
Probabilistic Graphical Models IFT6269 Project

Latent Dirichlet Allocation

Patrice B  chard ^{* 1} Th  ophile Gervet ^{* 2}

Abstract

TODO

1. Introduction and Related Works

1.1. Latent Dirichlet Allocation

1.1.1. GENERATIVE PROCESS

Latent Dirichlet allocation (LDA) is a generative probabilistic model of a corpus of D documents $\{\mathbf{w}_1, \dots, \mathbf{w}_D\}$. Each document $\mathbf{w}_d = (w_{d,1}, \dots, w_{d,N_d})$ is a sequence of N_d words from a vocabulary indexed by $\{1, \dots, V\}$.

The basic idea behind LDA is that each document is represented as a random mixture over K latent topics, where each topic is characterized by a multinomial distribution over words, parametrized by a vector β_k , $k \in 1, \dots, K$.

Formally, LDA assumes the following generative process for each document \mathbf{w}_d :

1. Sample a distribution over topics $\theta_d \sim \text{Dir}(\alpha)$
2. For each of the N_d words $w_{d,n}$ independently:
 - (a) Draw a topic from the distribution over topics $z_{d,n} \sim \text{Mult}(\theta_d)$
 - (b) Draw a word from this topic $w_{d,n} \sim \text{Mult}(\beta_{z_{d,n}})$

This generative process is illustrated in the directed graphical model of figure 1.

^{*}Equal contribution ¹Universit   de Montr  al, Montr  al, Canada ²McGill University, Montr  al, Canada. Correspondence to: Patrice B  chard <patrice.bechard@umontreal.ca>, Th  ophile Gervet <theophile.gervet@umontreal.ca>.

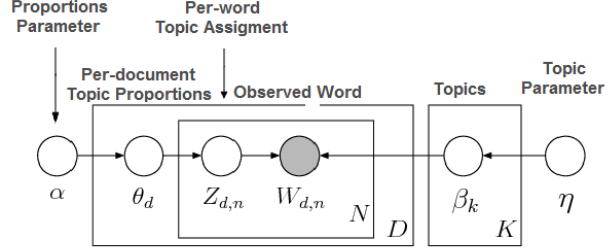


Figure 1. Graphical model corresponding to the generative process for LDA.

The joint distribution associated with this model is:

$$\prod_{k=1}^K p(\beta_k | \eta) \prod_{d=1}^D p(\theta_d | \alpha) \prod_{n=1}^{N_d} p(z_{d,n} | \theta_d) p(w_{d,n} | z_{d,n}, \beta_{1:K}) \quad (1)$$

where

$$\begin{aligned} z_{d,n} &\sim \text{Mult}(\theta_d), \text{ i.e } p(z_{d,n} | \theta_d) = \theta_{d,z_{d,n}} \\ w_{d,n} &\sim \text{Mult}(\beta_{z_{d,n}}), \text{ i.e } p(w_{d,n} | z_{d,n}, \beta_{1:K}) = \beta_{z_{d,n}, w_{d,n}} \\ \beta_k &\sim \text{Dir}(\eta) \\ \theta_d &\sim \text{Dir}(\alpha) \end{aligned}$$

This generative process is just a model of the structure we assume to be in our data. In reality, we only observe documents $\mathbf{w}_{1:D}$ and our goal is to infer the underlying topic structure: the topics $\beta_{1:K}$, the per document distributions over topics θ_d , and per document per word topic assignments $z_{d,n}$.

Since the generative process assumes independence between documents, the probability of the whole corpus decomposes as a product of terms for individual documents. In order to clarify the notation, from now on let us consider equations for a single document \mathbf{w} with N words.

Assuming the Dirichlet hyper-parameters α and η are

fixed, we want to infer the posterior distribution:

$$p(\beta, \theta, \mathbf{z} | \mathbf{w}, \alpha, \eta) = \frac{p(\beta, \theta, \mathbf{z}, \mathbf{w} | \alpha, \eta)}{p(\mathbf{w} | \alpha, \eta)} \quad (2)$$

Unfortunately, this distribution is intractable to compute because the normalization factor $p(\mathbf{w} | \alpha, \eta)$ cannot be computed exactly. We must use approximate inference to estimate the posterior over latent variables β , θ and \mathbf{z} . We will explore both variational inference and Gibbs sampling to solve this problem.

1.2. Related Works

TODO

2. Methods

2.1. Variational EM

To make our life easier, we derive a variational EM algorithm for a slightly simpler graphical model for LDA (the one presented in the original paper). We remove the Dirichlet prior over topics, this corresponds to removing the dependence of the β plate on the η random variable (see Figure 2, left). By doing so, we lose our ability to smooth topics through the η parameter.

Assuming the α parameter and topics β are fixed, the posterior for a single document \mathbf{w} now takes the form:

$$p(\theta, \mathbf{z} | \mathbf{w}, \alpha, \beta) = \frac{p(\theta, \mathbf{z}, \mathbf{w} | \alpha, \beta)}{p(\mathbf{w} | \alpha, \beta)} \quad (3)$$

This expression is still intractable because of the denominator $p(\mathbf{w} | \alpha, \beta)$.

Let us motivate a variational EM algorithm by looking at how this intractable posterior comes back to bite us when we try to maximize the likelihood of the observed data. Since the probability of the whole corpus decomposes as a product of terms for individual documents, the log likelihood of the whole corpus decomposes as a sum of terms for individual documents. The log likelihood of a single document \mathbf{w} takes the form:

$$\log p(\mathbf{w} | \alpha, \beta) = \log \int \sum_{\mathbf{z}} p(\theta, \mathbf{z}, \mathbf{w} | \alpha, \beta) d\theta \quad (4)$$

Because of the integral over θ and the sum over \mathbf{z} inside the log, we cannot maximize this objective directly. One common way to solve this problem is the Expectation Maximization (EM) algorithm. We make the following obser-

vation:

$$\log p(\mathbf{w} | \alpha, \beta) = \log \int \sum_{\mathbf{z}} \frac{p(\theta, \mathbf{z}, \mathbf{w} | \alpha, \beta) q(\theta, \mathbf{z})}{q(\theta, \mathbf{z})} d\theta \quad (5)$$

$$= \mathbf{E}_q[\log \int \sum_{\mathbf{z}} \frac{p(\theta, \mathbf{z}, \mathbf{w} | \alpha, \beta)}{q(\theta, \mathbf{z})} d\theta] \quad (6)$$

$$\geq \mathbf{E}_q[\log p(\theta, \mathbf{z}, \mathbf{w} | \alpha, \beta)] + \mathbf{E}_q[\log q(\theta, \mathbf{z})] \quad (7)$$

$$= \mathcal{L}(q, \alpha, \beta) \quad (8)$$

where we have introduced an arbitrary distribution $q(\theta, \mathbf{z})$ over the problematic latent variables θ and \mathbf{z} , and the last step follows from Jensen's inequality.

The EM algorithm consists in alternatively maximizing the lower bound \mathcal{L} on the log likelihood with respect to the distribution q (the E-step), and with respect to parameters α and β (the M-step).

One can easily show that maximizing \mathcal{L} with respect to q is equivalent to minimizing $KL(q(\theta, \mathbf{z}) || p(\theta, \mathbf{z} | \mathbf{w}, \alpha, \beta))$. A simple solution is thus to set $q(\theta, \mathbf{z})$ to be $p(\theta, \mathbf{z} | \mathbf{w}, \alpha, \beta)$, the true posterior over latent variables under the model.

In our case, things are not so simple: as we have seen, the posterior $p(\theta, \mathbf{z} | \mathbf{w}, \alpha, \beta)$ is intractable. This is where variational inference comes into play. The basic idea is to introduce a parametrized family of distributions over latent variables and phrase inference as an optimization problem.

We parametrize q as follows, according to the mean field independence assumption:

$$q(\theta, \mathbf{z} | \gamma, \phi) = q(\theta | \gamma) \prod_{n=1}^N q(z_n | \phi_n) \quad (9)$$

where the dirichlet parameters γ and the multinomial parameters ϕ_n , $n \in \{1, \dots, N\}$ are free variational parameters. The graphical model corresponding to this parameterization is illustrated in Figure 2, right.

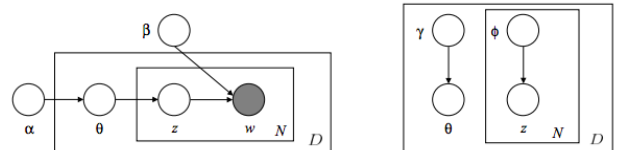


Figure 2. Graphical model corresponding to the generative process for LDA.

Now variational inference consists in solving the following

optimization problem:

$$(\gamma^*, \phi^*) = \arg \min_{\gamma, \phi} KL(q(\theta, \mathbf{z} | \gamma, \phi) || p(\theta, \mathbf{z} | \mathbf{w}, \alpha, \beta)) \quad (10)$$

Note that $q(\theta, \mathbf{z} | \mathbf{w}, \gamma^*, \phi^*)$ is actually a conditional distribution varying as a function of \mathbf{w} because the optimization objective depends on \mathbf{w} . It can thus be seen as an approximation to the posterior $p(\theta, \mathbf{z} | \mathbf{w}, \alpha, \beta)$.

With this tool at our disposal, we can come back and complete the EM algorithm described in section 1.1.2. We can expand our variational lower bound on the per-document log likelihood as follows:

$$\begin{aligned} \mathcal{L}(\gamma, \phi, \alpha, \beta) = & \mathbf{E}_q[\log p(\theta | \alpha)] + \mathbf{E}_q[\log p(\mathbf{z} | \theta)] + \\ & \mathbf{E}_q[\log p(\mathbf{w} | \mathbf{z}, \beta)] - \mathbf{E}_q[\log q(\theta | \gamma)] - \mathbf{E}_q[\log q(\mathbf{z} | \phi)] \end{aligned} \quad (11)$$

where the first three terms correspond to the expected complete log likelihood under q and the last two terms to the entropy of q .

2.1.1. E-STEP

In the E-step, we maximize \mathcal{L} with respect to the variational parameters γ and ϕ . This corresponds to applying the following update rules iteratively until convergence for each document:

$$\phi_{n,k} \propto \beta_{k,w_n} \exp\{\Psi(\gamma_k)\} \quad (12)$$

$$f\gamma_k = \alpha_k + \sum_{n=1}^N \phi_{n,k} \quad (13)$$

where $\Psi(x) = \frac{\Gamma'(x)}{\Gamma(x)}$ is the digamma function.

Theses updates make sense intuitively. Equation 12 tells us that $\phi_{n,k}$, the probability that word w_n is generated by latent topic k , is proportional to the probability of drawing w_n from topic k times a quantity proportional to how much the document likes topic k (the dirichlet parameter γ_k). Equation 13 tells us that the dirichlet parameter γ_k is equal to the pseudo-count α_k given by the dirichlet prior plus the expected count of words picking topic k .

2.1.2. M-STEP

For the M-step, we maximize the overall variational lower bound with respect to parameter β . We could also maximize with respect to α but we chose to keep it fixed to make things simpler. As said before, the overall log likelihood of the corpus is the sum of the log likelihoods for the individual documents. Likewise, the overall variational lower bound is the sum of the individual variational bounds.

Algorithm 1 Variational EM

Input: documents $\mathbf{w}_1, \dots, \mathbf{w}_D$
 Set $\alpha_k = 1/K$ for all k
 Randomly initialize variational parameters γ, ϕ
repeat
 for $d = 1$ **to** D **do**
 for $k = 1$ **to** K **do**
 for $n = 1$ **to** N_d **do**
 $\phi_{d,n,k} = \beta_{k,w_n} \exp\{\Psi(\gamma_{d,k})\}$
 end for
 normalize $\phi_{d,n}$ to sum to 1
 $\gamma_{d,k} = \alpha_k + \sum_{n=1}^{N_d} \phi_{d,n,k}$
 end for
 end for
 for $k = 1$ **to** K **do**
 for $j = 1$ **to** V **do**
 $\beta_{k,j} = \sum_{d=1}^D \sum_{n=1}^{N_d} \phi_{d,n,k} w_{d,n}^{(j)}$
 end for
 normalize β_k to sum to 1
 end for
until objective \mathcal{L} doesn't improve anymore

Maximizing with respect to β gives the following update:

$$\beta_{k,j} \propto \sum_{d=1}^D \sum_{n=1}^{N_d} \phi_{d,n,k} w_{d,n}^{(j)} \quad (14)$$

where we abuse notation and index word $w_{d,n}$ which was a scalar until here as a one hot vector.

Equation 14 also has an intuitive meaning: the probability of word j under topic k is proportional to the expected count under the variational distribution of the number of times word j is assigned topic k across the whole corpus.

2.1.3. OVERALL ALGORITHM

We implemented a variational "EM flavored" algorithm which we described in detail in the algorithm below. Note that to be strictly an EM algorithm, we should maximize \mathcal{L} with respect to ϕ and γ all the way at each E-step. Instead, we alternate a single step of improvement with respect to the variational parameters with one step of improvement with respect to β . Our algorithm has the same guarantee to increase the variational lower bound at each iteration, it is just more convenient to check for convergence only for the overall algorithm instead of for each E-step.

2.2. Gibbs Sampling

TODO

3. Experiments

TODO

4. Results

TODO