

**CS 589 Fall 2020**

**Latent semantic indexing**

**Latent Dirichlet allocation**

**Instructor: Susan Liu**

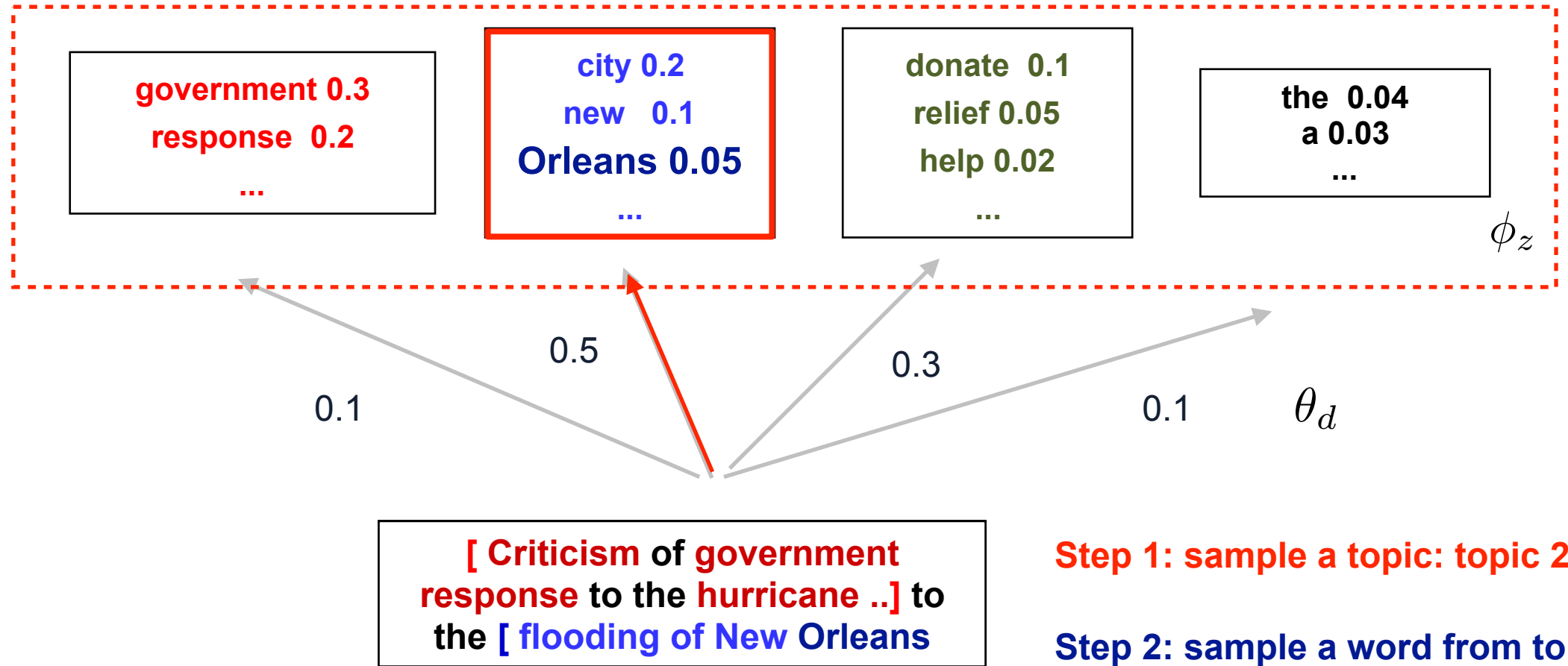
**TA: Huihui Liu**

**Stevens Institute of Technology**

*Most slides come from UIUC CS510 LDA lecture by Chase Geigle, ChengXiang Zhai*

# Review of PLSA

Phi  $T \times V$ ,  $V$ : vocabulary size  $\sim 50,000$ ,  $T$ : #topics,  $T=20$   
 theta  $D \times T$ ,  $D$ : #documents:  $\sim 10,000$ ,  $T$ : #topics,  $T=20$



# Today's lecture

- Latent semantic indexing
- Continue on topic model: Latent Dirichlet allocation
  - Bayesian inference of topic model
  - Variational inference for LDA
  - Gibbs sampling, Markov chain Monte-Carlo

# Vocabulary gap problem with vector space model

- Vector space model:

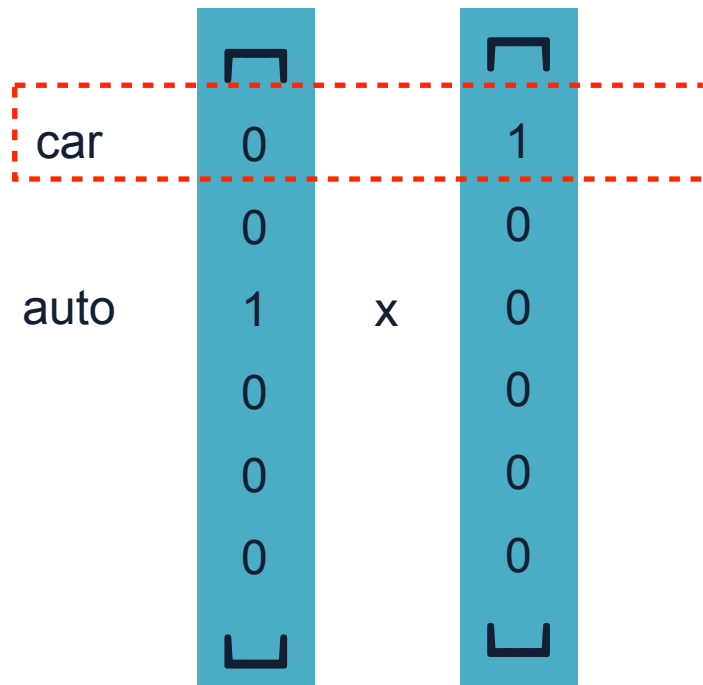
$$score(q, d) = \frac{q \cdot d}{\|q\| \cdot \|d\|}$$

- Challenges matching synonyms
  - e.g., auto vs. car
- Challenges matching polysemy
  - e.g., apple (fruit vs. company)

	doc1		doc2
	┌		┌
car	0		1
	0		0
auto	1	x	0
	0		0
	0		0
	└		└

# Low-dimensional, dense vector representation [Deerwester et al. 1998]

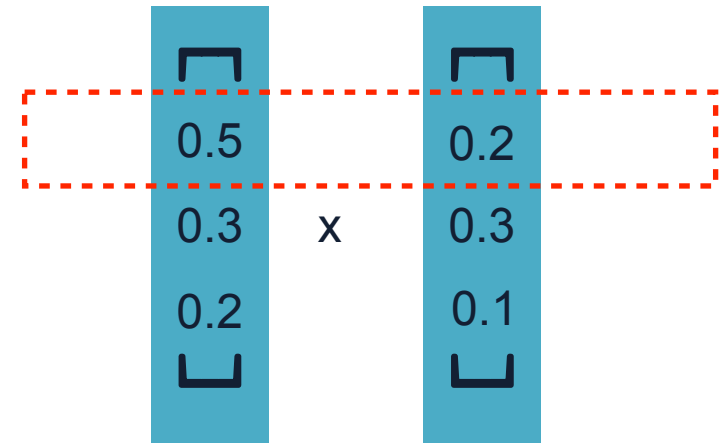
word indexing



latent semantic indexing



latent concept indexing

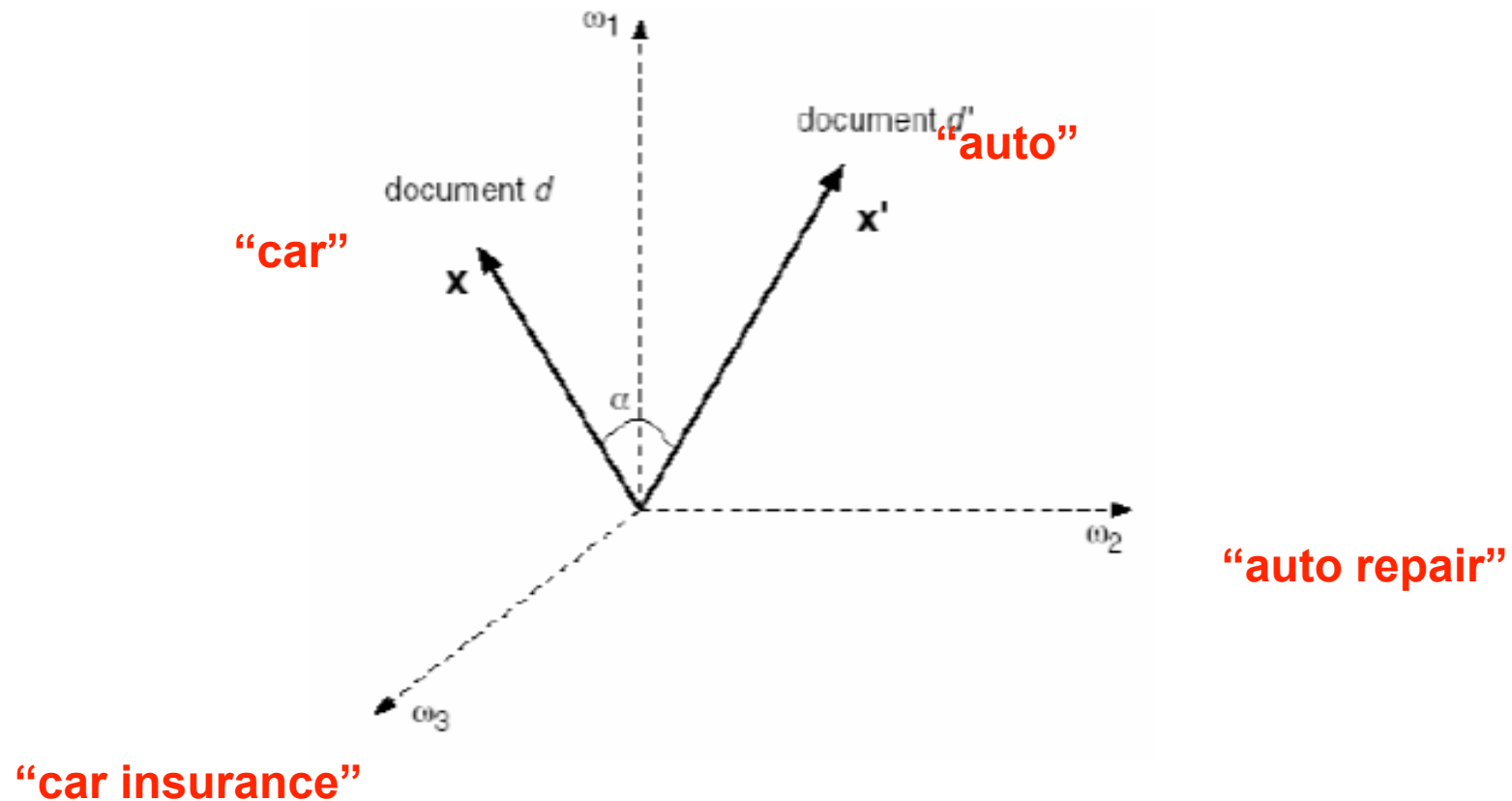


# Latent semantic indexing

- Uses statistically derived **conceptual indices** instead of individual word indexing for retrieval
- Assumes that there is some underlying or latent structure in word usage that is obscured by variability in word choice
- Key idea: instead of representing documents and queries as vectors in a **low dimensional space** of terms

# Low dimensional vector representation of words

- Axes are concepts, also called principal components (PCA)



# Singular value decomposition (matrix factorization)

- For a matrix  $A \in \mathbb{R}^{m \times n}$  of rank  $r$ , there exists a factorization (SVD) as follows:

$$A = U\Sigma V^T$$

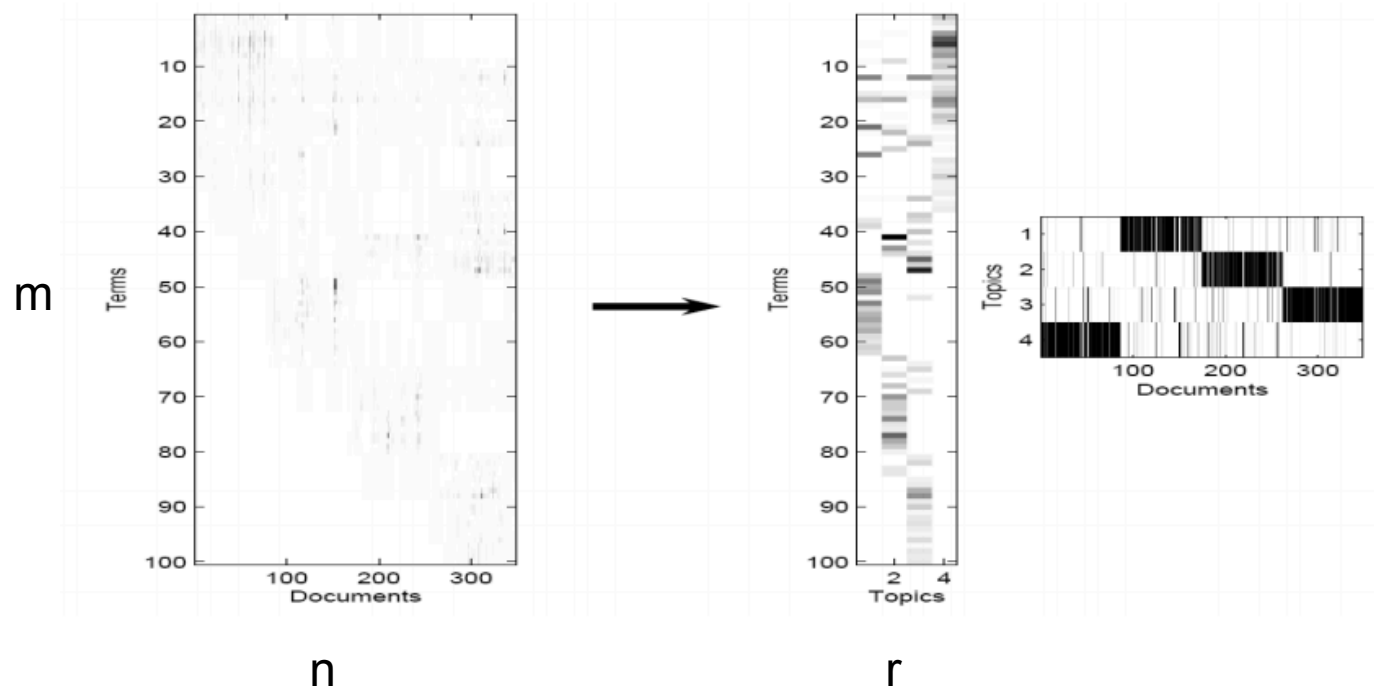
- $U$ : orthogonal eigen vectors of  $AA^T$
- $V$ : orthogonal eigen vectors of  $A^TA$
- $\Sigma$ : eigen values

$$\Sigma = \text{diag}(\sigma_1 \dots \sigma_r)$$



# Dimension reduction

- Map documents and queries to a low dimensional space
- Retrieval in this space may be superior to retrieval in the original space

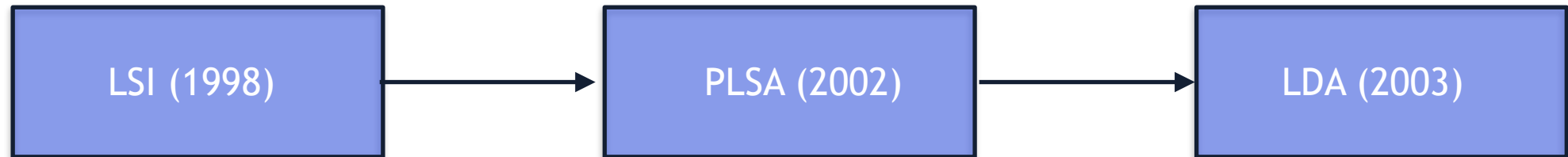


# What LSI can do

- LSI effectively does
  - Dimensionality reduction
  - Noise reduction
  - Exploitation of redundant data (linearity)
  - Correlation analysis and query expansion (with related words)
- Some of the individual effects can be achieved with simpler techniques (e.g. thesaurus construction). LSI does them together
- LSI handles **synonymy** well, not so much **polysemy** (vs word embedding)
- Challenge: SVD is complex to compute ( $O(n^3)$ ) – Needs to be updated as new documents are found/updated

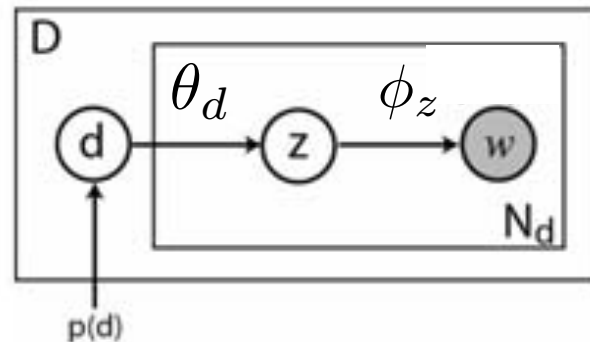
# Latent Dirichlet allocation [Blei et al. 03]

- A Bayesian topic model by considering the prior distribution of topics
  - That is, assuming the topic portions are randomly sampled from another distribution

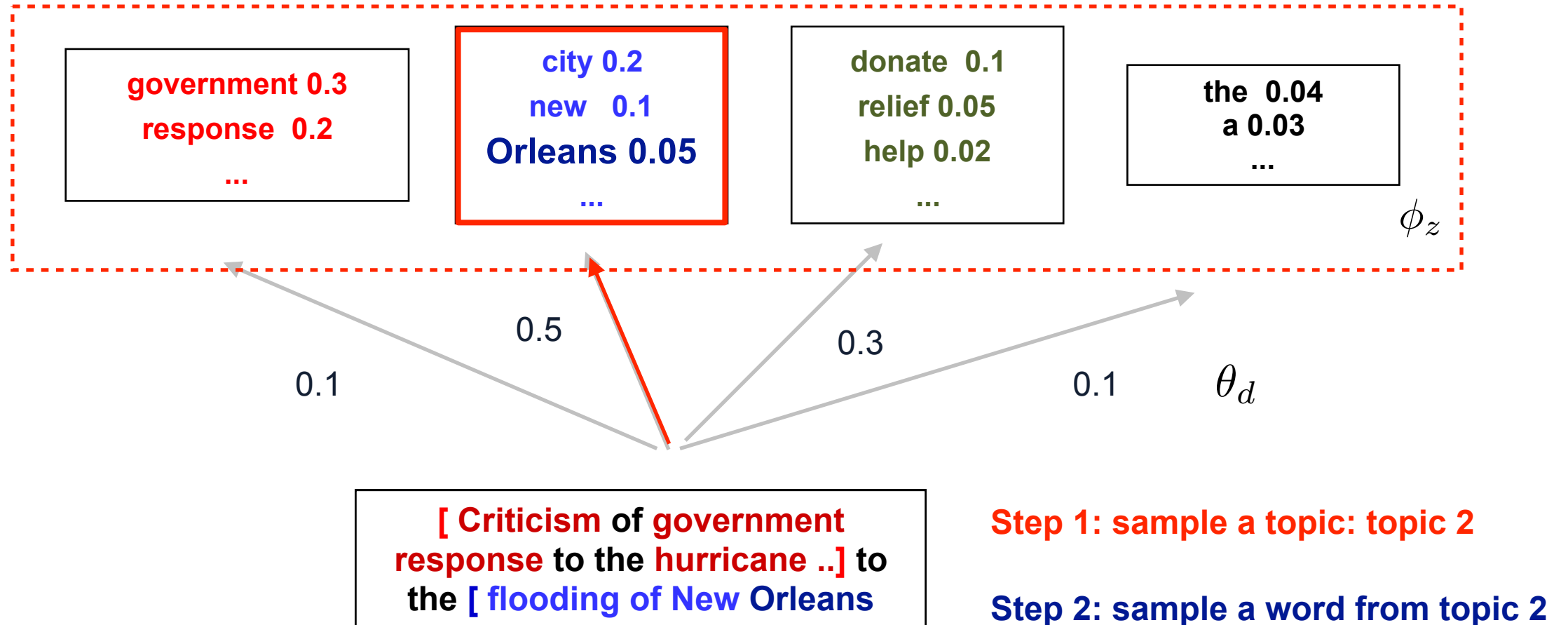


# Review of PLSA

- Documents are generated by sampling words from  $k$  latent topics
- For each document  $d$ :
  - For each token position  $i$ 
    - Choose a topic  $z \sim \text{Multinomial}(\theta_d)$  *is a fixed distribution*
    - Choose a term  $w \sim \text{Multinomial}(\phi_z)$

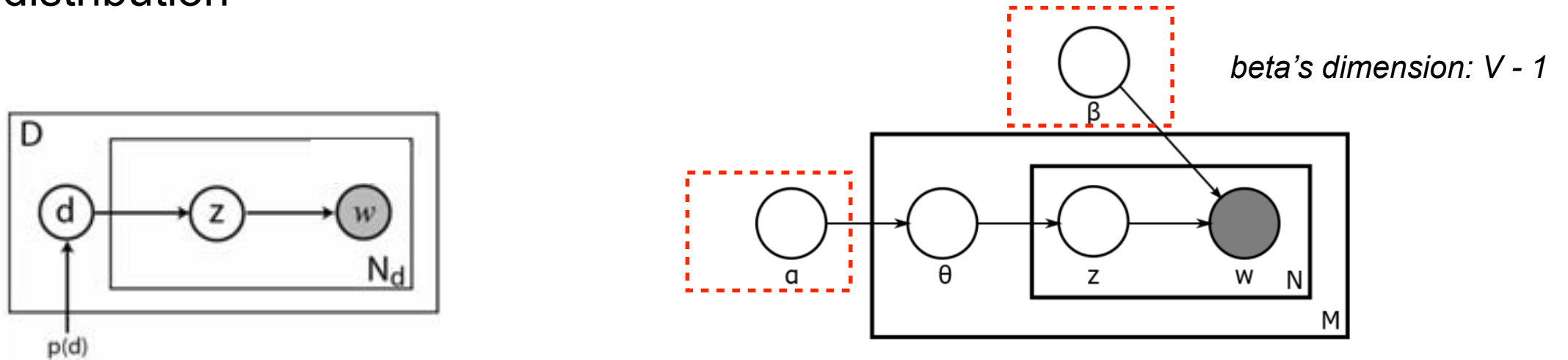


# Review of PLSA



# Latent Dirichlet allocation [Blei et al. 2003]

- A generative statistical topic model
- Adding a prior hyper parameter alpha to PLSA model
  - The document-topic probability is **randomly sampled** from the same prior distribution



*alpha's dimension:  $k - 1$*

# Notation table for LDA

- $d$  (or  $n$ ): index of document,  $d=1, \dots, D$
- $i$ : word inside document  $d$ ,  $i = 1, 2, \dots, N_d$
- $z$ : topic index,  $z = 1, 2, \dots, k$
- $w$  or  $v$ : word index
- $w_{\{d,i\}}$ : the  $i$ -th word in document  $d$
- $\theta_{\{d,z\}}$ : probability of the  $z$ -th topic in document  $d$
- $\phi_{\{z, w\}}$ : probability of the  $w$ -th word in topic  $z$
- $\alpha$ : prior distribution of  $\theta$ , dimension =  $k - 1$
- $\beta$ : prior distribution of  $\phi$ , dimension =  $V - 1$

# Bayes' rules

Chain rule: **joint distribution**

$$P(A, B) = P(A \cap B) = P(A|B)P(B) = P(B|A)P(A)$$

Bayes' rule:

**posterior**   **likelihood**   **prior**

$$P(A|B) = \frac{P(B|A)P(A)}{P(B)} = \left[ \frac{P(B|A)}{\sum_{X \in \{A, \bar{A}\}} P(B|X)P(X)} \right] P(A)$$

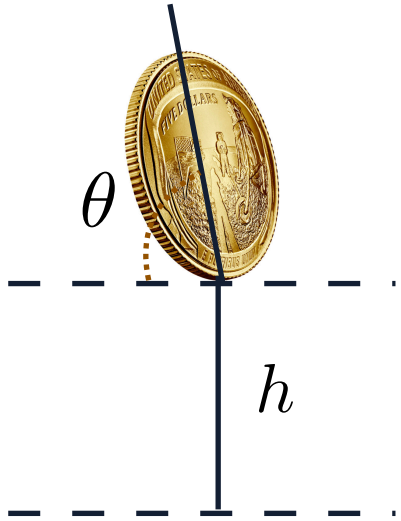
$$P(A|B) \propto P(B|A)P(A)$$
$$\sum_A P(A|B) = 1$$

**skipping estimating  $P(B)$**

**trick for estimating the posterior**



# Frequentist vs. Bayesian inference



- Frequentist treats parameters as **fixed**

sequence = 0, 1, 0, 0,  
1, 1, 0, 1, 0, 1, 0, 0

$$\mu = \frac{\#ups}{\#ups + \#downs}$$

$$x_i \sim \mu$$

- Bayesians treats parameters as **random**
  - Consider the force that causes the coin to be biased
  - The forces are explained a prior distribution:

$$x_i \sim \mu$$

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

# Maximum a Posterior (MAP) estimation

- Maximum likelihood estimation:

$$\theta_{MLE} = \arg \max_{\theta} p(x_1, \dots, x_n | \theta)$$

- Choose the parameter with the maximum posterior probability

$$\begin{aligned} \theta_{MAP} &= \arg \max_{\theta} p(\theta | x_1, \dots, x_n) \\ &= \arg \max_{\theta} p(x_1, \dots, x_n | \theta) p(\theta) \end{aligned}$$

**How many parameters in PLSA vs. LDA? PLSA:  $d*(k-1) + k*(V-1)$ , LDA:  $d*(k-1)+k*(V-1) + (k - 1) + (V - 1)$**

# PLSA -> LDA

- PLSA:

$$p_d(w | \{\theta_d\}, \{\phi_z\}) = \sum_{z=1}^k \theta_{d,z} \phi_{z,w}$$

$$\log p(D | \{\phi_z\}, \{\theta_d\}) = \sum_{d \in D} \sum_{w \in V} c(w, d) \log \left[ \sum_{z=1}^k \theta_{d,z} \theta_{z,w} \right]$$

- LDA:

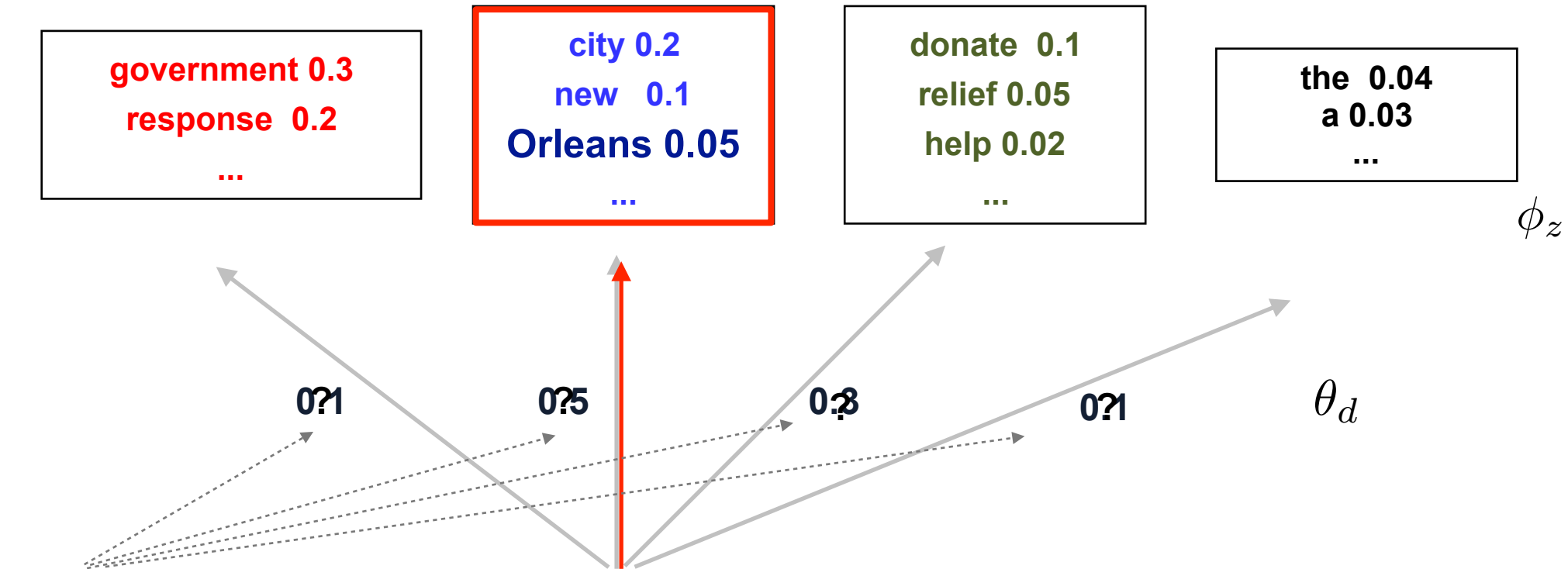
$$p_d(w | \{\theta_d\}, \{\phi_z\}) = \sum_{z=1}^k \theta_{d,z} \phi_{z,w}$$

$$\log p(d | \alpha, \{\phi_z\}) = \int \sum_{w \in V} c(w, d) \log \left[ \sum_{z=1}^k \theta_{d,z} \phi_{z,w} \right] p(\theta_d | \alpha) d\theta_d$$

**marginalized  
probability:**

$$\log p(D | \alpha, \beta) = \int \sum_{d \in D} \log p(d | \alpha, \{\phi_z\}) \prod_{z=1}^k p(\phi_z | \beta) d\phi_1 \dots d\phi_k$$

# PLSA -> LDA



[ Criticism of government response to the hurricane ..] to the [ flooding of New Orleans

Step 0: sample topic distribution

Step 1: sample a topic: topic 2

Step 2: sample a word from topic 2

# Solving Maximum a Posteriori inference for LDA

- Maximum likelihood estimation:

$$\begin{aligned} \mathcal{L} = \log p(\mathcal{W} \mid \mathbf{R}, \Phi, \Theta) &= \sum_d^D \sum_{d_i}^{N_d} \sum_z^T R_{(w_{d_i}, z)} (\log \phi_{(z, w_{d_i})} + \log \theta_{(d, z)}) \\ &+ \left[ \sum_{d=1}^D \lambda_d \left( 1 - \sum_{z=1}^T \theta_{(d, z)} \right) + \sum_{z=1}^T \sigma_k \left( 1 - \sum_{w=1}^V \phi_{z, w} \right) \right] \end{aligned}$$

**M step:**

$$\begin{aligned} \theta_{d, z} &\propto \sum_{v=1}^V R_{d, v, z} \text{count}(v, d) \\ \phi_{z, v} &\propto \sum_{d=1}^D [z_{d, v} == z] \text{count}(v, d) \end{aligned}$$

**E step:**

$$R_{(w_{d_i}, z)} \propto \phi_{z, v} \theta_{d, z}$$

# Solving Maximum a Posteriori inference for LDA

- Exact inference is intractable:

$$p(Z, \Phi, \Theta | D, \alpha, \beta) = \frac{p(D, Z, \Phi, \Theta | \alpha, \beta)}{p(D | \alpha, \beta)}$$

$$\log p(D | \alpha, \beta) = \int \sum_{d \in D} \log p(d | \alpha, \{\phi_z\}) \prod_{z=1}^k p(\phi_z | \beta) d\phi_1 \dots d\phi_k$$

- Equation (1) is computationally intractable due to the coupling of beta and phi in the denominator
- Question: **why do we need to compute the denominator?**

# Variational inference

- Key idea: use a **surrogate distribution** to approximate the posterior distribution of latent variables
  - Surrogate distribution is simpler to estimate than the true posterior
- Goal: finding the “best” surrogate distribution from a certain parametric family by **minimizing** the KL-divergence between the surrogate function (Q) to the true posterior (P)
- Typical surrogate distributions
  - Mean-field approximation [Blei et al. 03]
  - Expectation propagation [Minka et al. 02]
  - Collapsed variational Bayes [Teh et al. 07]

## Evidence lower bound (ELBO)

- Given that  $p(Z|D) = \frac{p(D, Z)}{p(D)}$ , we want to minimize the KL divergence between  $Q(Z)$  and  $p(Z|D)$ :

$$\begin{aligned} D_{\text{KL}}(q||p) &= \int_{\mathcal{Z}} q(Z) \left[ \log \frac{q(Z)}{p(Z, D)} + \log p(D) \right] \\ &= \int_{\mathcal{Z}} q(Z) [\log q(Z) - \log p(Z, D)] + \int_{\mathcal{Z}} q(Z) [\log p(D)] \\ &= \int_{\mathcal{Z}} q(Z) [\log q(Z) - \log p(Z, D)] + \log p(D) \end{aligned}$$

$$\Rightarrow \log p(D) = D_{\text{KL}}(q||p) - \mathbb{E}_q[\log q(Z) - \log p(Z, D)] = D_{\text{KL}}(q||p) + \mathcal{L}(q)$$



# Evidence lower bound (ELBO)

<https://stats.stackexchange.com/questions/205506/why-do-we-use-the-mean-field-approximation-for-variational-bayes>  
<https://www.cs.colorado.edu/~jbg/>

- $p(D)$  does not rely on the latent variable  $Z$ :

$$\Rightarrow \log p(D) = D_{\text{KL}}(q||p) - \mathbb{E}_q[\log q(Z) - \log p(Z, D)] = D_{\text{KL}}(q||p) + \mathcal{L}(q)$$

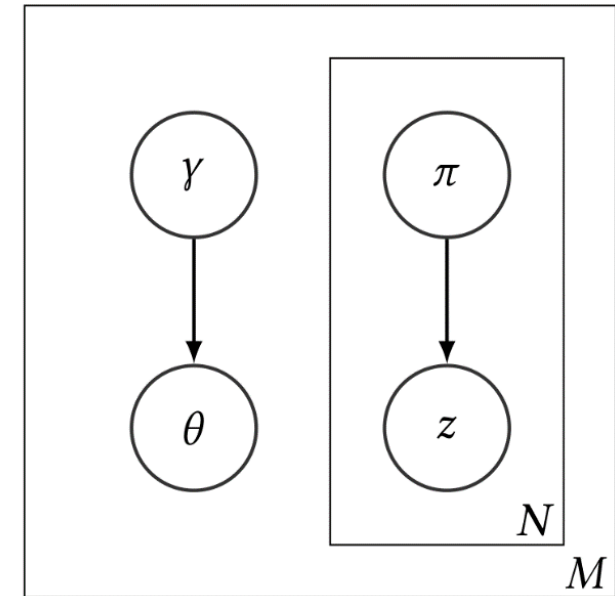
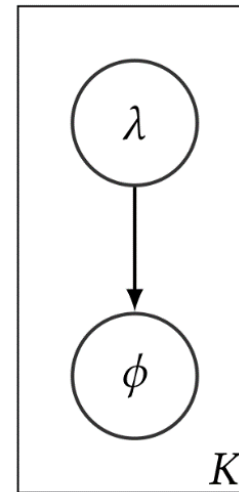
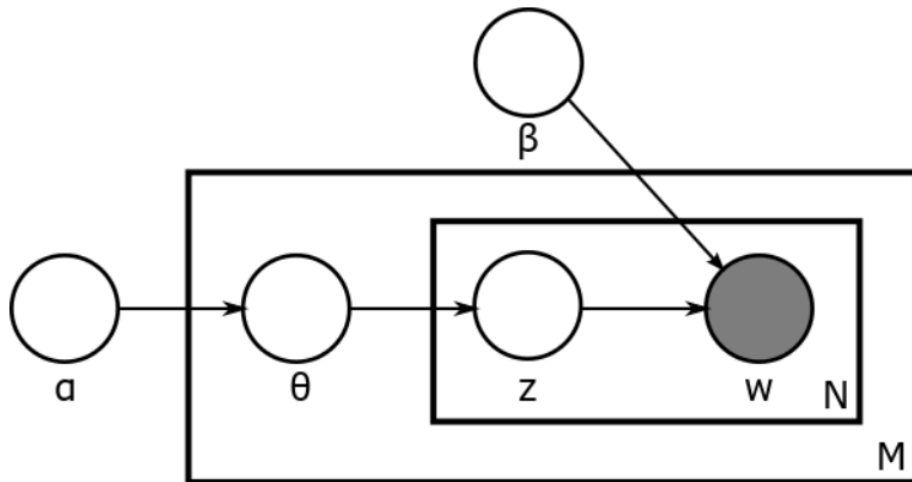
- Therefore minimizing KL divergence = maximizing  $L(q)$ , which we call evidence lower bound, as KL divergence is non-negative
- The lower bound is more tractable to optimize than the evidence:

$$L(q) = \int_Z q(Z) [\log q(Z) - \log p(Z, D)]$$

# LDA variational inference

- Assumption:  $q$  belongs to a family of distribution which can be factorized

$$q(\theta, \mathbf{z} \mid \gamma, \phi) = q(\theta \mid \gamma) \prod_{n=1}^N q(z_n \mid \phi_n)$$



# LDA variational inference

$$L(q) = \int_Z q(Z) [\log q(Z) - \log p(Z, D)]$$

$$p(Z, \Phi, \Theta | D, \alpha, \beta) = \frac{p(D, Z, \Phi, \Theta | \alpha, \beta)}{p(D | \alpha, \beta)}$$

- ELBO of LDA under the factorization assumption:

$$L(\gamma, \phi; \alpha, \beta) = \mathbb{E}_q[\log p(\theta | \alpha)] + \mathbb{E}_q[\log p(\mathbf{z} | \theta)] + \mathbb{E}_q[\log p(\mathbf{w} | \mathbf{z}, \beta)] \\ - \mathbb{E}_q[\log q(\theta)] - \mathbb{E}_q[\log q(\mathbf{z})]$$

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

← no  $\theta$ , why?

$$\mathbb{E}_q[\log(\theta_i) | \gamma] = \Psi(\gamma_i) - \Psi \left( \sum_{j=1}^k \gamma_j \right)$$

# LDA variational inference

- LDA variational inference algorithm: optimizing  $\gamma$ ,  $\phi$ ,  $\beta$
1. Randomly initialize variational parameters (can't be uniform)
  2. For each iteration:
    1. For each document, update  $\gamma$  and  $\phi$
    2. For corpus, update  $\beta$
    3. Compute L for diagnostics
  3. Return **expectation of variational parameters** for solution to latent variables

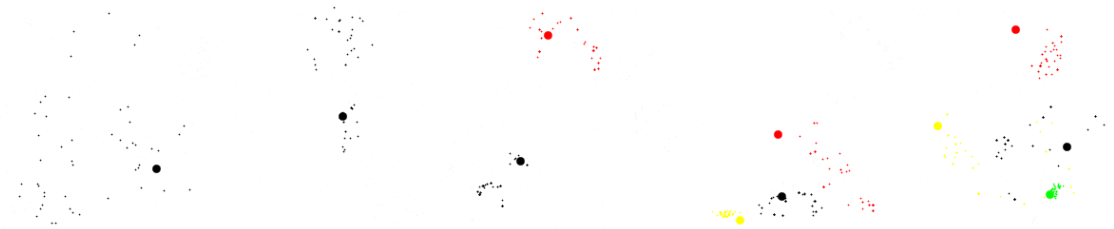
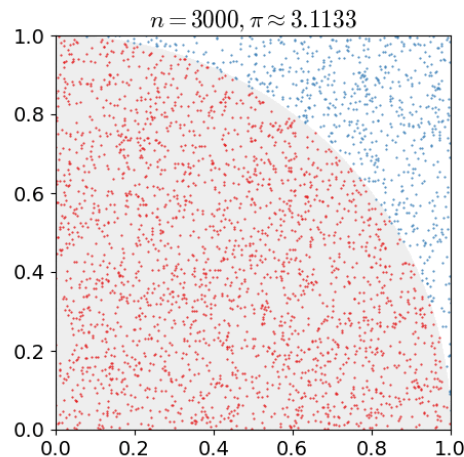
$$\gamma_i = \alpha_i + \sum_n \phi_{ni}$$
$$\phi_{ni} \propto \beta_{iv} \exp \left( \psi(\gamma_i) - \Psi \left( \sum_j \gamma_j \right) \right)$$
$$\beta_{ij} \propto \sum_d \sum_n \phi_{dni} w_{dn}^j$$
$$E_q [\log(\theta_i) | \gamma] = \Psi(\gamma_i) - \Psi \left( \sum_{j=1}^k \gamma_j \right)$$

# Pros and Cons of Variational inference

- Pros
  - Deterministic algorithm - easy to tell when converged
  - Embarrassingly parallelizable
- Cons:
  - Speed: make many calls to transcendental functions (no close form solution)
  - Quality is questionable due to the factorization assumption
  - Memory usage: requires  $O(MNK)$  to store the per-token variational distributions

# Monte Carlo for LDA inference

- Monte Carlo
  - A broad class of computational algorithms that rely on repeated random sampling to obtain numerical results
- Markov chain Monte Carlo (MCMC)
  - A Monte Carlo method with the Markov assumption, i.e., the next sample only rely on the previous state



# Conjugate prior distribution

- That is, the **posterior distributions** are in the same probability distribution family as the prior probability distribution
- The prior is called the **conjugate prior** of the likelihood probability
- Conjugate prior makes denominator easy to compute because you can integrate out theta
- e.g., Beta distribution is the conjugate prior of Bernoulli distribution

**Bernoulli:**  $p(x | \theta) = \theta^x (1 - \theta)^{1-x}.$

**Beta:**  $p(\theta | \alpha) = K(\alpha) \theta^{\alpha_1 - 1} (1 - \theta)^{\alpha_2 - 1}.$

**Demominator:** 
$$K(\alpha) = \left( \int \theta^{\alpha_1 - 1} (1 - \theta)^{\alpha_2 - 1} d\theta \right)^{-1}$$
$$= \frac{\Gamma(\alpha_1 + \alpha_2)}{\Gamma(\alpha_1) \Gamma(\alpha_2)}$$

# Collapsed variational inference [Teh et al. 07]

- If your priors are conjugate, they can be integrated out!
- Let  $n_{dzv}$  be the number of times that word  $v$  in document  $d$  has topic  $z$
- Let  $n_{.zv}$  be the number of times word  $v$  in any document has topic  $z$

$$p(\mathbf{z}, \mathbf{x} \mid \alpha, \beta) = \prod_d \frac{\Gamma(k\alpha)}{\Gamma(k\alpha + n_{d..})} \prod_z \frac{\Gamma(\alpha + n_{dz.})}{\Gamma(\alpha)} \prod_z \frac{\Gamma(V\beta)}{\Gamma(V\beta + n_{.z.})} \prod_v \frac{\Gamma(\beta + n_{.zv})}{\Gamma(\beta)}$$

**No theta or phi in the above formula**



# Collapsed variational inference [Teh et al. 07]

$$p(\mathbf{z}, \mathbf{x} \mid \alpha, \beta) = \prod_d \frac{\Gamma(k\alpha)}{\Gamma(k\alpha + n_{d..})} \prod_z \frac{\Gamma(\alpha + n_{dz.})}{\Gamma(\alpha)} \prod_z \frac{\Gamma(V\beta)}{\Gamma(V\beta + n_{.z.})} \prod_v \frac{\Gamma(\beta + n_{.zv})}{\Gamma(\beta)}$$

- How do we get back phi and theta? MAP estimates from  $\sigma$  and  $\delta$

$$\theta_{dz} = \frac{\alpha + n_{dk.}}{K\alpha + n_{d..}} \quad \phi_{zv} = \frac{\beta + n_{.zv}}{W\beta + n_{.z.}}$$

- How do we get  $n_{dzv}$  ?
  - Gibbs sampling: sample values of Z, count from those examples

# Collapsed Gibbs sampling [Griffths et al. 04; Teh et al. 07]

- **Key idea:** construct a well-behaved Markov chain such that
  1. the states of the chain represent an assignment of  $\mathbf{Z}$
  2. state transitions occur between states that differ in only one  $z_{\{d,i\}}$
  3. transition probabilities are based on the **full conditional**:

4.

$$p(z_{di} = k \mid \mathbf{z}^{-di}, \mathbf{x}, \alpha, \beta) = \frac{(\alpha + n_{dz^{\cdot}}^{-di}) (\beta + n_{\cdot zw_{di}}^{-di}) (W\beta + n_{\cdot z^{\cdot}}^{-di})^{-1}}{\sum_{z'=1}^k (\alpha + n_{dz'^{\cdot}}^{-di}) (\beta + n_{\cdot z'w_{di}}^{-di}) (W\beta + n_{\cdot z'^{\cdot}}^{-di})^{-1}}$$

5. where  $\mathbf{z}^{\{di\}}$  is the set of assignments for  $\mathbf{Z}$  without the  $i$ -th word in document  $d$

# Collapsed Gibbs sampling [Griffths et al. 04]

$i$	$w_i$	$d_i$	iteration	
			1	2
1	MATHEMATICS	1	2	?
2	KNOWLEDGE	1	2	
3	RESEARCH	1	1	
4	WORK	1	2	
5	MATHEMATICS	1	1	
6	RESEARCH	1	2	
7	WORK	1	2	
8	SCIENTIFIC	1	1	
9	MATHEMATICS	1	2	
10	WORK	1	1	
11	SCIENTIFIC	2	1	
12	KNOWLEDGE	2	1	
⋮	⋮	⋮	⋮	
⋮	⋮	⋮	⋮	
⋮	⋮	⋮	⋮	
50	JOY	5	2	

$$p(z_{di} = k \mid \mathbf{z}^{-di}, \mathbf{x}, \alpha, \beta) = \frac{(\alpha + n_{dz \cdot}^{-di}) (\beta + n_{z w_{di}}^{-di}) (W\beta + n_{z \cdot}^{-di})^{-1}}{\sum_{z'=1}^k (\alpha + n_{dz' \cdot}^{-di}) (\beta + n_{z' w_{di}}^{-di}) (W\beta + n_{z' \cdot}^{-di})^{-1}}$$

# Collapsed Gibbs sampling [Griffths et al. 04]

$i$	$w_i$	$d_i$	iteration	
			1	2
1	MATHEMATICS	1	2	2
2	KNOWLEDGE	1	2	?
3	RESEARCH	1	1	
4	WORK	1	2	
5	MATHEMATICS	1	1	
6	RESEARCH	1	2	
7	WORK	1	2	
8	SCIENTIFIC	1	1	
9	MATHEMATICS	1	2	
10	WORK	1	1	
11	SCIENTIFIC	2	1	
12	KNOWLEDGE	2	1	
⋮	⋮	⋮	⋮	
⋮	⋮	⋮	⋮	
⋮	⋮	⋮	⋮	
50	JOY	5	2	

$$p(z_{di} = k \mid \mathbf{z}^{-di}, \mathbf{x}, \alpha, \beta) = \frac{(\alpha + n_{dz}^{-di}) (\beta + n_{z'w_{di}}^{-di}) (W\beta + n_{z'}^{-di})^{-1}}{\sum_{z'=1}^k (\alpha + n_{dz'}^{-di}) (\beta + n_{z'w_{di}}^{-di}) (W\beta + n_{z'}^{-di})^{-1}}$$

# Collapsed Gibbs sampling [Griffths et al. 04]

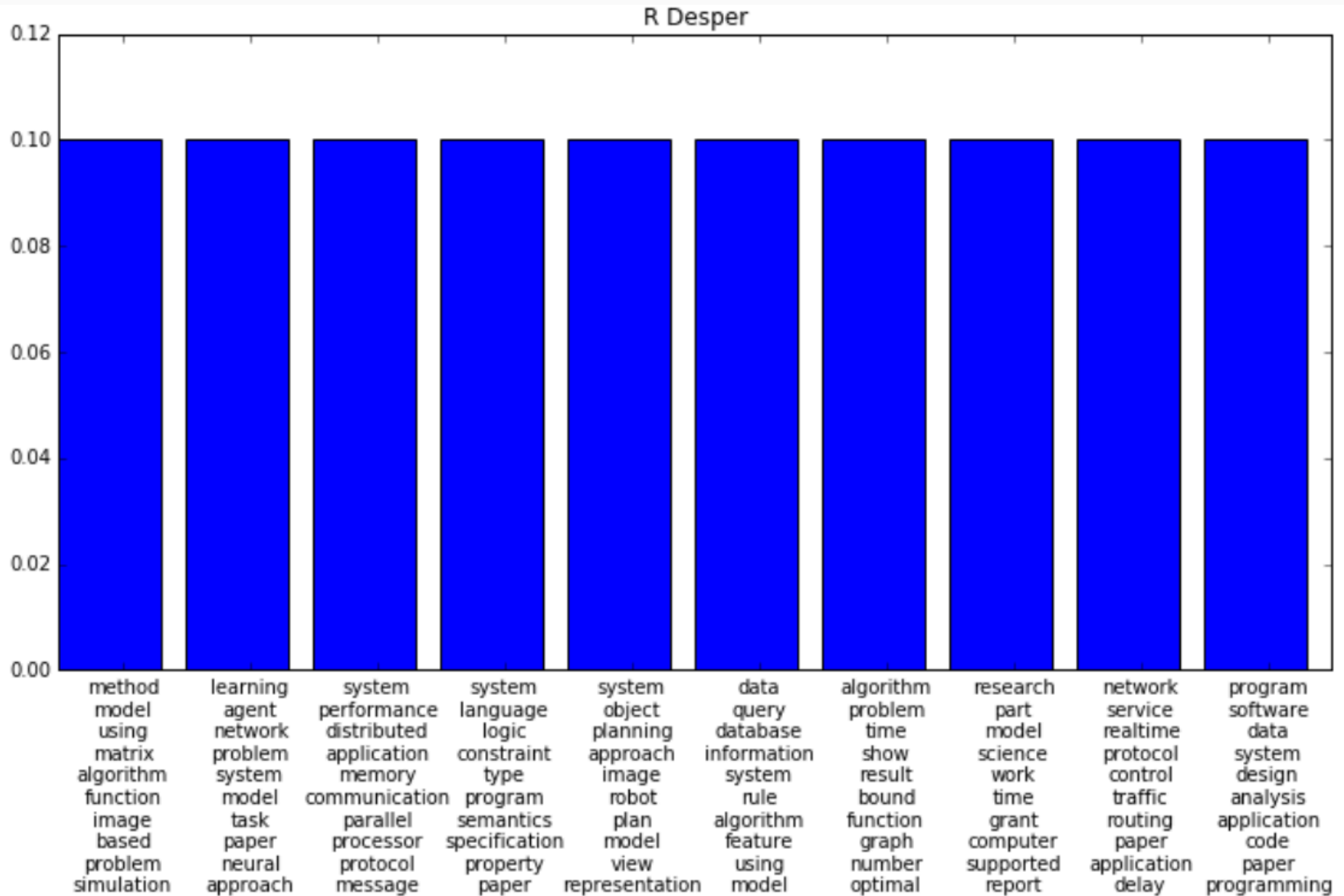
$i$	$w_i$	$d_i$	iteration	
			1	2
1	MATHEMATICS	1	2	2
2	KNOWLEDGE	1	2	1
3	RESEARCH	1	1	?
4	WORK	1	2	
5	MATHEMATICS	1	1	
6	RESEARCH	1	2	
7	WORK	1	2	
8	SCIENTIFIC	1	1	
9	MATHEMATICS	1	2	
10	WORK	1	1	
11	SCIENTIFIC	2	1	
12	KNOWLEDGE	2	1	
⋮	⋮	⋮	⋮	
⋮	⋮	⋮	⋮	
⋮	⋮	⋮	⋮	
50	JOY	5	2	

$$p(z_{di} = k \mid \mathbf{z}^{-di}, \mathbf{x}, \alpha, \beta) = \frac{(\alpha + n_{dz \cdot}^{-di}) (\beta + n_{z'w_{di}}^{-di}) (W\beta + n_{z \cdot}^{-di})^{-1}}{\sum_{z'=1}^k (\alpha + n_{dz' \cdot}^{-di}) (\beta + n_{z'w_{di}}^{-di}) (W\beta + n_{z' \cdot}^{-di})^{-1}}$$

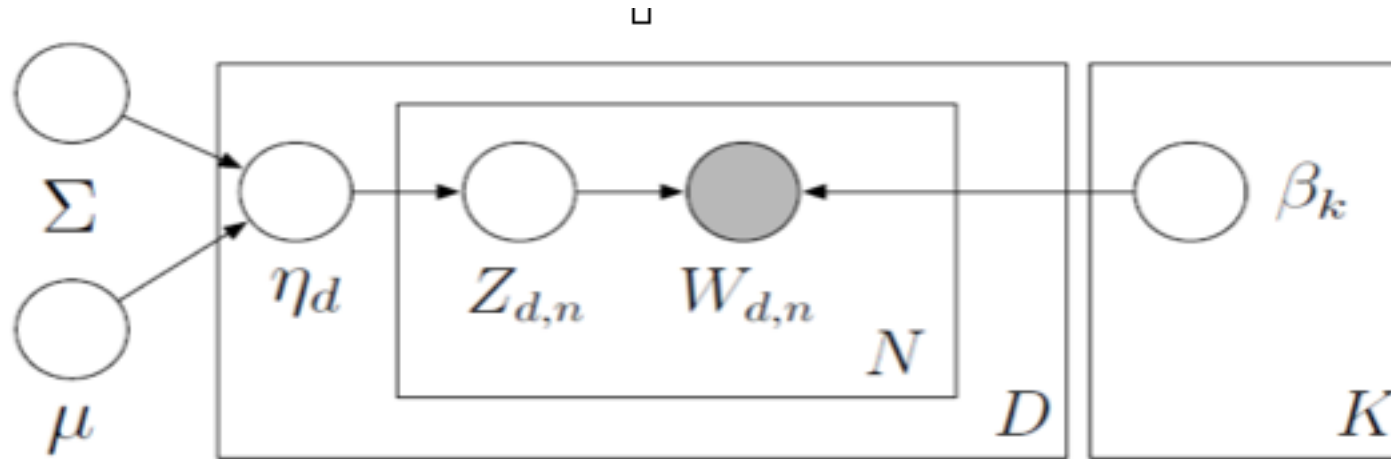
# Pros/Cons of Gibbs sampling

- Pros:
  - Ease of implementation
  - Fast iterations (no transcendental functions), fast convergence (at least relative to a full Gibbs sampler)
  - Low memory usage (only require  $O(MN)$  storage for the current values of  $z_{j,t}$ )
- Cons:
  - No obvious parallelization strategy (each iteration depends on previous)
  - Can be difficult to assess convergence

# LDA results



# Correlated topic models (CTM) [Blei et al. 2005]

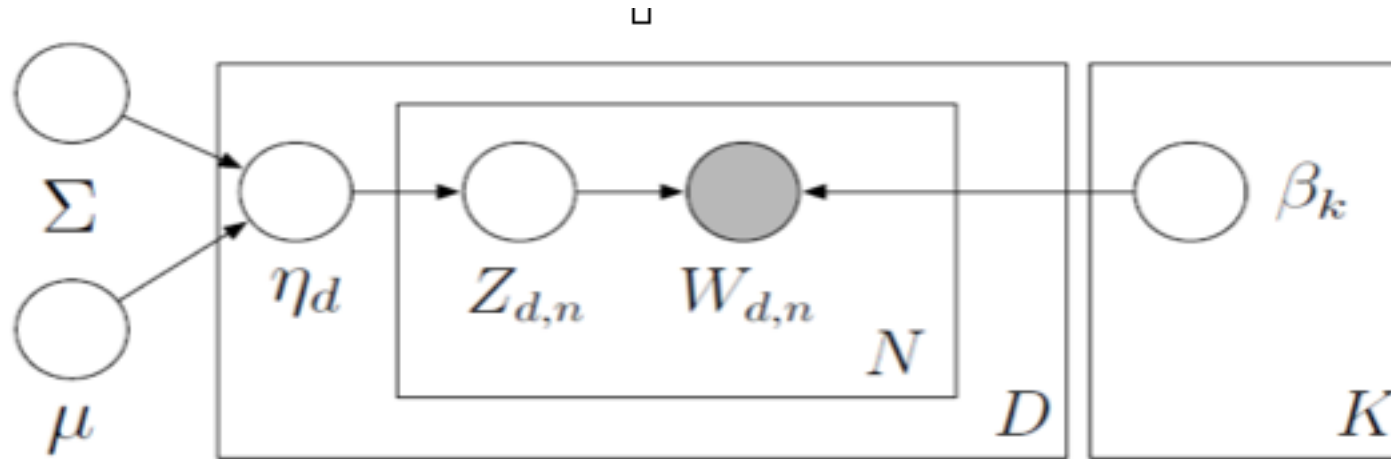


- Draw topics from a logistic normal, where topic occurrence can exhibit correlations



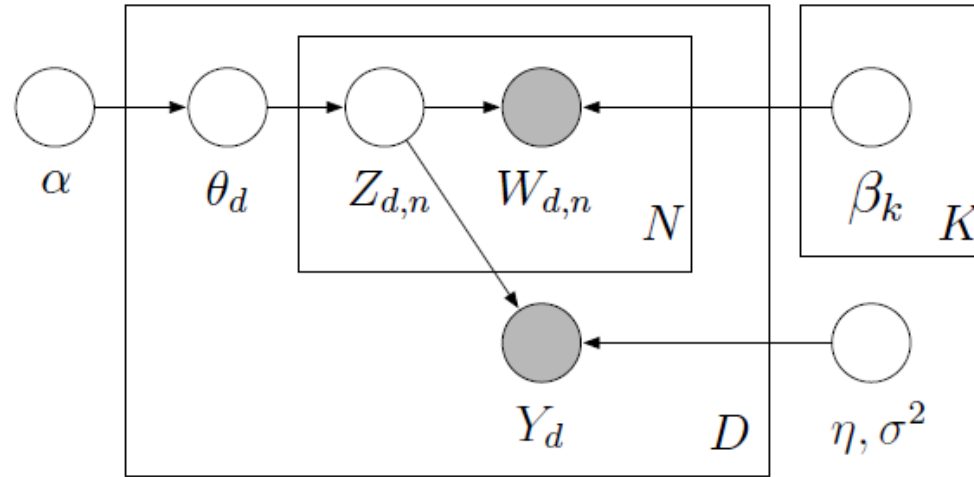


# Correlated topic models (CTM) [Blei et al. 2005]



- Draw topics from a logistic normal, where topic occurrence can exhibit correlations

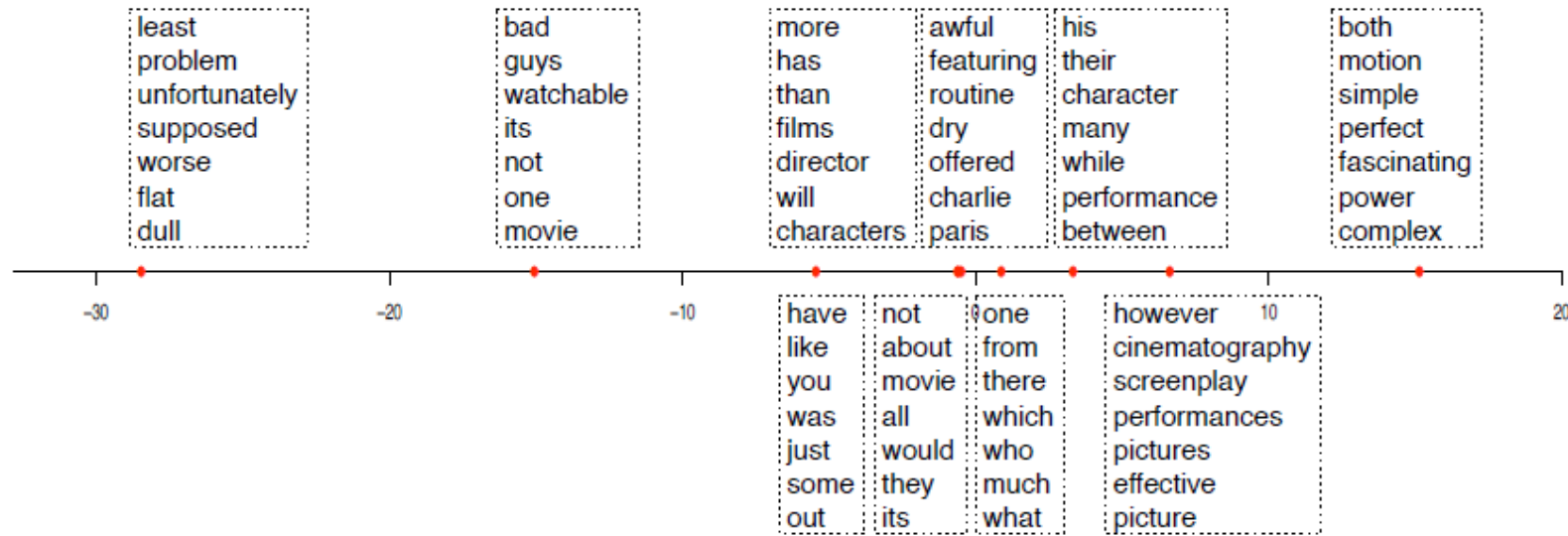
# Supervised LDA [Blei & McAuliffe 07]



- 1 Draw topic proportions  $\theta \mid \alpha \sim \text{Dir}(\alpha)$ .
- 2 For each word
  - Draw topic assignment  $z_n \mid \theta \sim \text{Mult}(\theta)$ .
  - Draw word  $w_n \mid z_n, \beta_{1:K} \sim \text{Mult}(\beta_{z_n})$ .
- 3 Draw response variable  $y \mid z_{1:N}, \eta, \sigma^2 \sim \text{N}(\eta^\top \bar{z}, \sigma^2)$ , where

$$\bar{z} = (1/N) \sum_{n=1}^N z_n.$$

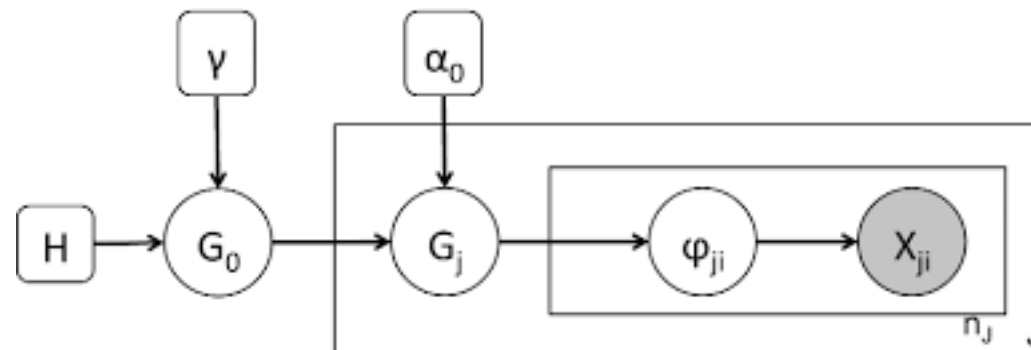
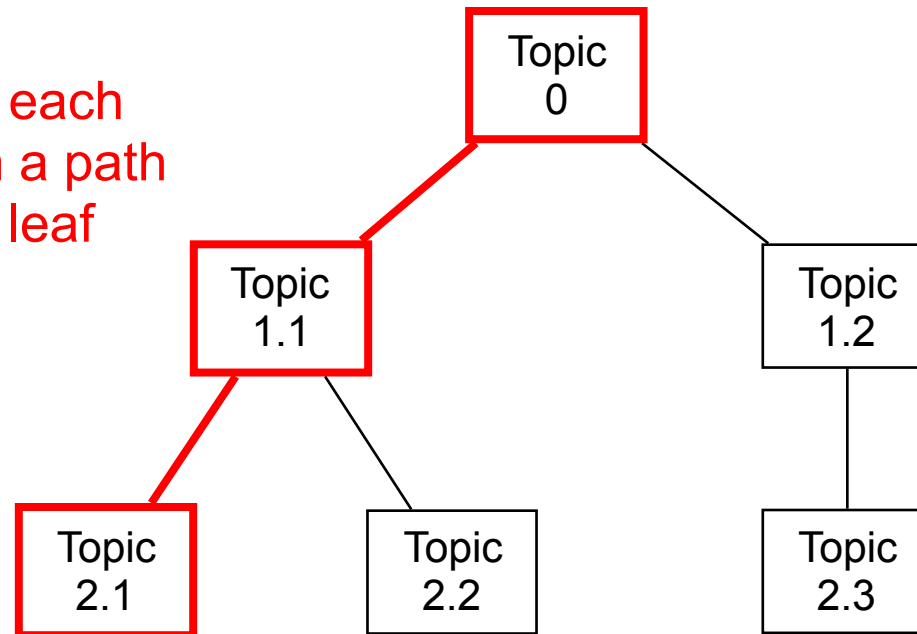
# Supervised LDA [Blei & McAuliffe 07]



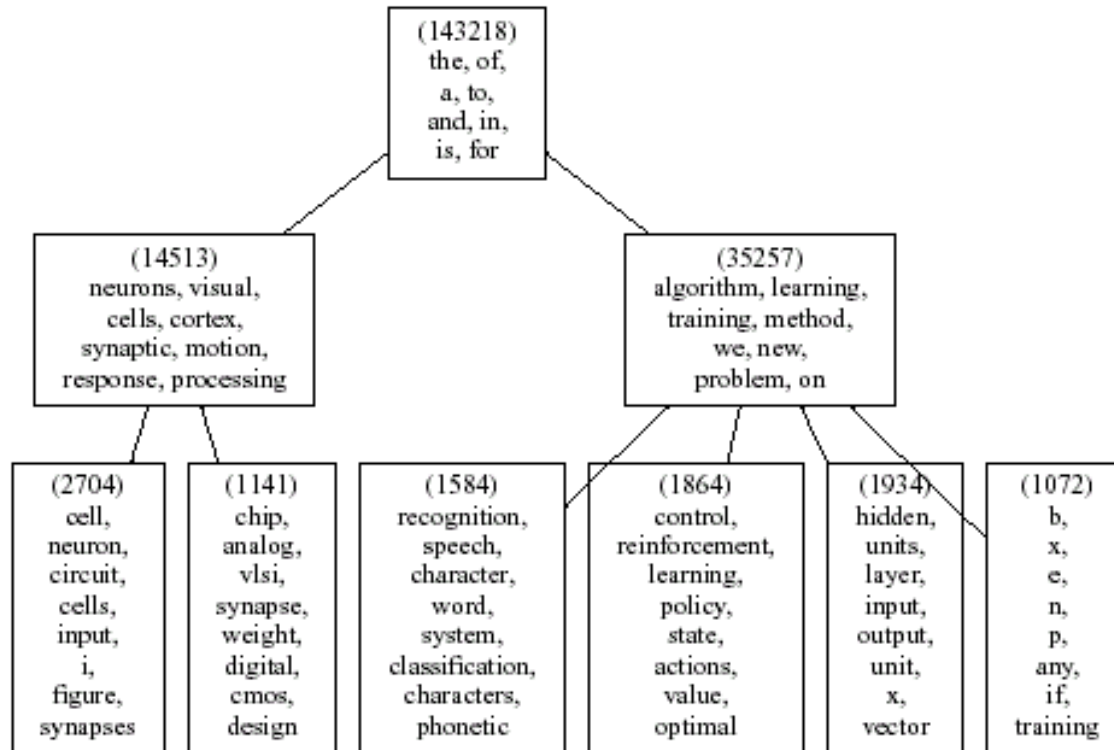
- 10-topic sLDA model on movie reviews (Pang and Lee, 2005).
- Response: number of stars associated with each review

# Learning topic hierarchies

The topics in each document form a path from root to leaf



# Learning topic hierarchies



# Today's lecture

- Latent semantic indexing
- Continue on topic model
  - Bayesian inference of topic model
  - Variational inference for LDA
  - Gibbs sampling, Markov chain Monte-Carlo