

20  
26  
22



University of Colorado **Boulder**

Department of Computer Science

CSCI 5622: Machine Learning

Chenhao Tan

Lecture 21: Variational Inference

Slides adapted from Jordan Boyd-Graber, Chris Ketelsen

# Learning Objectives

- Understand evaluations
- Understand the intuition behind variational inference

# Outline

- Evaluations of topic models
- Variational inference

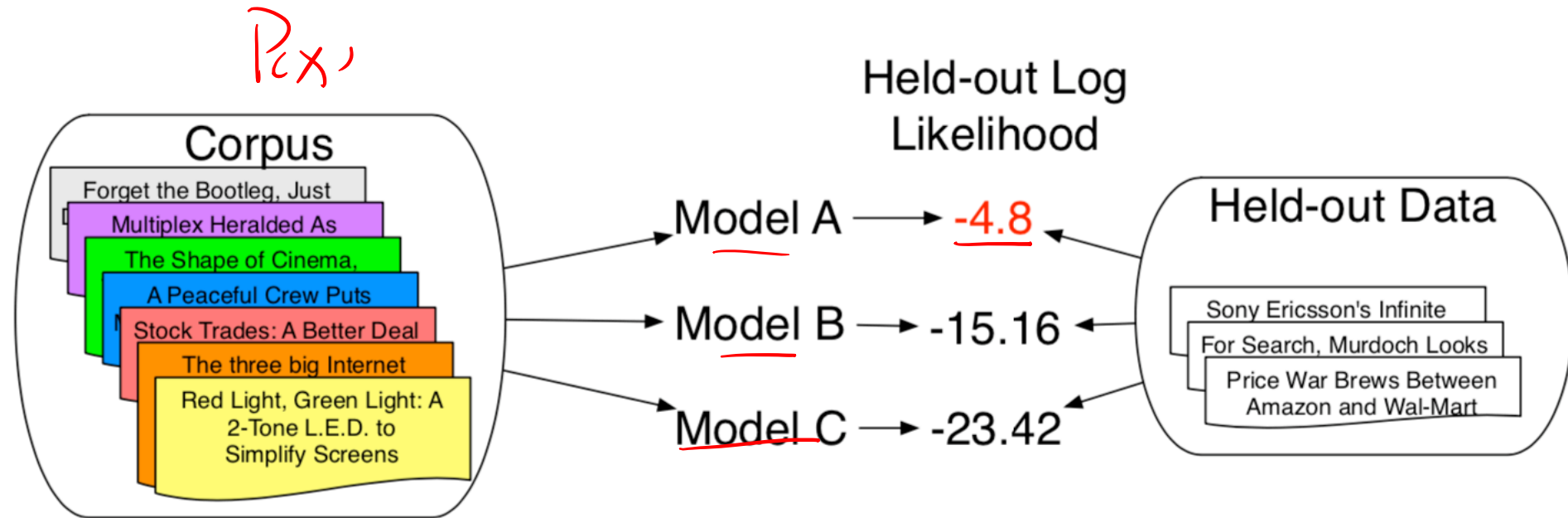
# Outline

- Evaluations of topic models
- Variational inference

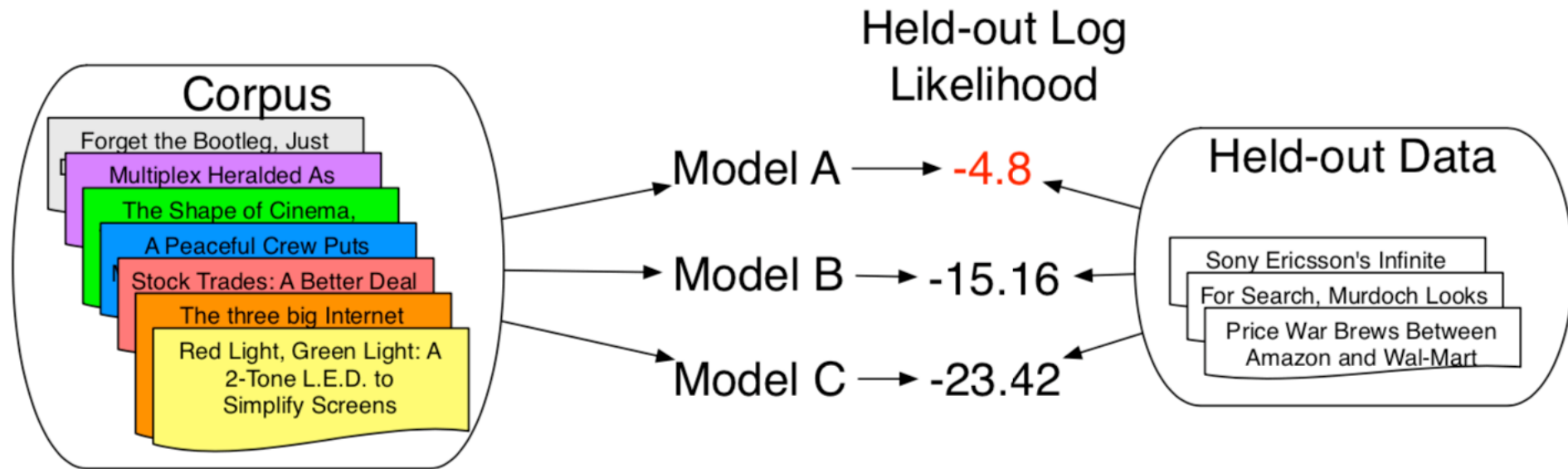
# Evaluating Topic Models

- Held-out log likelihood
- Word intrusion            *Jordan Boyd-Graber*

# Held-out log likelihood



# Held-out log likelihood



$$\begin{aligned} p(\underline{w} \mid \alpha, \beta) &= \int_{\theta, z} p(\theta, z, w \mid \alpha, \beta) \\ &= \prod_d \int_{\theta_d} p(\underline{\theta}_d \mid \alpha) \prod_n \sum_{z_{d,n}} p(\underline{z}_{d,n} \mid \theta_d) p(w_{d,n} \mid \beta, z_{d,n}) \end{aligned}$$

# Word Intrusion

## TOPIC 1

computer,  
technology,  
system,  
service, site,  
phone,  
internet,  
machine

## TOPIC 2

sell, sale,  
store, product,  
business,  
advertising,  
market,  
consumer

## TOPIC 3

play, film,  
movie, theater,  
production,  
star, director,  
stage



# Word Intrusion

1. Take the highest probability words from a topic

Original Topic

dog, cat, horse, pig, cow

# Word Intrusion

1. Take the highest probability words from a topic

Original Topic

dog, cat, horse, pig, cow

2. Take a high-probability word from another topic and add it

Topic with Intruder

dog, cat, **apple**, horse, pig, cow

# Word Intrusion

1. Take the highest probability words from a topic

Original Topic

dog, cat, horse, pig, cow

2. Take a high-probability word from another topic and add it

Topic with Intruder

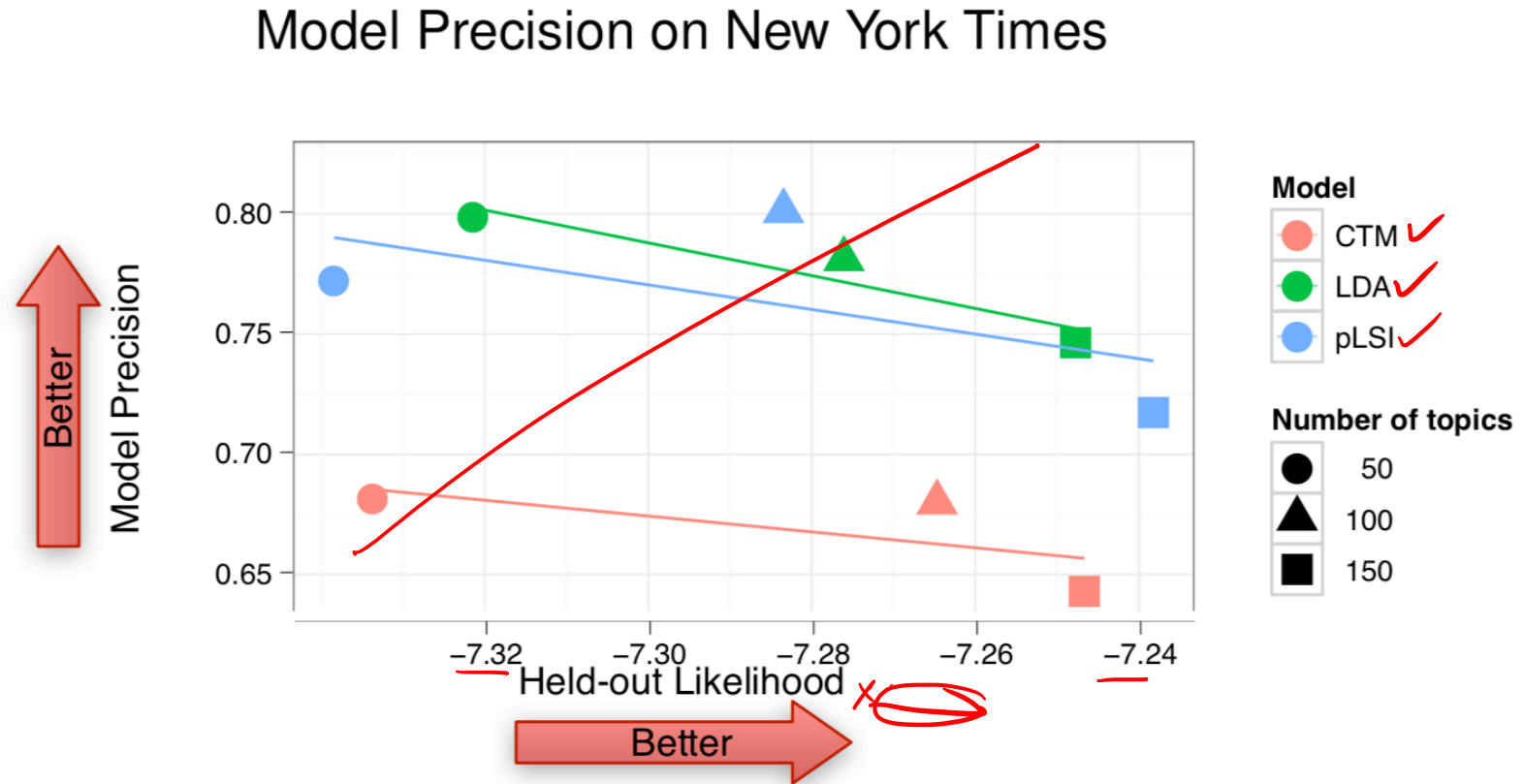
dog, cat, **apple**, horse, pig, cow

3. We ask users to find the word that doesn't belong

Hypothesis

If the topics are interpretable, users will consistently choose true intruder

# Interpretability and likelihood



within a model, higher likelihood  $\neq$  higher interpretability

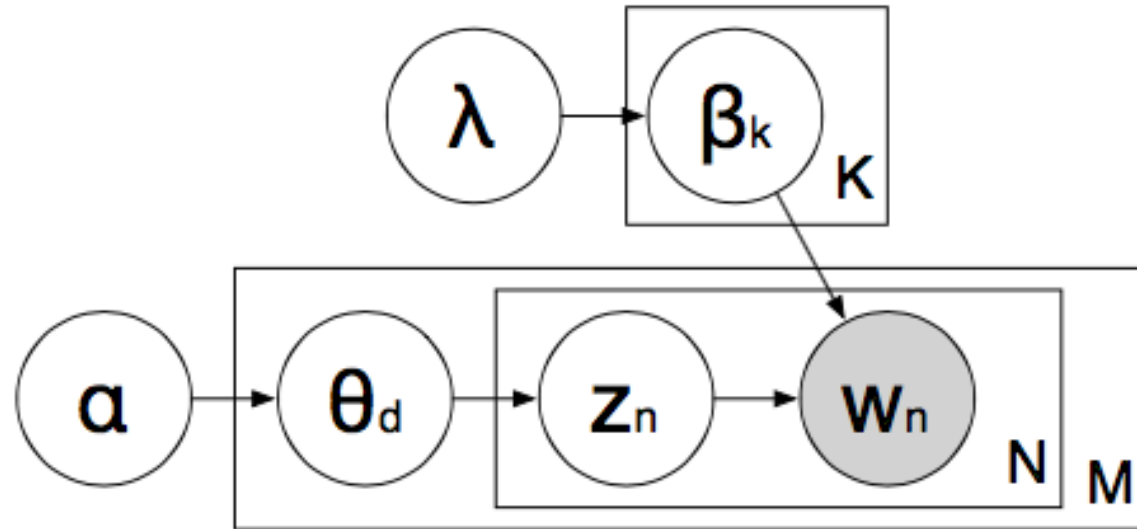
# Evaluation takeaway

- Measure what you care about
- If you care about prediction, likelihood is great
- If you care about a particular task, measure that

# Outline

- Evaluations of topic models
- Variational inference

# Joint distribution



$$\underline{p(\theta, z, w \mid \alpha, \beta)} = \prod_d p(\theta_d \mid \alpha) \prod_n p(z_{d,n} \mid \theta_d) p(w_{d,n} \mid \beta, z_{d,n})$$

# Joint distribution

$$p(\theta, z, w \mid \alpha, \beta) = \prod_d p(\theta_d \mid \alpha) \prod_n p(z_{d,n} \mid \theta_d) p(w_{d,n} \mid \beta, z_{d,n})$$

- $p(\theta_d \mid \alpha) = \frac{\Gamma(\sum_i \alpha_i)}{\prod_i \Gamma(\alpha_i)} \prod_k \theta_{d,k}^{\alpha_k - 1}$  (Dirichlet)
- $p(z_{d,n} \mid \theta_d) = \theta_{d,z_{d,n}}$  (Draw from Multinomial)
- $p(w_{d,n} \mid \beta, z_{d,n}) = \beta_{z_{d,n}, w_{d,n}}$  (Draw from Multinomial)



# Posterior distribution

Joint distribution:

$$p(\theta, z, w | \alpha, \beta) = \prod_d p(\theta_d | \alpha) \prod_n p(z_{d,n} | \theta_d) p(w_{d,n} | \beta, z_{d,n})$$

Posterior distribution:

$$p(\underline{\theta}, \underline{z} | \underline{w}, \alpha, \beta) = \frac{p(\theta, z, w | \alpha, \beta)}{p(\underline{w} | \alpha, \beta)}$$

$$\underline{p(w | \alpha, \beta)} = \int_{\theta, z} p(\theta, z, w | \alpha, \beta)$$

$$= \prod_d \int_{\theta_d} p(\theta_d | \alpha) \prod_n \sum_{z_{d,n}} p(z_{d,n} | \theta_d) p(w_{d,n} | \beta, z_{d,n})$$

# Variational inference

ELBO

Instead of estimating the posterior distribution directly, we use a distribution of simpler forms,  $q(\theta, z)$  to approximate  $P(\theta, z | w, \alpha, \beta)$ . We try to minimize the difference between  $p$  and  $q$

$$\ell(\theta, z) = \ell(\theta) + \ell(z)$$

$$\min \text{KL}(q || p) \equiv \mathbb{E}_q \left[ \log \frac{q(\theta, z)}{p(\theta, z | w)} \right] \Rightarrow \max \text{ELBO}$$

- If  $q$  and  $p$  are high, we're happy
- If  $q$  is high but  $p$  isn't, we pay a price
- If  $q$  is low, we don't care
- If  $\text{KL} = 0$ , then distribution are equal

$q$  is usually referred to as the variational distribution.

# KL divergence and evidence lower bound

$$KL(q||p) = E_q \log \frac{q(\theta, z)}{p(\theta, z|w)}$$

$$p(\theta, z|w) = \frac{p(\theta, z, w)}{p_w}$$

$$= E_q \log \frac{q(\theta, z) p_w}{p(\theta, z, w)}$$

$$= E_q \log q(\theta, z) - E_q \log p(\theta, z, w) + \underbrace{E_q \log p_w}_{= \sum_{\theta, z} q(\theta, z) \log p_w} = 1 \cdot \log p_w$$

$$= E_q \log q(\theta, z) - E_q \log p(\theta, z, w) + \log p_w \geq 0$$

$$\log p_w \geq \underbrace{E_q \log p(\theta, z, w)}_{\uparrow} - E_q \log q(\theta, z)$$

# KL divergence and evidence lower bound

- Conditional probability definition

$$p(z|x) = \frac{p(z, x)}{p(x)}$$

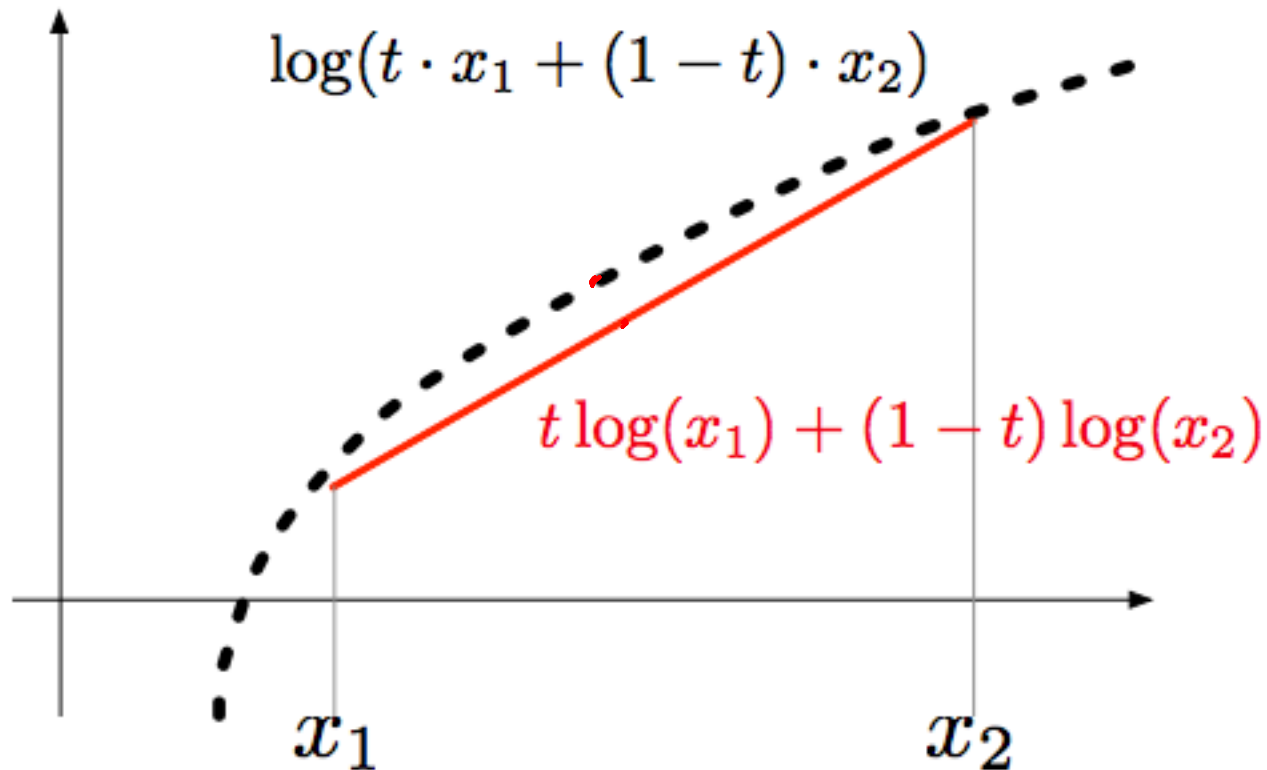
- Plug into KL divergence

$$\begin{aligned}\text{KL}(q(z) || p(z|x)) &= \mathbb{E}_q \left[ \log \frac{q(z)}{p(z|x)} \right] \\ &= \mathbb{E}_q [\log q(z)] - \mathbb{E}_q [\log p(z|x)] \\ &= \mathbb{E}_q [\log q(z)] - \mathbb{E}_q [\log p(z, x)] + \log p(x) \\ &= - (\mathbb{E}_q [\log p(z, x)] - \mathbb{E}_q [\log q(z)]) + \log p(x)\end{aligned}$$

- Negative of ELBO (plus **constant**); minimizing KL divergence is the same as maximizing ELBO

# A different way to get ELBO

- Jensen's inequality



When  $f$  is concave

$$\underbrace{f(\mathbb{E}[X]) \geq \mathbb{E}[f(X)]}$$

# Evidence Lower Bound

- Apply Jensen's inequality on log probability of data

$$\underline{\log p(x)} = \log \left[ \int_z p(x, z) \right]$$

$$= \log \int_z p(x, z) \frac{q(z)}{q(z)} dz$$

$$= \log \int_z q(z) \cdot \frac{p(x, z)}{q(z)} dz$$

$$\geq \int_z q(z) \log \frac{p(x, z)}{q(z)} dz$$

$$= \int_z q(z) \log p(x, z) dz - \int_z q(z) \log q(z) dz$$

$$= E_q \log p(x, z) - E_q \log q(z)$$

# Evidence Lower Bound

- Apply Jensen's inequality on log probability of data

$$\begin{aligned}\log p(x) &= \log \left[ \int_z p(x, z) \right] \\ &= \log \left[ \int_z p(x, z) \frac{q(z)}{q(z)} \right] \\ &= \log \left[ \mathbb{E}_q \left[ \frac{p(x, z)}{q(z)} \right] \right] \\ &\geq \mathbb{E}_q [\log p(x, z)] - \mathbb{E}_q [\log q(z)]\end{aligned}$$

Apply Jensen's equality and use log difference

# Variational inference

- Propose variational distribution  $q$   ~~$p_\theta(x)$~~
- Find ELBO (evidence lower bound) using  $q$
- Set derivatives to 0 and update variables



# Variational distribution for LDA

Joint distribution:

$$p(\theta, z, w | \alpha, \beta) = \prod_d p(\theta_d | \alpha) \prod_n p(z_{d,n} | \theta_d) p(w_{d,n} | \beta, z_{d,n})$$

---

Posterior distribution:

$$\underline{p(\theta, z | w, \alpha, \beta)} = \frac{p(\theta, z, w | \alpha, \beta)}{p(w | \alpha, \beta)} \quad \underline{\zeta(\theta, z)}$$

$$\begin{aligned} p(w | \alpha, \beta) &= \int_{\theta, z} p(\theta, z, w | \alpha, \beta) \\ &= \prod_d \int_{\theta_d} p(\theta_d | \alpha) \prod_n \sum_{z_{d,n}} p(z_{d,n} | \theta_d) p(w_{d,n} | \beta, z_{d,n}) \end{aligned}$$

# Variational distribution for LDA

$$q(\theta, z) = \prod_d q(\theta_d | \gamma_d) \prod_z q(z | \phi)$$

Handwritten notes in red ink:

- Below the first product:  $\frac{\prod_i \gamma_{di}}{\sum_i \gamma_{di}}$
- Below the second product:  $\prod_i \theta_{di}$
- To the right:  $q(z | \phi) = \phi_z$

- Variational document distribution over topics  $\gamma_d$ 
  - Vector of length  $K$  for each document
  - Non-negative
  - Doesn't sum to 1.0
- Variational token distribution over topic assignments  $\phi_{d,n}$ 
  - Vector of length  $K$  for every token
  - Non-negative, sums to 1.0

# Overall Algorithm

1. Randomly initialize variational parameters (can't be uniform)
2. For each iteration:
  - 2.1 For each document, update  $\underline{\gamma}$  and  $\underline{\phi}$  *variational document-topic*
  - 2.2 For corpus, update  $\underline{\beta}$  *topic-word*
  - 2.3 Compute  $\mathcal{L}$  for diagnostics
3. Return expectation of variational parameters for solution to latent variables

# Updates to Maximize ELBO

For each document  $d$ ,  $\gamma$  means  $\gamma_d$  here,  $\phi_n$  means  $\phi_{d,n}$  here:

$M \times N \times K$

$$\phi_{ni} \propto \beta_{iv} \exp \left( \Psi(\gamma_i) - \Psi \left( \sum_j \gamma_j \right) \right)$$

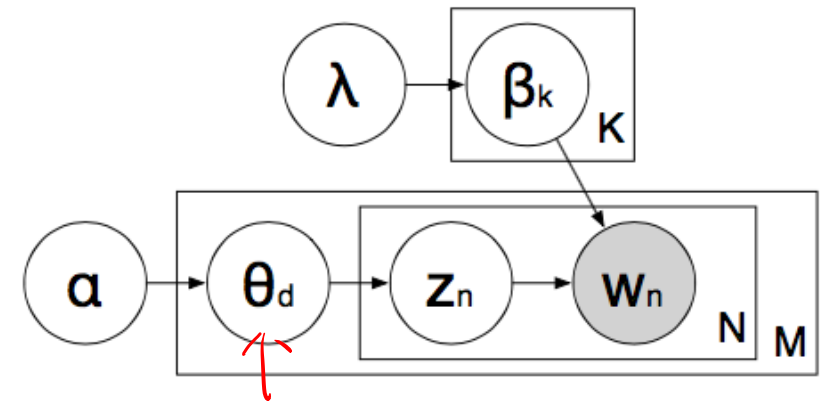
$K \times V$   $V = Wh$

$M \times K$

$$\gamma_i = \alpha_i + \sum_n \phi_{ni}$$

For the entire corpus,

$$\beta_{ij} \propto \sum_d \sum_n \phi_{dni} w_{dn}^j$$



## Example

- Three topics

$$\beta = \begin{matrix} \textcolor{red}{1} \\ \textcolor{red}{2} \\ \textcolor{red}{3} \end{matrix} \begin{bmatrix} \text{cat} & \text{dog} & \text{hamburger} & \text{iron} & \text{pig} \\ \textcolor{red}{.26} & \textcolor{red}{.185} & \textcolor{red}{.185} & \textcolor{red}{.185} & \textcolor{red}{.185} \\ \textcolor{red}{.185} & \textcolor{red}{.185} & \textcolor{red}{.26} & \textcolor{red}{.185} & \textcolor{red}{.185} \\ \textcolor{red}{.185} & \textcolor{red}{.185} & \textcolor{red}{.185} & \textcolor{red}{.26} & \textcolor{red}{.185} \end{bmatrix} = \textcolor{red}{I} \quad (4)$$

- Assume uniform  $\gamma$ : (2.0, 2.0, 2.0)
- Compute update for  $\phi$

$$\phi_{ni} \propto \textcolor{red}{\beta_{iv}} \textcolor{red}{\exp} \left( \textcolor{red}{\Psi(\gamma_i)} - \textcolor{red}{\Psi\left(\sum_j \gamma_j\right)} \right) \quad (5)$$

- For the first word (dog) in the document: **dog cat cat pig**

## Update $\phi$ for dog $\gamma=0$

$$\beta = \begin{bmatrix} \text{cat} & \text{dog} & \text{hamburger} & \text{iron} & \text{pig} \\ .26 & .185 & .185 & .185 & .185 \\ .185 & .185 & .26 & .185 & .185 \\ .185 & .185 & .185 & .26 & .185 \end{bmatrix} \begin{matrix} \phi_{ni} \propto \\ \beta_{iv} \exp \left( \Psi(\gamma_i) - \Psi \left( \sum_j \gamma_j \right) \right) \end{matrix}$$

- $\gamma = (2.000, 2.000, 2.000)$
- $\phi(0) \propto \underline{0.185} \times \underline{\exp(\Psi(2.000) - \Psi(2.000 + 2.000 + 2.000))} = 0.051$
- $\phi(1) \propto 0.185 \times \exp(\Psi(2.000) - \Psi(2.000 + 2.000 + 2.000)) = 0.051$
- $\phi(2) \propto 0.185 \times \exp(\Psi(2.000) - \Psi(2.000 + 2.000 + 2.000)) = 0.051$
- After normalization:  $\{0.333, 0.333, 0.333\}$

## Update $\phi$ for pig $n=1$

---

$$\beta = \begin{bmatrix} \text{cat} & \text{dog} & \text{hamburger} & \text{iron} & \text{pig} \\ .26 & .185 & .185 & .185 & .185 \\ .185 & .185 & .26 & .185 & .185 \\ .185 & .185 & .185 & .26 & .185 \end{bmatrix} \begin{matrix} \phi_{ni} \propto \\ \beta_{iv} \exp \left( \Psi(\gamma_i) - \Psi \left( \sum_j \gamma_j \right) \right) \end{matrix}$$

- $\gamma = (2.000, 2.000, 2.000)$
- $\phi(0) \propto 0.185 \times \exp(\Psi(2.000) - \Psi(2.000 + 2.000 + 2.000)) = 0.051$
- $\phi(1) \propto 0.185 \times \exp(\Psi(2.000) - \Psi(2.000 + 2.000 + 2.000)) = 0.051$
- $\phi(2) \propto 0.185 \times \exp(\Psi(2.000) - \Psi(2.000 + 2.000 + 2.000)) = 0.051$
- After normalization:  $\{0.333, 0.333, 0.333\}$

## Update $\phi$ for cat

$$\beta = \begin{matrix} \begin{matrix} 0 \\ 1 \\ 2 \end{matrix} \end{matrix} \begin{bmatrix} \text{cat} & \text{dog} & \text{hamburger} & \text{iron} & \text{pig} \\ \begin{matrix} \rightarrow \\ \downarrow \end{matrix} \begin{bmatrix} .26 \\ .185 \\ .185 \end{bmatrix} & \begin{bmatrix} .185 \\ .185 \\ .185 \end{bmatrix} & \begin{bmatrix} .185 \\ .26 \\ .185 \end{bmatrix} & \begin{bmatrix} .185 \\ .185 \\ .26 \end{bmatrix} & \begin{bmatrix} .185 \\ .185 \\ .185 \end{bmatrix} \end{bmatrix} \begin{matrix} \phi_{ni} \propto \\ \beta_{iv} \exp \left( \Psi(\gamma_i) - \Psi \left( \sum_j \gamma_j \right) \right) \\ \beta_{0, \text{cat}} \\ \beta_{1, \text{cat}} \\ \beta_{2, \text{cat}} \end{matrix}$$

- $\gamma = (2.000, 2.000, 2.000)$
- $\phi(0) \propto \underline{0.260} \times \exp(\underline{\Psi(2.000) - \Psi(2.000 + 2.000 + 2.000)}) = 0.072$
- $\phi(1) \propto 0.185 \times \exp(\Psi(2.000) - \Psi(2.000 + 2.000 + 2.000)) = 0.051$
- $\phi(2) \propto 0.185 \times \exp(\Psi(2.000) - \Psi(2.000 + 2.000 + 2.000)) = 0.051$
- After normalization: {0.413, 0.294, 0.294}



## Update $\gamma$

---

- Document: dog cat cat pig
- Update equation

$$\gamma_i = \alpha_i + \sum_n \phi_{ni} \quad (6)$$

- Assume  $\alpha = (.1, .1, .1)$

	$\phi_0$	$\phi_1$	$\phi_2$	
dog	.333	.333	.333	
cat	.413	.294	.294	x2
pig	.333	.333	.333	
$\alpha$	0.1	0.1	0.1	
sum	1.592	1.354	1.354	

- Note: **do not normalize!**

## Update $\beta$

---

- Count up all of the  $\phi$  across all documents
- For each topic, divide by total
- Corresponds to maximum likelihood of expected counts

# Recap

- Topic models: a neat way to model discrete count data
- Variational inference converts intractable optimization to maximizing ELBO