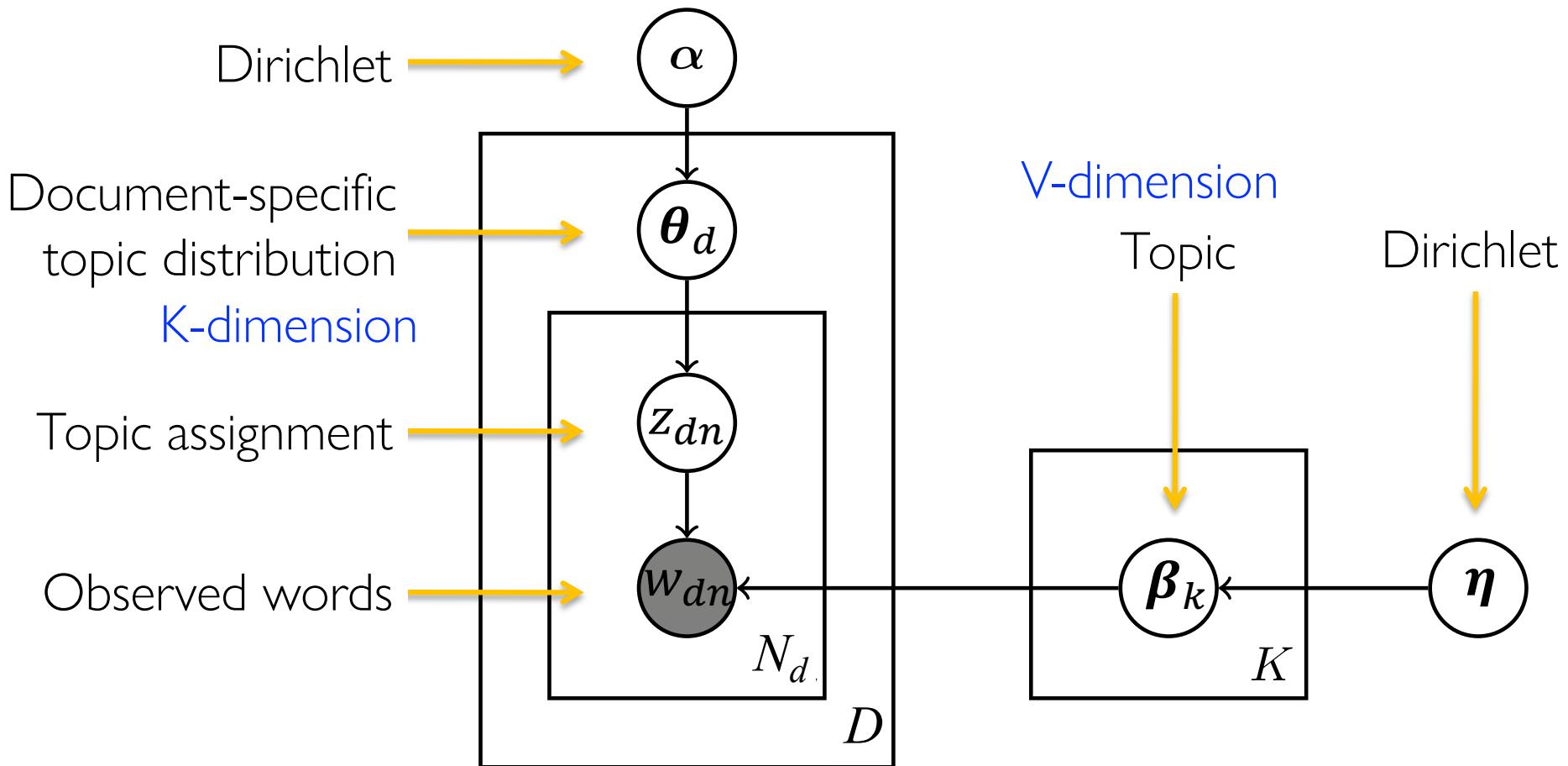


# Outline

- Probabilistic Topic Model
  - Dirichlet-Multinomial Model
  - Latent Dirichlet Allocation
  - Document Generation Process
- Variational Inference
  - Mean Field Variational Inference
  - Variational Inference for LDA
  - Parameter Estimation for LDA



# Latent Dirichlet Allocation (LDA)



Can you explain graph of LDA now?

$$p(\beta, \theta, z, w | \alpha, \eta) = p(\beta | \eta) p(\theta | \alpha) p(z | \theta) p(w | z, \beta)$$

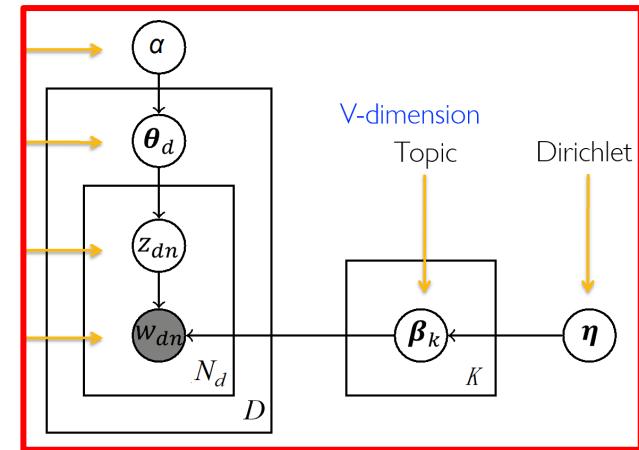
# Maximum Likelihood Estimation

- If we know all the latent variables:
  - (We do not know them in fact.)
  - Let's try to compute MLE as an exercise!

$$\log p(\beta, \theta, z, w | \alpha, \eta)$$

$$= \log p(\beta | \eta) + \log p(\theta | \alpha) + \log p(z | \theta) + \log p(w | z, \beta)$$

$$\begin{aligned}
 &= \sum_{k=1}^K \log p(\vec{\beta}_k | \eta) + \sum_{d=1}^D \log p(\vec{\theta}_d | \alpha) + \sum_{d=1}^D \sum_{n=1}^{N_d} \log p(z_{d,n} | \vec{\theta}_d) \\
 &\quad + \sum_{d=1}^D \sum_{n=1}^{N_d} \log p(w_{d,n} | z_{d,n}, \vec{\beta}_{1:K})
 \end{aligned}$$



PGM to simplify learning and inference!

# Complex MLE Objective

- Recall the form of Dirichlet Distribution:  $p(\vec{\theta}|\boldsymbol{\alpha}) = \frac{1}{B(\boldsymbol{\alpha})} \prod_{k=1}^K \theta_k^{\alpha_k - 1}$   
词表上有V个词

$$= \sum_{k=1}^K \left( \sum_{v=1}^V (\eta_v - 1) \log \beta_{kv} - \log B(\boldsymbol{\eta}) \right) + \sum_{d=1}^D \sum_{k=1}^K (\alpha_k - 1) \log \theta_{dk} - \log B(\boldsymbol{\alpha})$$

beta\_z: V维

$$+ \sum_{d=1}^D \sum_{n=1}^{N_d} \log \theta_{d,z_{d,n}} + \sum_{d=1}^D \sum_{n=1}^{N_d} \log \beta_{z_{d,n} w_{d,n}}$$

- We will meet some more complex terms later ☺
- We need to solve for parameters  $\boldsymbol{\alpha}, \boldsymbol{\eta}$ .

Newton Method

$$\operatorname{argmax}_{\boldsymbol{\eta}} \sum_{k=1}^K \left( \sum_{v=1}^V (\eta_v - 1) \log \beta_{kv} - \log B(\boldsymbol{\eta}) \right)$$



# ELBO and EM Algorithm

- Evidence Lower Bound (ELBO):

$$\begin{aligned}\mathcal{L}(q, \theta) &= \sum_z q(z) \log \left( \frac{p(x, z | \theta)}{q(z)} \right) \\ &= -\text{KL}[q(z) || p(z|x, \theta)] + \log p(x|\theta)\end{aligned}$$

- EM Algorithm: 必考

- Choose initial  $\theta^{\text{old}}$

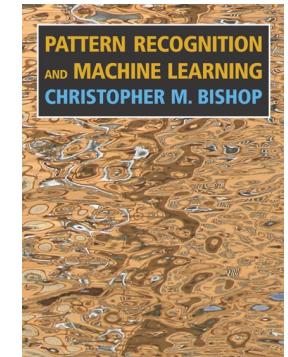
- E-step: Let  $q^* = \underset{q}{\operatorname{argmax}} \mathcal{L}(q, \theta^{\text{old}})$

Minimize

Optimal choice for E-step  
 $q^*(z) = p(z|x, \theta^{\text{old}})$

- M-step: Let  $\theta^{\text{new}} = \underset{\theta}{\operatorname{argmax}} \mathcal{L}(q^*, \theta)$

- Go to step 2, until converged.



PRML

Chapter 10

# Intractable E-Step

- We need to solve: Optimal choice for E-step  $q^*(z) = p(z|x, \theta^{\text{old}})$

$$\text{In LDA: } p(\theta, z, \beta | w, \alpha, \eta) = \frac{p(\theta, z, \beta, w | \alpha, \eta)}{p(w | \alpha, \eta)}$$

- But the denominator is **intractable!**

$$p(w|\alpha, \eta) = \int \int \sum_z p(\theta, z, \beta, w | \alpha, \eta) d\theta d\beta$$

Inference is hard  
for general  
Bayesian Models!

- We need **approximation** here!
  - There are two ways to solve E-step approximately:
    - Variational Inference (this lecture)
    - Sampling – Markov Chain Monte Carlo (next lecture) 会考

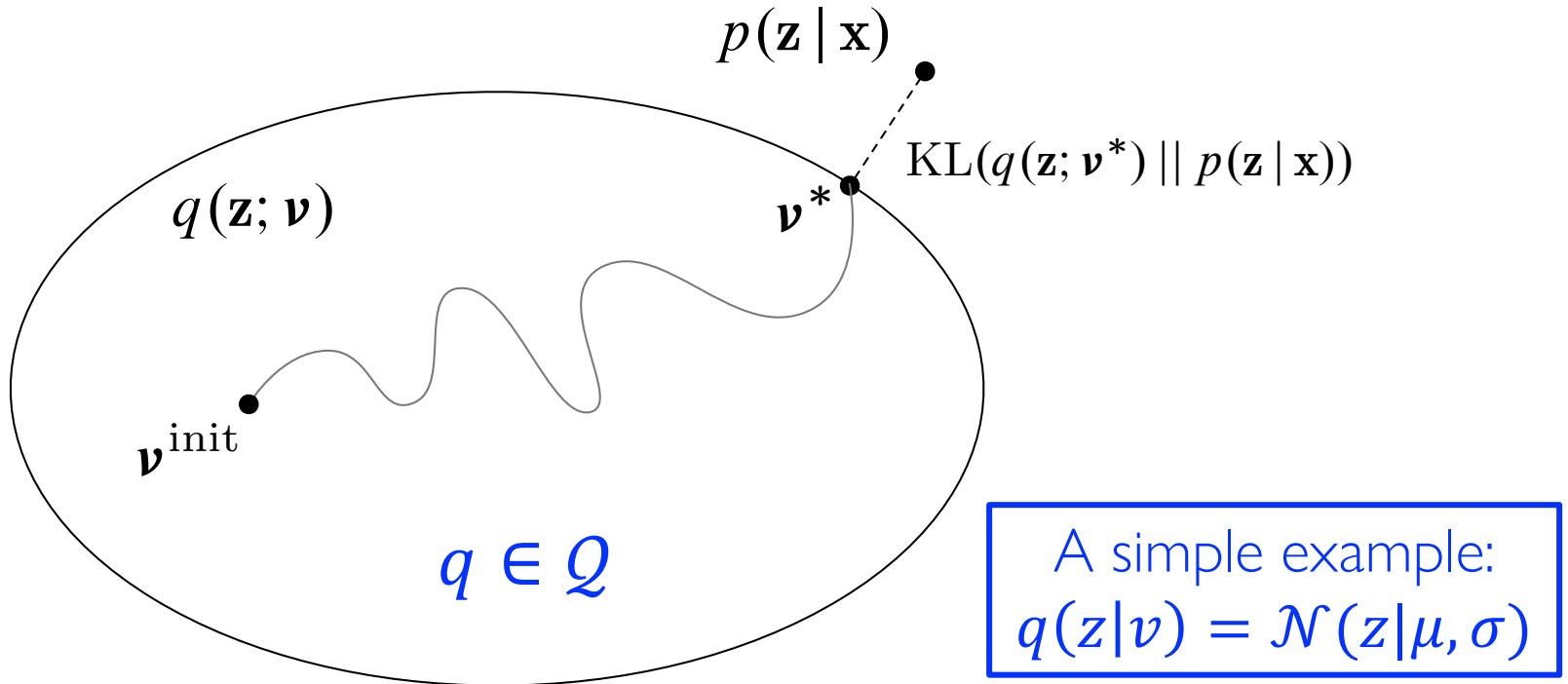


# Variational Methods

- Suppose E-step is difficult:
  - Hard to take expectation w.r.t.  $q^*(z) = p(z|x, \theta^{\text{old}})$
  - For example, topic models.
- Solution: Restrict to distributions  $\mathcal{Q}$  that are easy to work with.
- The ELBO now **looser**:
  - Intractable!
  - We compute this
$$q^* = \underset{q \in \mathcal{Q}}{\operatorname{argmin}} \text{KL} [q(z), p(z|x, \theta^{\text{old}})] = \underset{q \in \mathcal{Q}}{\operatorname{argmax}} \mathcal{L}(q, \theta^{\text{old}})$$
- Find an easy-to-work **variational distribution  $q^*$**  to approximate the inference distribution  $p(z|x, \theta^{\text{old}})$ .
  - This group of methods are called **Variational Methods**.



# Variational Inference (VI)



- A variational distribution  $q(z; \nu)$  along with its parameter  $\nu$  is used.
  - Objective in VI: Find  $\nu$  to maximize the ELBO.
- $\nu$  may have no relation with  $\theta^{\text{old}}$ , which makes VI easier to solve.

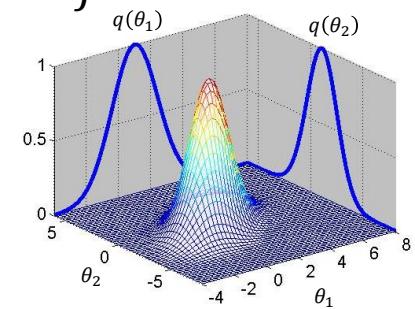


# Mean Field Theory

不要求

- Assume  $\mathcal{Q}$  a distribution  $q(z)$  that factorizes  $q(z) = \prod_{i=1}^m q_i(z_i)$ .
- Variational inference of  $q(z)$  corresponds to the mean field theory.
- Seek all factors  $\{q_i = q_i(z_i)\}$  to maximize ELBO:  $q^* = \underset{q \in \mathcal{Q}}{\operatorname{argmax}} \mathcal{L}(q)$

$$\begin{aligned}\mathcal{L}(q) &= \int \prod_i q_i \log \frac{p(x, z)}{\prod_i q_i} dz = \int \prod_i q_i \left\{ \log p(x, z) - \sum_i \log q_i \right\} dz \\ &= \int q_j \left\{ \int \log p(x, z) \prod_{i \neq j} q_i dz_i \right\} dz_j - \int q_j \log q_j dz_j + C \\ &= \int q_j \{ \log \tilde{p}(x, z_j) \} dz_j - \int q_j \log q_j dz_j + C\end{aligned}$$

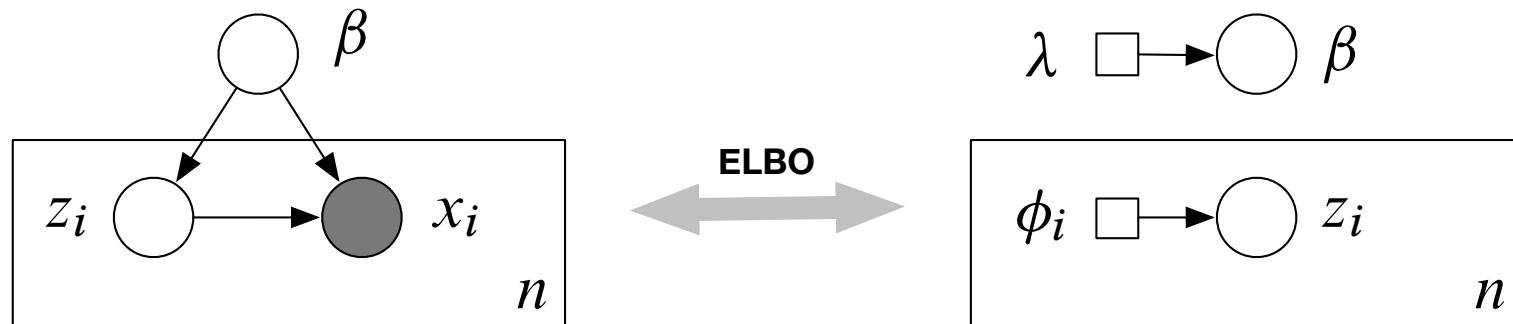




# Mean Field Theory

- ELBO:  $\mathcal{L}(q) = \int q_j \{\log \tilde{p}(x, z_j)\} dz_j - \int q_j \log q_j dz_j + C$ 
  - Where  $\log \tilde{p}(x, z_j) = \int \log p(x, z) \prod_{i \neq j} q_i dz_i \triangleq \mathbb{E}_{i \neq j}[\log p(x, z)]$
- We keep  $\{q_{j \neq i}\}$  fixed and maximize  $\mathcal{L}(q)$  for all  $q_j(z_j)$ :
  - Note that  $\mathcal{L}(q) = -\text{KL}(q_j(z_j) \| \tilde{p}(x, z_j)) + C$ 
    - The maximizer of  $\mathcal{Q}$  to ELBO is  $q_j(z_j) = \tilde{p}(x, z_j)$
  - A general expression for the optimal solution  $q_j^*(z_j)$ :
 
$$\log q_j^*(z_j) = \mathbb{E}_{i \neq j}[\log p(x, z)] + C$$
- In physics, the mean field theory studies the behavior of high-dimensional stochastic models by studying a simpler model that approximates the original one by averaging over degrees of freedom.

# Mean Field Variational Inference

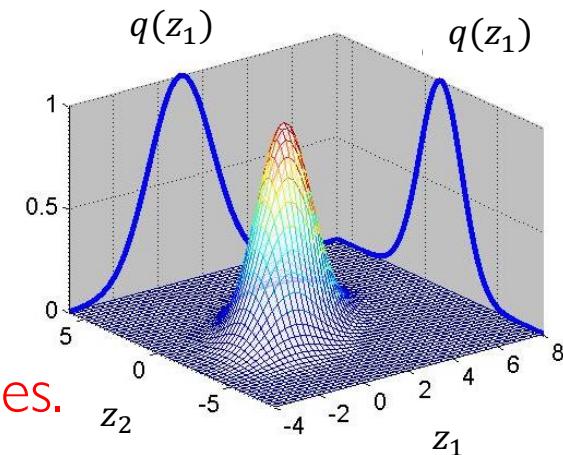


- Consider a simpler graphical model.

$$p(\beta, \mathbf{z}, \mathbf{x}) = p(\beta) \prod_{i=1}^n p(z_i, x_i | \beta)$$

- The **difficulty** to tackle this distribution:

- Latent variable depends on other latent variables.



- Mean Field Assumption tries to disentangle latent variables to simplify:

$$q(\beta, \mathbf{z}; \lambda, \phi) = q(\beta; \lambda) \prod_{i=1}^n q(z_i; \phi_i)$$

可以假设beta  
是高斯分布

No dependency!

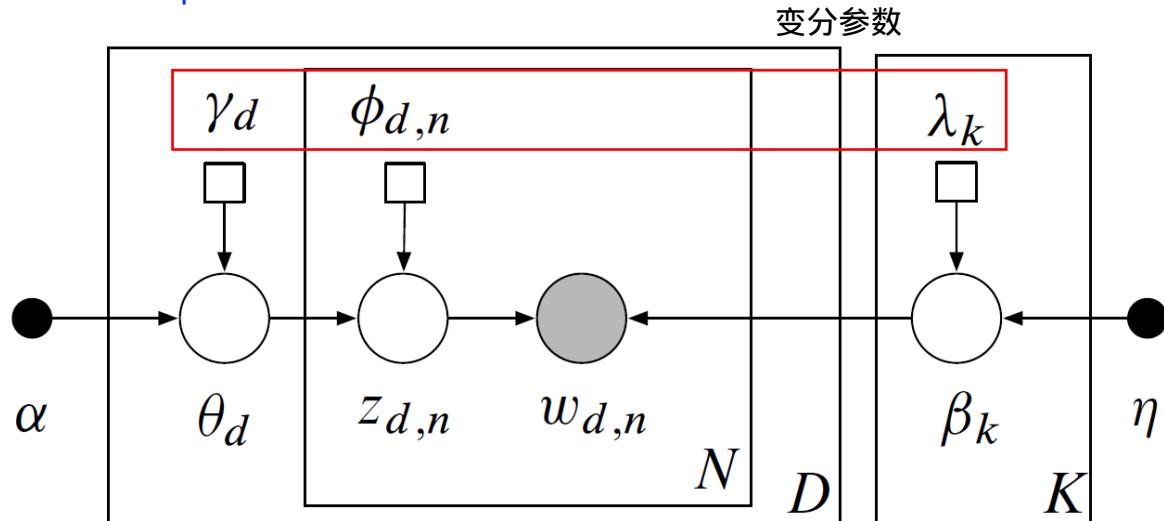
# Outline

- Probabilistic Topic Model
  - Dirichlet-Multinomial Model
  - Latent Dirichlet Allocation
  - Document Generation Process
- Variational Inference
  - Mean Field Variational Inference
  - **Variational Inference for LDA**
  - Parameter Estimation for LDA



# Variational Inference for LDA

- Mean field assumption in LDA:



$$q(\vec{\theta}_{1:D}, \mathbf{z}_{1:D, 1:N_d}, \vec{\beta}_{1:K}) = \prod_{k=1}^K q(\vec{\beta}_k | \vec{\lambda}_k) \prod_{d=1}^D q(\vec{\theta}_d | \vec{\gamma}_d) \prod_{d=1}^D \prod_{n=1}^N q(z_{d,n} | \vec{\Phi}_{d,n})$$

Dirichlet                      Dirichlet                      Multinomial  
Distribution                      Distribution                      Distribution

# Variational Inference for LDA

- Recall the form of ELBO:

$$\mathcal{L}(q, \theta) = \sum_z q(z) \log \left( \frac{p(x, z | \theta)}{q(z)} \right) = \mathbb{E}_{z \sim q} \log \left( \frac{p(x, z | \theta)}{q(z)} \right)$$

- In LDA, we have known the form of  $q$ , and what left is:

$$p(\beta, \theta, z, w | \alpha, \eta) = p(\beta | \eta) p(\theta | \alpha) p(z | \theta) p(w | z, \beta)$$

- The objective function  $\mathcal{L}$  in VI of LDA (Solvable!) turns out to be:

$$\begin{aligned} \mathcal{L} = & \sum_{k=1}^K E_q [\log p(\vec{\beta}_k | \eta)] + \sum_{d=1}^D E_q [\log p(\vec{\theta}_d | \vec{\alpha})] + \\ & \sum_{d=1}^D \sum_{n=1}^N E_q [\log p(z_{d,n} | \vec{\theta}_d)] + \\ & \sum_{d=1}^D \sum_{n=1}^N E_q [\log p(w_{d,n} | z_{d,n}, \vec{\beta}_{1:K})] - E_q \log q \end{aligned}$$



# Variational Inference for LDA

- We solve one of the terms as an example:

$$E_q[\log p(\boldsymbol{\theta}|\boldsymbol{\alpha})]$$

All latent variables  
are gone!

$$\begin{aligned}&= E_q\left[\log \left( \exp \left\{ \left( \sum_{i=1}^K (\alpha_i - 1) \log \theta_i \right) + \log \Gamma(\sum_{i=1}^K \alpha_i) - \sum_{i=1}^K \log \Gamma(\alpha_i) \right\} \right) \right], \\&= \left( \sum_{i=1}^K (\alpha_i - 1) E_q[\log \theta_i] \right) + \log \Gamma(\sum_{i=1}^K \alpha_i) - \sum_{i=1}^K \log \Gamma(\alpha_i), \\&= \left( \sum_{i=1}^K (\alpha_i - 1) (\Psi(\gamma_i) - \Psi(\sum_j \gamma_j)) \right) + \log \Gamma(\sum_{i=1}^K \alpha_i) - \sum_{i=1}^K \log \Gamma(\alpha_i).\end{aligned}$$

- $\Psi$  is the first derivative of the log Gamma function  $\Gamma$ .
- About the last step please refer to Section A.1 of the original paper ☺

# Variational Inference for LDA

- The objective function  $\mathcal{L}$  in VI of LDA turns out to be something like:



$$\begin{aligned} L(\gamma, \phi; \alpha, \beta) = & \log \Gamma\left(\sum_{j=1}^k \alpha_j\right) - \sum_{i=1}^k \log \Gamma(\alpha_i) + \sum_{i=1}^k (\alpha_i - 1) \left( \Psi(\gamma_i) - \Psi\left(\sum_{j=1}^k \gamma_j\right) \right) \\ & + \sum_{n=1}^N \sum_{i=1}^k \phi_{ni} \left( \Psi(\gamma_i) - \Psi\left(\sum_{j=1}^k \gamma_j\right) \right) \\ & + \sum_{n=1}^N \sum_{i=1}^k \sum_{j=1}^V \phi_{ni} w_n^j \log \beta_{ij} \\ & - \log \Gamma\left(\sum_{j=1}^k \gamma_j\right) + \sum_{i=1}^k \log \Gamma(\gamma_i) - \sum_{i=1}^k (\gamma_i - 1) \left( \Psi(\gamma_i) - \Psi\left(\sum_{j=1}^k \gamma_j\right) \right) \\ & - \sum_{n=1}^N \sum_{i=1}^k \phi_{ni} \log \phi_{ni}, \end{aligned}$$

You may read the original paper to find the details

- Bayesians face something like this everyday 😞
- Although this is scaring, Just use Calculus or SGD!  $\nabla L = 0$

# Variational Inference for LDA

- Iterate the following steps (Descend each variable iteratively):

(1) For each topic  $k$  and term  $v$ :

$$\lambda_{k,v}^{(t+1)} = \eta + \sum_{d=1}^D \sum_{n=1}^N \mathbb{1}(w_{d,n} = v) \phi_{n,k}^{(t)}. \quad \text{隐变量对应的变分分布}$$

(2) For each document  $d$ :

(a) Update  $\gamma_d$ :

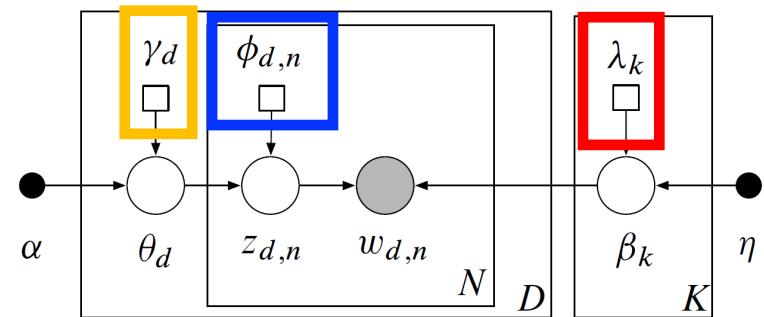
$$\gamma_{d,k}^{(t+1)} = \alpha_k + \sum_{n=1}^N \phi_{d,n,k}^{(t)}.$$

(b) For each word  $n$ , update  $\vec{\phi}_{d,n}$ :

$$\phi_{d,n,k}^{(t+1)} \propto \exp \left\{ \Psi(\gamma_{d,k}^{(t+1)}) + \Psi(\lambda_{k,w_n}^{(t+1)}) - \Psi(\sum_{v=1}^V \lambda_{k,v}^{(t+1)}) \right\},$$

where  $\Psi$  is the digamma function, the first derivative of the log  $\Gamma$  function.

$$\sum_{k=1}^V \phi_{d,n,k} = 1$$



Digamma function  $\Psi$  can only be computed by approximation, e.g. Taylor

# Outline

- Probabilistic Topic Model
  - Dirichlet-Multinomial Model
  - Latent Dirichlet Allocation
  - Document Generation Process
- Variational Inference
  - Mean Field Variational Inference
  - Variational Inference for LDA
  - Parameter Estimation for LDA



# EM Algorithm

- ELBO:

$$\begin{aligned}\mathcal{L}(q, \theta) &= \sum_z q(z) \log \left( \frac{p(x, z | \theta)}{q(z)} \right) \\ &= -\text{KL}[q(z) || p(z|x, \theta)] + \log p(x|\theta)\end{aligned}$$

- EM Algorithm with Variational Inference:

- Choose initial  $\theta^{\text{old}}$

- Let  $q^* = \underset{q}{\operatorname{argmax}} \mathcal{L}(q, \theta^{\text{old}})$  We are here

- Let  $\theta^{\text{new}} = \underset{\theta}{\operatorname{argmax}} \mathcal{L}(q^*, \theta)$  What we left

- In LDA:  
Solve  $\alpha$  and  $\eta$

- Go to step 2, until converged.



# Parameter Estimation for LDA

- **Minimize** the part of ELBO that relates to  $\alpha$ :

$$L_{[\alpha]} = \sum_{d=1}^D \left( \log \Gamma \left( \sum_{j=1}^K \alpha_j \right) - \sum_{i=1}^K \log \Gamma (\alpha_i) + \sum_{i=1}^K \left( (\alpha_i - 1) \left( \Psi(\gamma_{di}) - \Psi \left( \sum_{j=1}^K \gamma_{dj} \right) \right) \right) \right)$$

- We can't find an analytical solution.

- Use **Newton method**.

Newton Iteration Function

$$x_{n+1} = x_n - \frac{f'(x_n)}{f''(x_n)}$$

- Compute **Derivatives** and **Hessian**:

$$\frac{\partial L}{\partial \alpha_i} = D \left( \Psi \left( \sum_{j=1}^K \alpha_j \right) - \Psi (\alpha_i) \right) + \sum_{d=1}^D \left( \Psi (\gamma_{di}) - \Psi \left( \sum_{j=1}^K \gamma_{dj} \right) \right)$$

$$\frac{\partial L}{\partial \alpha_i \alpha_j} = \delta(i, j) D \Psi' (\alpha_i) - \Psi' \left( \sum_{j=1}^K \alpha_j \right)$$

- The problem of  $\eta$  is as the same form of  $\alpha$ 's problem.



# LDA: Learning

- Iterate between:

- E-step: Iterative optimization

- (1) For each topic  $k$  and term  $v$ :

$$\lambda_{k,v}^{(t+1)} = \eta + \sum_{d=1}^D \sum_{n=1}^N 1(w_{d,n} = v) \phi_{n,k}^{(t)}.$$

- (2) For each document  $d$ :

- (a) Update  $\gamma_d$ :

$$\gamma_{d,k}^{(t+1)} = \alpha_k + \sum_{n=1}^N \phi_{d,n,k}^{(t)}.$$

- (b) For each word  $n$ , update  $\vec{\phi}_{d,n}$ :

$$\phi_{d,n,k}^{(t+1)} \propto \exp \left\{ \Psi(\gamma_{d,k}^{(t+1)}) + \Psi(\lambda_{k,w_n}^{(t+1)}) - \Psi(\sum_{v=1}^V \lambda_{k,v}^{(t+1)}) \right\},$$

- where  $\Psi$  is the digamma function, the first derivative of the log  $\Gamma$  function.

- M-step:

- Solve  $\alpha^{\text{new}}, \eta^{\text{new}} := \underset{\alpha, \eta}{\operatorname{argmax}} \mathcal{L}(q^*, \alpha, \eta)$  (Newton or SGD)

$$\frac{\partial L}{\partial \alpha_i \alpha_j} = \delta(i, j) D \Psi'(\alpha_i) - \Psi' \left( \sum_{j=1}^K \alpha_j \right)$$



# LDA: Inference

- After computing all parameters in the mean field approximation:

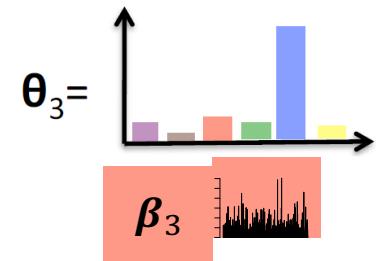
$$q(\vec{\theta}_{1:D}, \mathbf{z}_{1:D, 1:N_d}, \vec{\beta}_{1:K}) = \prod_{k=1}^K q(\vec{\beta}_k | \vec{\lambda}_k) \prod_{d=1}^D q(\vec{\theta}_d | \vec{\gamma}_d) \prod_{d=1}^D \prod_{n=1}^{N_d} q(z_{d,n} | \vec{\Phi}_{d,n})$$

- How do we compute  $\vec{\theta}_{1:D}, \mathbf{z}_{1:D, 1:N_d}, \vec{\beta}_{1:K}$ ? Argmax!

$$-\vec{\theta}_d = \operatorname{argmax}_{\vec{\theta}} q(\vec{\theta} | \vec{\gamma}_d), \vec{\beta}_k = \operatorname{argmax}_{\vec{\beta}} q(\vec{\beta} | \vec{\lambda}_k)$$

$$-z_{d,n} = \operatorname{argmax}_z q(z | \vec{\Phi}_{d,n})$$

- Variational Inference does not only enable EM algorithm.
  - The essential goal of VI is to do inference.



{hockey}

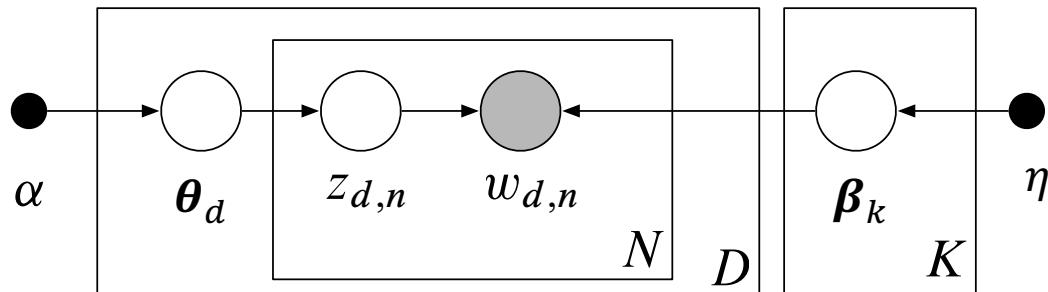
team, season,  
hockey, player,  
penguins, ice,  
canadiens,  
puck, montreal,  
stanley, cup



# LDA: Evaluation

不要求

- How to evaluation LDA?
  - First, we need a **held out corpora** for testing.
  - Which component of the model can **generalize to unseen data?**



- Of course, parameter  $\boldsymbol{\alpha}, \boldsymbol{\eta}$  can generalize.
  - $\boldsymbol{\beta}_{1:K}$  can also generalize to related corpus!
- Try to compute **the marginal likelihood of words** on a new dataset:

$$\text{perplexity } (D_{\text{test}}) = \exp \left\{ - \frac{\sum_{d=1}^{D_{\text{te}}} \log p(\mathbf{w}_d | \boldsymbol{\alpha}, \boldsymbol{\eta}, \boldsymbol{\beta})}{\sum_{d=1}^{D_{\text{te}}} N_d} \right\}$$

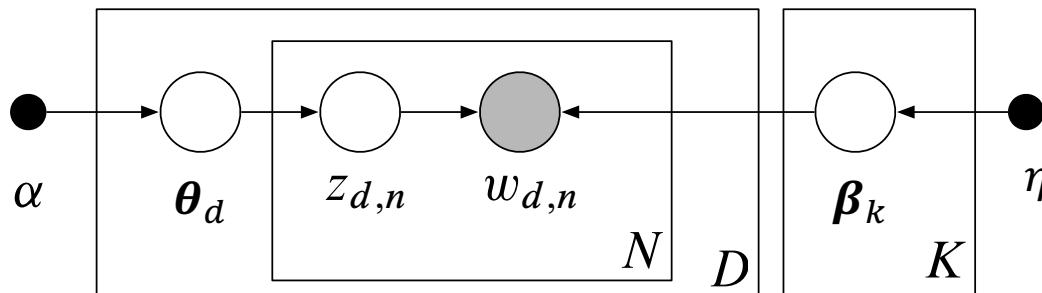
Lower is better





# Perplexity Metric

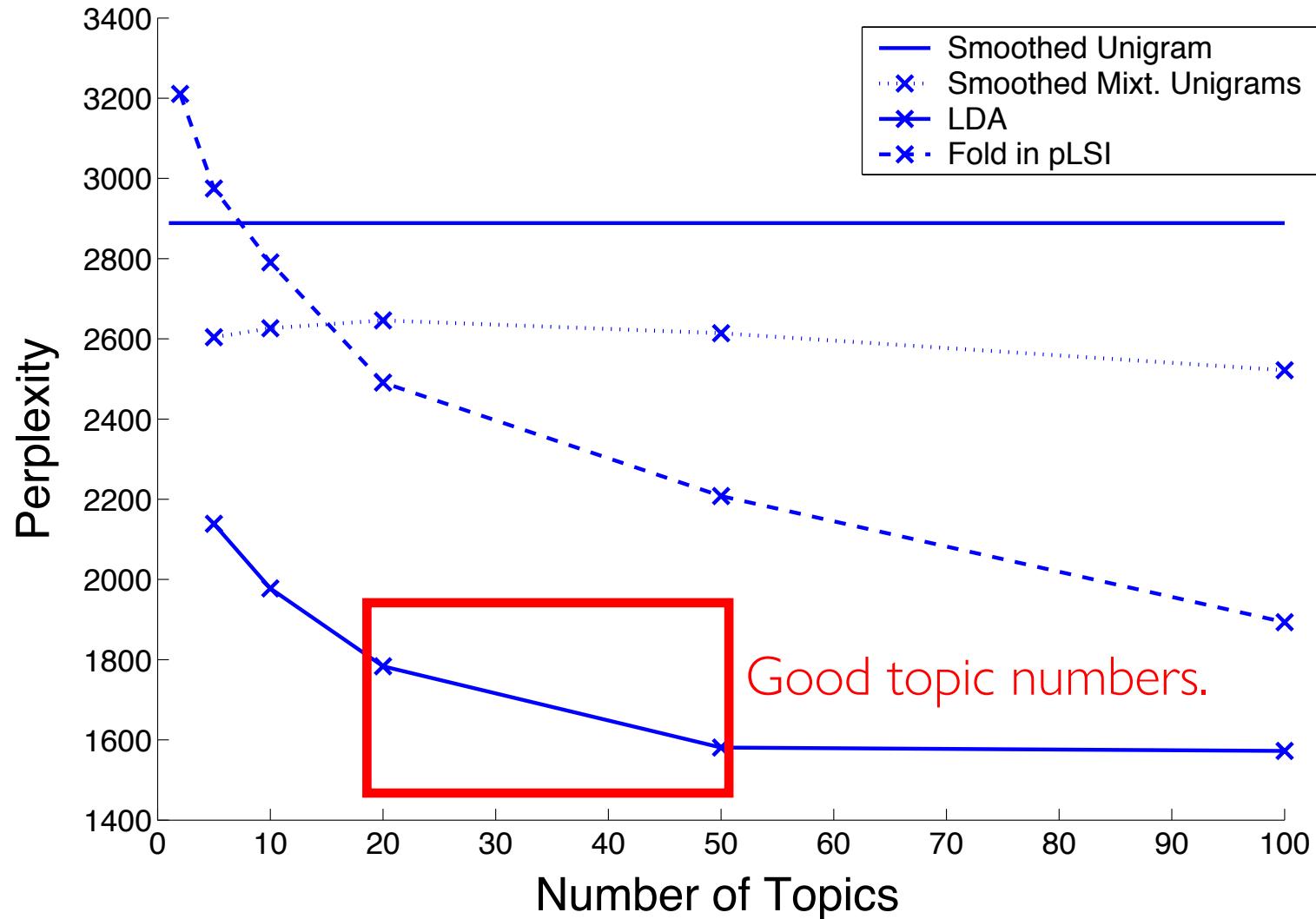
- With  $\alpha, \eta, \beta_{1:K}$ , we can use similar method to infer  $\theta, z$  on test data.



- After knowing all the latent variables, the marginal likelihood of words can be computed as:

$$\text{perplexity } (D_{\text{test}}) = \exp \left\{ - \frac{\sum_{d=1}^{D_{te}} \log p(\mathbf{w}_d | \mathbf{z}_d, \beta) p(\mathbf{z}_d | \boldsymbol{\theta}_d) p(\boldsymbol{\theta}_d | \alpha)}{\sum_{d=1}^{D_{te}} N_d} \right\}$$

# Performance Evaluation



# Thank You Questions?

Mingsheng Long

[mingsheng@tsinghua.edu.cn](mailto:mingsheng@tsinghua.edu.cn)

<http://ise.thss.tsinghua.edu.cn/~mlong>

答疑：东主楼11区413室