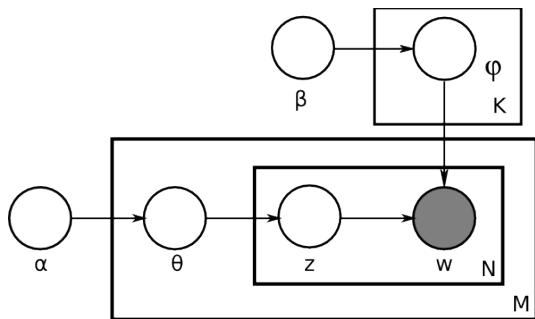


# Lecture 19: Probabilistic topic models II: LDA (part 2)

William Webber ([william@williamwebber.com](mailto:william@williamwebber.com))

COMP90042, 2014, Semester 1, Lecture 19

# LDA diagram and process



1. Choose term probabilities for each topic:  $\Phi_i \sim \mathcal{D}(\beta)$
2. Choose topic probabilities for each document:  $\Theta_d \sim \mathcal{D}(\alpha)$
3. Choose the topic of each token:  $z_{dn} \sim \mathcal{M}(\theta_d)$
4. Choose the token:  $w_{dn} \sim \mathcal{M}(\phi_{z_{dn}})$

# The multinomial distribution

- ▶ Let  $X$  be a multinomial random variable
- ▶ A “realization” of  $X$  takes on  $k$  distinct values,  $\{X_1, \dots, X_i, \dots, X_k\}$
- ▶  $X$  has  $k + 1$  parameters:
  - ▶  $n > 0$ , the number of “trials”
  - ▶  $\{p_1, \dots, p_k\}$ , the probability of each distinct value at each trial
    - ▶  $0 \leq p_i \leq 1$
    - ▶  $\sum_{i=1}^k p_i = 1$
    - ▶ i.e.  $\mathbf{p}$  is a probability distribution over  $k$  values
- ▶  $X_j \in \{0, 1, \dots, n\}$
- ▶  $\sum X_i = n$
- ▶ Intuition
  - ▶ Roll a biased  $k$ -sided dice  $n$  times
  - ▶ Count number of times each face turns up
  - ▶  $X_i$  is the number of times the  $i$ 'th face turns up
- ▶ There is a formula but we don't need to worry about it!

# The Dirichlet distribution

- ▶ Let  $\psi$  be a Dirichlet random variable
- ▶ A “realization” of  $\psi$  takes on  $k$  values,  $\{\psi_1, \dots, \psi_i, \dots, \psi_k\}$ 
  - ▶  $0 \leq \psi_i \leq 1$
  - ▶  $\sum_{i=1}^k \psi_i = 1$
  - ▶ i.e. a realization of  $\psi$  is itself a probability distribution
- ▶  $\psi$  has  $k$  parameters,  $\alpha = \{\alpha_1, \dots, \alpha_k\}$ , with  $\alpha_i > 0$
- ▶ Let  $A = \sum \alpha_i$
- ▶ The expected value of  $\psi_i$  is  $\alpha_i/A$  (written  $\mathbb{E}[\psi_i] = \alpha_i/A$ )
- ▶ The greater  $A$ , the closer  $\psi_i$  is likely to be to  $\alpha_i/A$
- ▶ A realization of  $\psi$  can give us the parameters  $\{p_1, \dots, p_k\}$  for a multinomial variable  $X$
- ▶ If  $\alpha_i = \alpha_j \forall i, j$ , we say that  $\psi$  is symmetric
- ▶ There is a formula but we don't need to worry about it!

# The probability equations of LDA

$$P(\mathbf{w}_1, \dots, \mathbf{w}_n | \alpha, \beta) = \prod_{i=1}^n P(\mathbf{w}_i | \alpha, \beta)$$

$$P(\mathbf{w}_i | \alpha, \beta) = \int P(\mathbf{w}_i | \alpha, \phi) P(\phi | \beta) d\beta$$

$$P(\mathbf{w}_i | \alpha, \phi) = \int P(\mathbf{w}_i, \theta_i | \alpha, \phi) d\theta_i$$

$$P(\mathbf{w}_i, \theta_i | \alpha, \phi) = P(\mathbf{w}_i | \theta_i, \phi) P(\theta_i | \alpha)$$

$$P(\mathbf{w}_i | \theta_i, \phi) = \prod_{j=1}^m P(w_{i,j} | \theta_i, \phi)^{c_{i,j}}$$

$$P(w_{i,j} | \theta_i, \phi) = \sum_{k=1}^K P(z_{i,j} = k | \theta_i) P(w_{i,j} | \phi_k)$$

$K$  number of topics

$\mathbf{w}_i$  bag of words (terms and  $f_{d,t}$ ) for document  $i$

## The components: $\theta_i$

- ▶  $\theta_i$  is the multinomial distribution over topics for document  $i$
- ▶ There are  $K$  topics (where  $K$  is semi-arbitrarily chosen by us)
- ▶ Therefore  $\theta_i$  has  $K$  parameters
- ▶  $\theta_{i,k}$  is the probability that an arbitrary word in document  $i$  will belong to topic  $k$
- ▶  $\alpha$  is the prior to  $\theta$ 
  - ▶ That is, the probabilities  $\{\theta_{i,1}, \dots, \theta_{i,K}\}$  are a “realization” of a Dirichlet random variable with parameters  $\{\alpha_1, \dots, \alpha_K\}$
- ▶ The same Dirichlet RV is prior to all  $\theta_i$
- ▶  $\alpha$  is asymmetric, meaning we allow certain topics to be a *priori* more likely than others

## The components: $\phi_k$

- ▶  $\phi_k$  is the multinomial distribution over terms for topic  $k$
- ▶ There are  $m = |V|$  terms in the vocabulary
- ▶ Therefore  $\phi_k$  has  $m$  parameters
- ▶  $\phi_{k,i}$  is the probability that an arbitrary word produced by topic  $k$  will be  $V_i$
- ▶  $\beta$  is the Dirichlet prior to  $\phi$ 
  - ▶ That is, the probabilities  $\{\phi_{i,1}, \dots, \phi_{i,m}\}$  are a “realization” of a Dirichlet random variable with parameters  $\{\beta_1, \dots, \beta_m\}$
- ▶  $\beta$  is symmetric, meaning that all words are *a priori* as likely for each topic
  - ▶ This does *not* mean that
    - ▶ the posterior distribution  $\phi_k$  over terms for topics will be flat
    - ▶ each  $\phi_k$  will give the same distribution over terms

## Symmetric and asymmetric priors

Why is  $\beta$  symmetric, but  $\alpha$  asymmetric? Following Wallach et al., 2009

- ▶ Asymmetric  $\alpha$  leads to more stable results
- ▶ In particular, models are more stable to choice of number of topics  $K$ 
  - ▶ Think back to LSI, where some topics are “more important” than others
  - ▶ And model for topic  $k \leq K$  is independent of choice of  $K$
- ▶ However, asymmetric  $\beta$  does not improve stability
- ▶ Asymmetric requires more parameters to fit than symmetric
- ▶ Therefore, only employ asymmetric priors if they provide some advantage

(In the full/pure Bayesian model, we apply yet another prior to  $\alpha$ , known as a gamma prior; in practice, this is approximated using empirical methods.)



## Deciphering the formulae: $P(w_{i,j}|\theta_i, \phi)$

$$P(w_{i,j}|\theta_i, \phi) = \sum_{k=1}^K P(z_{i,j} = k|\theta_i)P(w_{i,j}|\phi_k) \quad (1)$$

- ▶  $P(w_{i,j})$  is the prob that an arbitrary term in doc  $i$  is the term  $j$ 
  - ▶ We don't care about the word position in the doc
  - ▶ ... a standard assumption of the unigram term model
- ▶  $\phi_k$  is the probability distribution over terms for topic  $k$
- ▶ Therefore,  $P(w_{i,j}|\phi_k)$  is just  $P(j|\phi_k)$
- ▶  $z_{i,j}$  is the topic that generates term  $j$  of document  $i$
- ▶  $P(z_{i,j} = k|\theta_i)$  is the probability that this topic is  $k$
- ▶  $\theta_i$  is the prob dist over topics for document  $i$
- ▶ Therefore,  $P(z_{i,j} = k|\theta_i)$  is just  $P(k|\theta_i)$

## Deciphering the formulae: $P(w_{i,j}|\theta_i, \phi)$ (cont.)

$$P(w_{i,j} = t|\theta_i, \phi) = \sum_{k=1}^K P(z_{i,j} = k|\theta_i)P(w_{i,j} = t|\phi_k) \quad (2)$$

- ▶  $P(z_{i,j} = k|\theta_i)P(w_{i,j}|\phi_k)$  is the probability that the term is  $j$  and the topic is  $k$
- ▶ The term must come from exactly one topic
- ▶ Therefore, we sum these probabilities over all  $K$  topics
- ▶ And this gives us  $P(w_{i,j}|\theta_i, \phi)$ 
  - ▶ That is, the probability that an arbitrary term in doc  $i$  is  $j$

## Deciphering the formulae: $P(\mathbf{w}_i|\theta_i, \phi)$

$$P(\mathbf{w}_i|\theta_i, \phi) = \prod_{j=1}^m P(w_{i,j}|\theta_i, \phi)^{c_{i,j}} \quad (3)$$

- ▶  $c_{i,j}$  is the number of times that term  $j$  occurs in document  $i$
- ▶ We assume (unigram model) that these occurrences are independent
- ▶ Therefore the probability that term  $j$  occurs  $c_{i,j}$  times is the probability of each occurrence, raised to the  $c_{i,j}$ 'th power
- ▶ We also assume (unigram model) that the occurrence of different terms is independent
- ▶ Therefore, the probability of the bag-of-words representation  $\mathbf{w}_i$  of document  $i$  is just the product of all the individual probabilities

## Deciphering the formulae: $P(\mathbf{w}_i, \theta_i | \alpha, \phi)$

$$P(\mathbf{w}_i, \theta_i | \alpha, \phi) = P(\mathbf{w}_i | \theta_i, \phi) P(\theta_i | \alpha) \quad (4)$$

- ▶  $P(\mathbf{w}_i, \theta_i | \alpha, \phi)$  is the (joint) probability of observing:
  - ▶ the bag-of-words representation  $\mathbf{w}_i$  of document  $i$
  - ▶ document  $i$ 's distribution over topics  $\theta_i$
- ▶ We have previously figured out  $P(\mathbf{w}_i | \theta_i, \phi)$
- ▶ The multinomial  $\theta_i$  is an “observation” of the Dirichlet RV  $\alpha$
- ▶ So  $P(\theta_i | \alpha)$  is the prob of the multinomial  $\theta_i$  given the prior  $\alpha$ 
  - ▶ There is a formula for this (that we won't worry about!)
- ▶ We assume conditional independence of  $\mathbf{w}_i$  and  $\theta_i$
- ▶ So their joint probability is just the product of their individual (marginal) probabilities

## Deciphering the formulae: $P(\mathbf{w}_i|\alpha, \phi)$

$$P(\mathbf{w}_i|\alpha, \phi) = \int P(\mathbf{w}_i, \theta_i|\alpha, \phi) d\theta_i \quad (5)$$

- ▶  $P(\mathbf{w}_i|\alpha, \phi)$  is
  - ▶ the probability of the bag-of-words document  $i$
  - ▶ given our prior for document distributions over topics  $\alpha$
  - ▶ and the (list of  $k$ ) topic distributions over words  $\phi$
- ▶  $\theta_i$  is the document distribution over topics for doc  $i$
- ▶ We already have  $P(\mathbf{w}_i, \theta_i|\alpha, \phi)$  (from previous step)
- ▶ We can “remove”  $\theta_i$  by integrating over all  $\theta_i$ 's
  - ▶ If you're unfamiliar with calculus, think of the integral as analogous to a sum over a continuous variable

## Deciphering formulae: $P(\mathbf{w}_i|\alpha, \beta)$

$$P(\mathbf{w}_i|\alpha, \beta) = \int P(\mathbf{w}_i|\alpha, \phi)P(\phi|\beta) d\beta \quad (6)$$

- ▶  $\beta$  is our (symmetric) prior for topic distributions over terms
- ▶ We've already calculated  $P(w_i|\alpha, \phi)$
- ▶ We can write:

$$P(\mathbf{w}_i|\alpha, \beta) = \int P(\mathbf{w}_i|\alpha, \beta, \phi)P(\phi|\alpha, \beta) d\beta \quad (7)$$

by the law of total probability

- ▶ Analogous to:

$$P(A) = \sum_b P(A|B = b)P(B = b) \quad (8)$$

- ▶  $w_i$  is independent of  $\beta$ , given  $\theta$ , and  $\phi$  is independent of  $\alpha$
- ▶ Therefore, Equation 7 simplifies to Equation 6

## Deciphering the formulae: $P(\mathbf{w}_1, \dots, \mathbf{w}_n | \alpha, \beta)$

$$P(\mathbf{w}_1, \dots, \mathbf{w}_n | \alpha, \beta) = \prod_{i=1}^n P(\mathbf{w}_i | \alpha, \beta) \quad (9)$$

- ▶ We assume that documents are probabilistically independent
- ▶ Therefore the probability of generating a set of documents  $\{\mathbf{w}_1, \dots, \mathbf{w}_n\}$
- ▶ ... is the product of the probability of generating each individual document

# The probability equations of LDA

$$P(\mathbf{w}_1, \dots, \mathbf{w}_n | \alpha, \beta) = \prod_{i=1}^n P(\mathbf{w}_i | \alpha, \beta)$$

$$P(\mathbf{w}_i | \alpha, \beta) = \int P(\mathbf{w}_i | \alpha, \phi) P(\phi | \beta) d\phi$$

$$P(\mathbf{w}_i | \alpha, \phi) = \int P(\mathbf{w}_i, \theta_i | \alpha, \phi) d\theta_i$$

$$P(\mathbf{w}_i, \theta_i | \alpha, \phi) = P(\mathbf{w}_i | \theta_i, \phi) P(\theta_i | \alpha)$$

$$P(\mathbf{w}_i | \theta_i, \phi) = \prod_{j=1}^m P(w_{i,j} | \theta_i, \phi)^{c_{i,j}}$$

$$P(w_{i,j} | \theta_i, \phi) = \sum_{k=1}^K P(z_{i,j} = k | \theta_i) P(w_{i,j} | \phi_k)$$



# Solving LDA

Directly solving LDA would involve finding parameters that maximize the empirical likelihood  $\mathcal{L}$  of the observed documents  $\{\mathbf{w}_1, \dots, \mathbf{w}_n\}$ :

$$\mathcal{L} = \prod_{j=1}^m \prod_{i=1}^n P(w_{i,j} | z_{i,j}, \phi) P(z_{i,j} | \theta_i) P(\theta_d | \alpha) P(\phi | \beta) \quad (10)$$

- ▶ Note: parameters found are not  $\alpha, \beta$
- ▶ Rather, they are parameters to  $\phi, \theta$

These parameters cannot be directly solved as  $z_{i,j}$  not observed.

Instead, an approximation method must be used.

# Gibbs sampling

A common approach is to use *Gibbs sampling*

- ▶ A Monte-Carlo Markov Chain method from statistical physics
  - ▶ “Monte Carlo” means based on random simulation
  - ▶ “Markov Chain” describes a random process in which each state depends only on the previous state
- ▶ Basic idea is in a complex model with many dependent variables:
  - ▶ Sequentially sample each variable, dependent upon state of all other variables
  - ▶ Observe averages over very large number of samples as probability estimates

# Collapsed Gibbs sampling

Collapsed Gibbs method developed by Griffiths and Steyvers, 2006:

- ▶ Marginalize out  $\theta, \phi$
- ▶ Instead, estimate  $P(\mathbf{z}|\mathbf{w})$  (that is,  $P(z_{i,j}|w_{i,j})$  for all  $i, j$ )
- ▶ They derive the approximation:

$$P(z_{i,j}|\bar{z}_{i,j}) \propto (N_{i,z} + \alpha_z)(N_{z,j} + \beta) \quad (11)$$

where

- $\bar{z}_{i,j}$  all other topic assignments to words
- $N_{i,z}$  number of times topic  $z$  has been assigned to words in document  $i$
- $N_{z,j}$  number of times word  $j$  has been assigned to topic  $z$

- ▶ Iterate many, many times
- ▶ Count how many times each word assigned to each topic
- ▶ Normalize these counts to estimate  $\theta_i, \phi_k$

## Further reading

- ▶ Blei, Ng, and Jordan, “Latent Dirichlet Allocation”, JMLR, 2003
- ▶ Crain, Zhou, Yang, and Zha, “Dimensionality Reduction and Topic Modeling”, Chapter 5 of Aggarwal and Zhai (ed.), *Mining Text Data*, 2012 (brief summary of Gibbs sampling).
- ▶ Sun, Deng, and Han, “Probabilistic Models for Text Mining”, Chapter 8 of Aggarwal and Zhai (ed.), *Mining Text Data*, 2012 (gives  $P(\mathbf{w}_1, \dots, \mathbf{w}_n | \alpha, \beta)$ ).
- ▶ Griffiths and Steyvers, “Finding Scientific Topics”, *PNAS*, 2004 (collapsed Gibbs sampling for solving LDA)