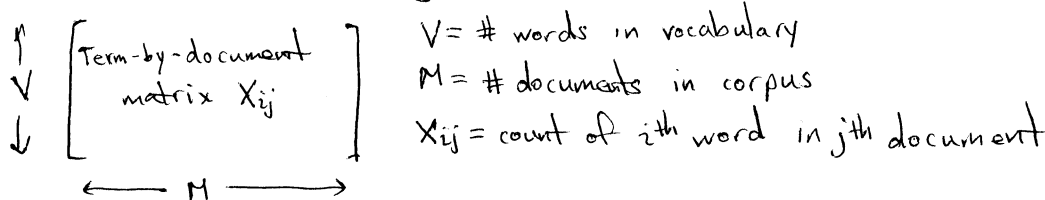


Review

• Document modeling

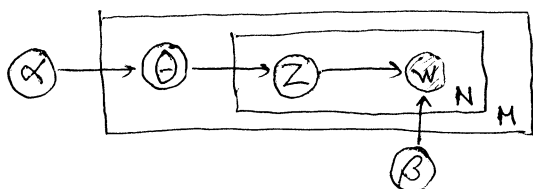


• Notation

Document $\vec{w} = (w_1, w_2, \dots, w_N)$ sequence of N words

Corpus $D = (\vec{w}_1, \vec{w}_2, \dots, \vec{w}_M)$ collection of M documents

• Latent Dirichlet allocation (LDA)



• Generative model

For each document \vec{w} in corpus D :

1) Choose $N = \# \text{ words}$

2) Choose topic weights θ from $P(\theta | \alpha) = \frac{\Gamma(\sum_{k=1}^K \alpha_k)}{\prod_k \Gamma(\alpha_k)} \prod_{k=1}^K \theta_k^{\alpha_k - 1}$

3) Do N times:

a) choose topic $z_n \in \{1, 2, \dots, K\}$ from $P(z_n | \theta) = \frac{\theta_{z_n}}{\sum_{k=1}^K \theta_k} = \prod_{k=1}^K \theta_k^{I(z_n, k)}$ ↑ indicator function

b) choose word $w_n \in \{1, 2, \dots, V\}$ from $P(w_n = j | z_n = k) = \beta_{kj}$

How to learn parameters $\{\alpha_k\}$ and $\{\beta_{kj}\}$ from data X ?

Learning from "complete data" (warm-up)

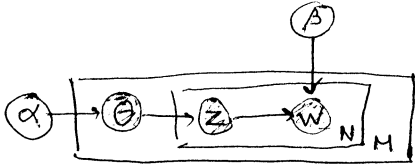
• Joint probabilities

$$P(\theta, \vec{z}, \vec{w} | \alpha, \beta) = P(\theta | \alpha) \prod_{n=1}^N P(z_n | \theta) P(w_n | z_n, \beta)$$

↑ topic weights
 ↑ topic sequence
 ↑ word sequence
 ↑ parameters
 (per document)

$$P(\{\theta_m\}_{m=1}^M, \{\vec{z}_m\}, D | \alpha, \beta) = \prod_{m=1}^M P(\theta_m, \vec{z}_m, \vec{w}_m | \alpha, \beta)$$

(over corpus)



• Log-likelihood

$$\mathcal{L}(\alpha, \beta) = \underbrace{\sum_m \log P(\theta_m | \alpha)}_{\mathcal{L}(\alpha)} + \sum_{nm} \log P(z_{nm} | \theta_m) + \underbrace{\sum_{nm} \log P(w_{nm} | z_{nm}, \beta)}_{\mathcal{L}(\beta)}$$

• Maximum likelihood (ML) estimation

For multinomial parameters:

$$\beta^* = \operatorname{argmax}_{\beta} \mathcal{L}(\beta)$$

$$\frac{\partial \mathcal{L}}{\partial \beta} = 0 \longrightarrow \beta_{kj}^* = \frac{\sum_{n,m} \mathbb{I}(w_{nm}, j) \mathbb{I}(z_{nm}, k)}{\sum_{n,m} \mathbb{I}(z_{nm}, k)}$$

simple ratio of counts

For dirichlet parameters:

$$\alpha^* = \operatorname{argmax}_{\alpha} \mathcal{L}(\alpha)$$

$$\mathcal{L}(\alpha) = \sum_m \left\{ \log \Gamma\left(\sum_k \alpha_k\right) - \sum_k \log \Gamma(\alpha_k) + \sum_k (\alpha_k - 1) \log \theta_{km} \right\}$$

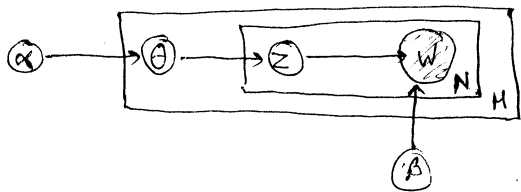
↑ not obvious, but concave α
 ↙ observed topic weights

Define: $\psi(x) = \frac{d}{dx} [\log \Gamma(x)]$ "digamma function"

$$\frac{\partial \mathcal{L}}{\partial \alpha_k} = 0 \rightarrow \frac{1}{M} \sum_{m=1}^M \log \theta_{km} = \Psi(\alpha_k) - \Psi\left(\sum_{\ell=1}^k \alpha_\ell\right)$$

These nonlinear equations have a unique solution for α^* , which can be computed by (say) Newton's method.

Incomplete data



• Inference

$$P(\theta, \vec{z} | \vec{w}, \alpha, \beta) = \frac{P(\theta, \vec{z}, \vec{w} | \alpha, \beta)}{P(\vec{w} | \alpha, \beta)}$$

How to compute denominator?

• Marginal probability of document

$$P(\vec{w} | \alpha, \beta) = \int d\theta \sum_{\vec{z}} P(\theta, \vec{z}, \vec{w} | \alpha, \beta)$$

↙ product over N words in document

$$= \int d\theta P(\theta | \alpha) \left\{ \prod_{n=1}^N \left[\sum_{z_n=1}^K P(z_n | \theta) P(w_n | z_n, \beta) \right] \right\}$$

$$P(\vec{w}_j | \alpha, \beta) = \int_0^1 d\theta_1 \int_0^{1-\theta_1} d\theta_2 \dots \int_0^{1-\theta_1-\dots-\theta_{k-1}} d\theta_k \frac{\Gamma(\sum \alpha_k)}{\prod_k \Gamma(\alpha_k)} \prod_k \theta_k^{\alpha_k - 1} \left\{ \prod_{i=1}^V \left(\sum_{k=1}^K \theta_k \beta_{kj} \right)^{x_{ij}} \right\}$$

↙ product over V words in vocabulary, weighted by word-document count x_{ij}

• Log-likelihood of corpus

$$\log P(D | \alpha, \beta) = \sum_{j=1}^M \log P(\vec{w}_j | \alpha, \beta)$$

How to maximize what we cannot compute?

• EM algorithm

- E-step compute statistics of posterior distribution:

$$\text{e.g. } P(z_n=k|\vec{w}, \alpha, \beta), E[\log \theta_k | \vec{w}, \alpha, \beta]$$

- M-step updates α, β based on these statistics.

- common approximations for intractable E-step:

(1) Markov chain Monte Carlo (MCMC)

(2) Variational method

• Review of MCMC

To estimate $P(\theta|\vec{w}, \alpha, \beta)$ and $P(z_n|\vec{w}, \alpha, \beta)$:

- Fix words w_1, w_2, \dots, w_n to observed values

- Initialize hidden nodes $\{\theta_k\}_{k=1}^K, \{z_n\}_{n=1}^N$ at random values

- Repeat S times: ↖ (constrained to sum to 1)

• pick hidden node $X \in \{\theta, z_1, \dots, z_n\}$ at random

• use Bayes rule to compute $P(X | \text{all other nodes at current values})$

• resample X from this distribution

• Notation

denote hidden configurations as

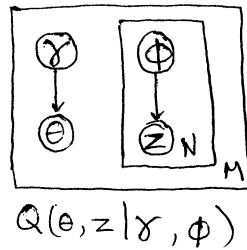
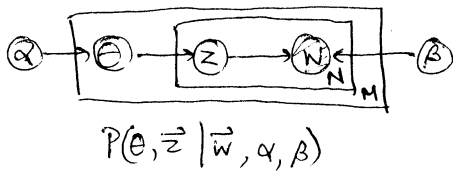
$$\{\theta^{(s)}, z_1^{(s)}, \dots, z_n^{(s)}\}_{s=1}^S$$

$$\cdot \text{Estimate } P(z_n=k|\vec{w}, \alpha, \beta) = \frac{1}{S} \sum_{s=1}^S \mathbb{I}(z_n^{(s)}, k)$$

$$\cdot \text{Estimate } P(\theta|\vec{w}, \alpha, \beta) = \frac{1}{S} \sum_{s=1}^S \delta(\theta - \theta^{(s)})$$

• Variational method

Approximate intractable $P(\theta, \vec{z} | \vec{w}, \alpha, \beta)$ by a tractable distribution $Q(\theta, z | \gamma, \vec{\phi})$



if θ and z are peaked, the approx may have the same mode as the original

• Explicit form of approximation:

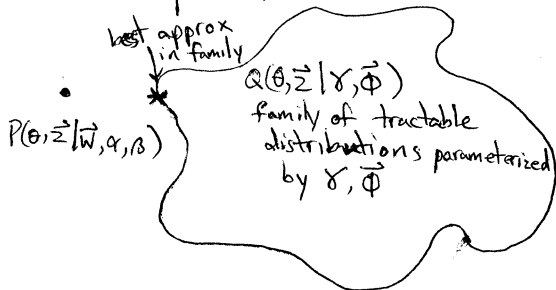
$$Q(\theta, \vec{z} | \gamma, \vec{\phi}) = \underbrace{\frac{\Gamma(\sum \gamma_k)}{\prod_{k=1}^K \Gamma(\gamma_k)}}_{\text{dirichlet dist.}} \prod_k \theta_k^{\gamma_k - 1} \underbrace{\prod_{n=1}^N \prod_{k=1}^K \phi_{kn}^{\mathbb{I}(z_n, k)}}_{\text{distribution over the } K \text{ topics over } N \text{ words}}$$

• Variational inference

• choose "variational" parameters $\gamma, \vec{\phi}$ to obtain best approximation:

$$\operatorname{argmin}_{\gamma, \vec{\phi}} [KL(Q(\theta, \vec{z} | \gamma, \vec{\phi}), P(\theta, \vec{z} | \vec{w}, \alpha, \beta))]$$

• Cartoon picture



• tractable to compute $KL(Q, P)$ up to constant independent of $\{\gamma, \vec{\phi}\}$

• Useful lower bound

shorthand: drop dependence on $\alpha, \beta, \gamma, \vec{\phi}$

$$\begin{aligned}\log P(\vec{w}) &= \log \left[\frac{P(\vec{w}, \theta, \vec{z})}{P(\theta, \vec{z} | \vec{w})} \right] \text{ for all } \theta, \vec{z} \\ &= \int d\theta \sum_{\vec{z}} Q(\theta, \vec{z}) \log \left[\frac{P(\vec{w}, \theta, \vec{z})}{P(\theta, \vec{z} | \vec{w})} \right] \text{ for any distribution } Q \\ &= \int d\theta \sum_{\vec{z}} Q(\theta, \vec{z}) \log \left[\frac{P(\vec{w}, \theta, \vec{z})}{Q(\theta, \vec{z} | \vec{w})} \right] + \underbrace{\text{KL}(Q(\theta, \vec{z}), P(\theta, \vec{z} | \vec{w}))}_{\geq 0} \\ &\geq \int d\theta \sum_{\vec{z}} Q(\theta, \vec{z}) \log \left[\frac{P(\vec{w}, \theta, \vec{z})}{Q(\theta, \vec{z})} \right]\end{aligned}$$