

# Variational Bayesian Inference for Unsupervised Lexicon Discovery

Elias Stengel-Eskin, BA&Sc, Honours Cognitive Science

August 24, 2017

## 1 Introduction

### 1.1 Background

Spoken language is a defining characteristic of our species, and allows us to communicate effectively with each other. This makes its study of great interest both to researchers **attempts to further understand human cognition** and developers engineering better interactive systems. By studying the computational processes **behind** speech, we can gain a better understanding of the phenomenon as a whole while developing improved technologies. Many current state-of-the-art in spoken language systems focus on supervised learning—that is, training a system on vast amount of labelled data. **The labelling process is costly and time-consuming**, meaning that sufficient data exist only for a small fraction of the world’s languages—even then, the amount of quality data can be insufficient. This often leads to the underrepresentation of many world languages and language families, particularly those spoken in the developing world. In addition, **it** does not offer a faithful model of human language capabilities. Human learning is unsupervised in the machine-learning sense; this means that we infer linguistic structure and rules implicitly. When considered together, these factors motivate an unsupervised approach which incorporates existing theories about language and cognition in order to eliminate the need for labelled training data. For speech, this is particularly useful, as an incredibly large number of speech recordings exist but the majority of these recordings are not adequately labelled to be used in supervised learning settings.

### 1.2 The Model

The specific language phenomenon we model is lexicon discovery, for which we implement an unsupervised learning algorithm. Our **framework**, which is similar to the state-of-the-art system presented in Lee et al. (2015), constructs a complete hierarchy of linguistic units (phonemes, morphemes, and words) directly from an acoustic input. The **“unsupervised lexicon discovery”** model (ULD) presented by Lee et al. (2015) was the first to jointly model the induction of phonemes, morphemes, and words, combining and extending earlier work in phoneme discovery and structural learning. ULD is composed of three main components functioning in tandem: a Dirichlet Process Hidden Markov Model (DPHMM) (Lee and Glass, 2012) for segmenting continuous audio input and hypothesizing a sequence of reusable phone-like units (PLUs), an adaptor grammar (AG) (Johnson et al., 2007) which recognizes and

stores frequently reused composite units based on an underlying grammatical structure—in this case, grouping PLUs into morphemes and words—and finally a noisy channel model, which allows substitutions, deletions, and insertions to occur between the inputs and outputs of the DPHMM and A components, approximating a phonological system. By allowing the DPHMM and A to constrain and update each other, the noisy channel is crucial to the *joint learning* aspect of the model, which has been shown to improve model accuracy (Johnson, 2008).

ULD uses nonparametric Bayesian inference to implement the unsupervised learning model. In such an inference model, we posit unobserved (latent) variables which are conditionally dependent on the data (which we observe). Defining the model in this way allows us to dynamically update our latent evidence (our hypothesis) according to the data (our evidence) while incorporating prior theoretical assumptions about the problem, yielding a powerful method for unsupervised learning. A Bayesian model is defined as follows: let  $Z$  be the set of latent variables, let  $X$  be the set of observed data, and let  $\Phi$  be a set of static model hyperparameters specified by the user. Then by Bayes' rule we have:

$$P(Z|X, \Phi) = \frac{P(X|Z, \Phi)P(Z, \Phi)}{\sum_{z \in Z} P(X|Z, \Phi)P(Z, \Phi)}$$

The numerator is often called the *generative model*, and gives a joint distribution on data and latent variables. It consists of the conditional probability of the data given our latent variables (often called the *likelihood*), which define our model, multiplied by the prior probability of our model (the *prior*). This gives us an unnormalized likelihood—generating the data from our model. However, in order to obtain a proper probability measure (something that sums to 1) we need to divide by the normalizing constant—the probability of the data. We obtain this probability by summing over all possible ways of obtaining the data—that is, all possible latent variable values. It is easy to see why, given a sufficiently complex model and a large amount of data, this sum becomes computationally intractable.

In an absence of a way to directly compute the marginal probability of the data, there are two general approaches to doing Bayesian inference. The first, used in the original ULD model, is sampling. The high-level intuition here is that, given a generative model and a way of randomly sampling from it, we can approximate the *posterior* (the left-hand side of Bayes' rule) by repeatedly taking random samples from the generative model. This technique has played a major role in Bayesian inference, but is difficult to parallelize across multiple threads and processes. This makes it nearly impossible to scale sampling algorithms to the types of large speech datasets available (Blei et al., 2017). The second method for avoiding the intractable sum required to obtain the marginal likelihood is known as variational inference. In our implementation of the ULD model, we re-implement the joint learning framework with this alternative method, which allows for parallelization and faster runtimes.

### 1.3 Variational Inference

Variational inference re-casts the challenge of computing the posterior distribution on latent variables as an optimization problem by allowing the iterative optimization of a simplified approximation of the posterior—finding itself well to parallelization across multiple cores

and interfacing with tools such as the MapReduce framework for cluster computation (Zhai et al., 2012). In order to approximate our posterior, we will introduce a family of variational distributions  $q_\nu(Z)$  which have the same support as the posterior ( $p(Z|X, \Phi)$ ), where  $\nu$  is a parameter which can be used to adjust the distribution  $q$ .

### 1.3.1 Computing the ELBO

Our ultimate goal is to find the  $q(Z)$  which minimizes the Kullback-Liebler (KL) divergence between  $q(Z)$  and  $p(Z|X)$ , or  $D_{KL}(q(Z) || p(Z|X))$ , where KL-divergence is a measure of the difference between two probability distributions. KL-divergence is given by

$$\begin{aligned} D_{KL}(q_\nu(Z) || p(Z | X, \Phi)) &= \mathbb{E}_q[\log \frac{q_\nu(Z)}{p(Z | X, \Phi)}] \\ &= \mathbb{E}_q[\log q_\nu(Z)] - \mathbb{E}_q[\log p(Z, X | \Phi)] + \log p(X | \Phi) \end{aligned} \quad (1)$$

where  $\mathbb{E}_q$  indicates taking the expected value with respect to  $q$ . Unfortunately, the second value  $\mathbb{E}_q[\log p(Z, X | \Phi)]$  again requires computing our intractable marginal probability, so we cannot directly compute KL divergence. However, this equation does give us a valuable result: a lower bound on  $\log p(X)$  called the evidence lower bound (ELBO). To obtain this, note that because of what KL divergence represents, it can never be negative. This gives us:

$$\begin{aligned} 0 &\leq \mathbb{E}_q[\log q(Z)] - \mathbb{E}_q[\log p(Z, X)] + \log p(X) \\ -\log p(X) &\leq \mathbb{E}_q[\log q(Z)] - \mathbb{E}_q[\log p(Z, X)] \\ \log p(X) &\geq \mathbb{E}_q[\log p(Z, X)] - \mathbb{E}_q[\log q(Z)] \end{aligned} \quad (2)$$

Thus

$$ELBO(q) = \mathbb{E}_q[\log p(Z, X)] - \mathbb{E}_q[\log q(Z)] \quad (3)$$

$$= \mathbb{E}_q[\log p(Z, X)] + H(q) \quad (4)$$

where  $H(q)$  is the entropy of the distribution  $q$ . This derivation yields an important fact:

$$\log p(X) - D_{KL}(q(Z) || p(Z | X)) = ELBO(q) \quad (5)$$

This explains why maximizing the ELBO allows us to minimize the KL divergence—the maximal ELBO is  $\log p(X)$ —they are equal, the KL-divergence must be 0 (Blei et al., 2017). For an expanded derivation, as well as another equivalent, see Appendix A.

### 1.3.2 Mean Field Approximation

One of the main reasons why we cannot compute the posterior directly is the presence of conditional dependencies between latent variables in it. One way to approximate the posterior with a simpler distribution is to make a mean field assumption—that is, to assume that the variational distribution  $q$  has none of these conditional dependencies, i.e.

$$q(Z) = \prod_{z \in Z} q(z)$$

This is a very powerful assumption, because it allows us to optimize each variational distribution iteratively. That is to say, while holding all other variational distributions constant, we find the variational parameters for  $q_{\nu_i}(z_i)$  that maximize the marginal likelihood. Through using the chain rule, we can derive the following lower bound for each variational distribution:

$$\mathcal{L}_i = \mathbb{E}_q[\log p(z_i|Z_{-i}, X, \Phi)] - \mathbb{E}_q[\log q_{\nu_i}(z_i)] \quad (6)$$

where  $Z_{-i}$  indicates all latent variables in  $Z$  which are not  $z_i$ .

### 1.3.3 Updates

Recall that our goal is to adjust each  $\nu_i$  in order to maximize our lower bound  $\mathcal{L}_i$  for each latent variable. Finding the value for  $\nu_i$  that locally maximizes this function can be done by setting the first derivative with respect to  $\nu_i$  equal to 0 and solving for  $\nu_i$ . Taking the derivative of the objective function can be costly, especially since it must be done every time we want to update our parameters and for every parameter in the factorized representation of the variational distribution. Using exponential family random variables allow us to leverage some convenient mathematical facts and avoid this computation. When each  $q_i$  and each distribution in the generative model are in the exponential family, we obtain the following closed-form update for each  $\nu_i$ :

$$\nu_i = \mathbb{E}_q[g_i(Z_{-i}, X, \Phi)] = \mathbb{E}_q \left[ \frac{\sum_{z_n \in Z_{-i}} t(x_n, z_n)}{\phi_2 + N} \right] \quad (7)$$

where  $g_i(Z_{-i}, X, \Phi)$  is a function which gives the natural parameters of the exponential family distribution in the posterior,  $\phi_1$  and  $\phi_2$  are the parameters for the exponential family distribution in the prior, and  $t(x_n, z_n)$  is the sufficient statistic for the prior distribution—in many cases, this is simply the count of occurrences of  $z_n$ . For a more in-depth explanation of this result, see Appendix B.

### 1.3.4 Coordinate Ascent

We now have a way of optimizing each variational distribution by setting the parameters to the expected value of the natural parameters in the posterior, conditioned on the other latent variables and the data. If our objective function could be formulated as a strictly convex function, then one update of the variational parameters would be sufficient to find a solution, since a local optimum in a convex function is a global one. However, given the complex nature of most Bayesian models, this will rarely be the case. Instead, we use Coordinate Ascent Variational Inference (CAVI) in order to iteratively find local maxima, with the ultimate goal of converging on the global maximum. It is worth noting that the CAVI algorithm is a generalization of the well-known Expectation-Maximization algorithm (Dempster et al., 1977), but while the latter gives a point estimate of the posterior, the former returns an approximation of the full distribution. In the algorithm we will alternate between computing our objective function (analogous to the expectation step) and updating our variational parameters (analogous to the maximization step) (Neal and Hinton, 1998).

```

initialize each  $\nu_i$ 
while not ELBO converged do
  for each variational parameter  $\nu_i$  do
     $\nu_i = \mathbb{E}_q[g_i(Z_{-i}, X, \Phi)]$ 
  re-compute ELBO  $\mathcal{L}(q) = \mathbb{E}_q[\log p(Z, X)] + H(q)$ 

```

Figure 1: The CAVI algorithm

The initialization step for each  $\nu_i$  can be random, but this is not required. Often, we let  $q_{\nu_i}(z_i)$  be a distribution of the same type as  $p(z_i)$  in the generative model, and initialize it with uniform or random parameters. However, we can choose the initial parameters more deliberately and encode some bias in the variational distribution, with the caveat that different initializations can lead to convergence on different local optima (Blei et al., 2017)

## 2 Model Definition

The model can be broken up into roughly three parts: the adaptor grammar, the noisy channel, and the Dirichlet-process hidden Markov model. From the perspective of the generative model, the adaptor grammar builds a syntactic tree whose yield is a set of top-level phone-like units (PLUs). The tree groups these into morphemes and words. The noisy channel, using a sequence of edit operations (insertion, deletion, and substitution) maps these top-level PLUs to bottom-level PLUs, modeling some of the phonological processes which occur during speech production. Finally, the Dirichlet-process hidden Markov model takes these bottom-level PLUs and finds an acoustic signal that could have generated them.

### 2.1 Adaptor grammars

First developed by (Goldman, 1999), adaptor grammars take as input some context-free grammar and some strings which can be parsed by that grammar. Using a non-parametric distribution, they store parse trees while biasing the reuse of frequently occurring trees which reveals patterns in the linguistic structure of the data. By increasing the likelihood of reusing a tree according to its frequency, adaptor grammars instantiate a “rich get richer” dynamic where common trees become more likely than uncommon ones. In ULD, we use adaptor grammars to group discovered phones into morphemes and words. Before formally defining adaptor grammars, we need to define context free grammars, probabilistic context free grammars, and the Pitman-Yor Process.

#### 2.1.1 Context-free Grammars

A context-free grammar (CFG) is a tuple  $(N, E, R, S)$  where  $N$  is a set of nonterminals symbols,  $E$  is a set of terminal symbols disjoint from  $N$ ,  $R$  is a set of rules of the form  $A \rightarrow \beta$  where  $A \in N$  and  $\beta \in (N \cup E)^*$  (i.e. any concatenation of symbols in  $N$  and  $E$ ). For this implementation, we constrain our CFGs to be in Chomsky normal form, where

all the rules are either of form  $A \rightarrow a$  where  $a \in E$  or  $A \rightarrow BC$  where  $B \in N \cup E$  and  $C \in N \cup E$ . Note that any CFG without epsilon productions (rules that go to the empty string) can be rewritten in Chomsky normal form. (Hill III, 1979)

### 2.1.2 Probabilistic Context-free Grammars

Similarly to a CFG a probabilistic context-free grammar (PCFG) is a tuple  $(N, E, R, S, \theta)$  where  $N, E, R, S$  are the same as in a CFG, and  $\theta$  is a set of probability vectors such that  $\sum_{A \rightarrow \beta \in R_A} \theta_{A \rightarrow \beta} = 1$ , where  $R_A$  is the set of rules which have nonterminal  $A$  on the left-hand side.

### 2.1.3 Pitman-Yor Process

The Pitman-Yor process (Pitman and Yor, 1997) can be thought of as a distribution on infinite-sided dice, or as generating a partition of integers. Perhaps more intuitively, a Pitman-Yor process defines a distribution over distributions—each draw from a Pitman-Yor process is itself an infinite distribution. There are multiple equivalent ways of defining a Pitman-Yor process—we will be using the stick-breaking construction, as it gives us an iterative definition. Given a scale parameter  $a$ , a discount factor  $b$  and a base distribution  $G_0$ , a Pitman-Yor process which partitions  $[0, 1]$  into countably infinite segments is defined by the following algorithm:

---

```

for  $i \in \{1, \dots\}$  do
    draw  $\nu_i \sim \text{Beta}(1 - b, a + ib)$ 
    sample atom  $z_i \sim G_0$ 
    define  $\pi_i \triangleq \nu_i \prod_{j=1}^{i-1} (1 - \nu_j)$ 
define  $G(z) \triangleq \sum_{i=1}^{\infty} \pi_i \delta(z_i, z)$  where  $\delta(z_i, z) = 1$  if  $z_i = z$ , 0 otherwise

```

---

Recall that a draw from a Beta distribution is a biased coin. Intuitively, each  $\nu_i$  is a coin that gives the probability of stopping at that stick, and  $1 - \nu_i$  is the probability of continuing to the next stick. Then  $\pi_i$ , the probability of being at stick  $i$ , is equivalent to the product of the probability of having passed sticks  $1, \dots, i - 1$  and the probability of stopping at stick  $i$ .

### 2.1.4 Adaptor grammar definition

With these components, we can formally define an adaptor grammar as a tuple  $(G, M, a, b, \alpha)$  where  $G$  is a CFG,  $M$  is a set of *adapted nonterminals*,  $a$  and  $b$  are Pitman-Yor process parameters, and  $\alpha$  is a set of Dirichlet distribution parameters indexed by each nonterminal in  $N$ . Let  $A_1, \dots, A_k$  be a reverse topological sorting of the adapted nonterminals in  $M$ , such that for all  $A_j \in M$  the children of  $A_j$  come before  $A_j$  in the ordering. The following two algorithms define an adaptor grammar.

---

**Result:** Building the grammar

*constructing the PCFG*

**foreach** *nonterminal*  $A$  **in**  $N$  **do**

draw rule weights  $\theta_A \sim \text{Dir}(\alpha_A)$  ;

*constructing the grammatons*  $G_A$

**for**  $A \in A_1, \dots, A_k$  **do**

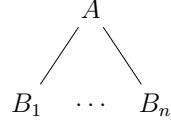
draw  $\pi_A \sim \text{PYP}(a_A, b_A)$

*construct tree*  $z_{A,i}$

**for**  $i \in \{1, \dots\}$  **do**

draw rule  $A \rightarrow B_1 \dots B_n$  from  $R_A$

set  $z_{A,i} =$



**while**  $z_{A,i}$  *has nonterminals in leaves* **do**

choose a  $B$  from  $B_1 \dots B_n$

**if**  $B$  *is non-adapted nonterminal* **then**

└ expand  $B$  using the PCFG

**else**

└ expand  $B$  using  $G_B$

└ *guaranteed to exist because of topological ordering*

**for**  $i \in \{1, \dots\}$  **do**

└ set  $G_A(z_{A,i}) = \pi_{A_i}$

**Result:** generating data

*generating derivation trees*  $z_i$

**for**  $i \in \{1, \dots, |X|\}$  **do**

draw  $\pi_A \sim \text{PYP}(a_A, b_A)$

*construct tree*  $z_{A,i}$

**for**  $i \in \{1, \dots\}$  **do**

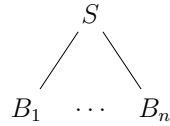
**if**  $S$  *is adapted nonterminal* **then**

└ draw  $z_i \sim G_S$

**else**

└ draw  $S \rightarrow B_1 \dots B_n$  from  $R_S$

set  $z_i =$



**while**  $z_i$  *has nonterminals in leaves* **do**

choose a  $B$  from  $B_1 \dots B_n$

**if**  $B$  *is non-adapted nonterminal* **then**

└ expand  $B$  using the PCFG

**else**

└ expand  $B$  using  $G_B$


└ set  $x_i$  to the yield of tree  $z_i$

---

Figure 2: Algorithm 2: Defining an adaptor grammar

This definition of adaptor grammars gives us the following latent variables:

- $z_i$ : the full derivational trees that yielded the data.
- $z_{A,i}$ : the stored sub-trees headed by adapted non-terminals
- $\nu$ : the set of stick-weight proportions for the Pitman-Yor process
- $\theta$ : the set of PCFG rule probabilities

Our inference problem can be formalized as finding the posterior distribution on full derivational trees  $z_i$ —these depend on all the other latent variables, and reveal the inferred underlying linguistic structure. However, this inference is over an extremely large set of latent variables. In fact, in the current formalization, we cannot do inference over this set of latent variables, since some are countably infinite. To make this problem finite, we use a truncated stick-breaking representation, where after a sufficiently large  $i$  we let  $\nu_i = 1$ , so that the probability of continuing past that stick is 0. Beyond the large number of latent variables, we need to take into account the large number of potential parses for each sentence given the grammar. Indeed, averaging rule probabilities over all of these parses  be the most costly portion of the algorithm (see section TODO: INSERT I/O section).

## 2.2 Dirichlet-process hidden Markov model



# A Deriving the ELBO

## A.1 The problem

Given our generative model and our data, we would like to find a posterior distribution:  $P(Z | X)$ . Using Bayes Rule, we get:  $P(Z | X) = \frac{P(X|Z)P(Z)}{P(X)} = \frac{P(X|Z)P(Z)}{\sum_{\forall Z} P(X|Z)P(Z)}$  We call the numerator of the fraction on the right the generative model. It is composed of the product of the likelihood \*of the hypothesis\* ( $P(X | Z)$ ) and the prior probability of the hypothesis ( $P(Z)$ ). Note that the former is not a probability but a measure of how well our hypothesis fits the data. The denominator is the “marginal likelihood” of the data,  $P(X)$ . To find this, we need to marginalize out (sum over) all possible hypotheses. Because the hypotheses are the range of values for all of the latent variables in our model, this summation is computationally intractable. In order to obtain a posterior, we are forced to use some approximate means of finding this denominator. Often, a sampling approach is used. However, sampling can be very slow to converge and is not easily parallelizable across multiple cores. The variational Bayesian approach, on the other hand, treats the problem of finding an appropriate marginal distribution as an optimization problem.

- Let  $Z$  be our set of hidden variable collections:
- Let  $\Phi$  be the collection of all model parameters (Pitman-Yor parameters  $a, b$  and Dirichlet distribution parameter  $\alpha$ ).
- Let  $X$  be the set of observations. In the case of word segmentation, for example, these would be each string of unsegmented phonemes.
- Note that our goal is to find  $P(Z | X)$ , the posterior (where  $Z$  is the set of latent variables)
- recall  $P(Z | X) = \frac{P(X|Z,\Phi)P(Z|\Phi)}{\sum_{\forall Z} P(X|Z,\Phi)P(Z|\Phi)}$

## A.2 Important formulae

### A.2.1 Jensen’s inequality

Jensen’s inequality states that for a convex function  $f$  and random variable  $X$ :

$$f(\mathbb{E}[X]) \leq \mathbb{E}[f(X)] \quad (8)$$

We are using the logarithm of the probability, so the function is actually concave. Jensen’s inequality works both ways, meaning we switch the direction of the inequality:

$$\log(\mathbb{E}[X]) \geq \mathbb{E}[\log(X)] \quad (9)$$

### A.2.2 Expected value

Note that for discrete random variables (which we are dealing with in this case)

$$\mathbb{E}_q(f(x)) = \sum_{\forall x} q(x)f(x) \quad (10)$$


### A.2.3 Logarithms

Throughout this derivation (and the variational literature as a whole) the logarithm of the probability is used. There are various reasons to do this. Firstly, taking the logarithm allows us to use information-theoretic measures such as entropy. Furthermore, logarithms allow us to transform expensive multiplication and division into cheaper addition and subtraction, and helps when working with probabilities below the floating-point precision bound. Recall some facts about logarithms:

- $\lim_{n \rightarrow 0} \log n = -\infty$
- $\log AB = \log A + \log B$
- $\log \frac{A}{B} = \log A - \log B$

### A.3 Derivation of variational bound

The value we are looking to approximate is our denominator, which is the likelihood of the data integrated over all hypotheses. Recall that our inference problem lies in finding the denominator to the Bayesian equation  $P(Z | X) = \frac{P(X|Z)P(Z)}{\int_{\forall Z} P(X|Z)P(Z)dZ}$ . Our hypotheses in this case are possible values for the latent variables in the model. This integral (or in the discrete case, summation) is often computationally intractable, so we use a variational approximation for it. One way we can do this is by using the Kullback Leibler (KL) divergence between this intractable integral and some variational distribution  $q$ .

1. Let  $q_\nu(Z)$  be a family of variational distributions with variational parameter  $\nu$ .
2. to get the marginal likelihood ( $\log p(X | \Phi)$ ) we  take the KL divergence between  $q_\nu(Z)$  and  $p(Z | X, \Phi)$ .
3. KL divergence is given by:

$$\begin{aligned}
 D_{KL}(q_\nu(Z) || p(Z | X, \Phi)) &= \mathbb{E}_q[\log \frac{q_\nu(Z)}{p(Z | X, \Phi)}] \\
 &= \mathbb{E}_q[\log q_\nu(Z) - \log p(Z | X, \Phi)] \\
 &= \mathbb{E}_q[\log q_\nu(Z) - \log \frac{p(Z, X | \Phi)}{p(X | \Phi)}] \\
 &= \mathbb{E}_q[\log q_\nu(Z) - (\log p(Z, X | \Phi) - \log p(X | \Phi))] \\
 &= \mathbb{E}_q[\log q_\nu(Z)] - \mathbb{E}_q[\log p(Z, X | \Phi)] + \log p(X | \Phi)
 \end{aligned} \tag{11}$$

(Blei and Jordan, 2006)

If we think about what KL divergence represents, we can intuitively understand why it cannot be negative. From here, we can see how minimizing this equation is the same as maximizing the lower bound on  $\log p(X | \Phi)$ :

$$\begin{aligned}
0 &\leq \mathbb{E}_q[\log q_\nu(Z)] - \mathbb{E}_q[\log p(Z, X | \Phi)] + \log p(X | \Phi) \\
&\quad - \log p(X | \Phi) \leq \mathbb{E}_q[\log q_\nu(Z)] - \mathbb{E}_q[\log p(Z, X | \Phi)] \\
&\quad \log p(X | \Phi) \geq \mathbb{E}_q[\log p(Z, X | \Phi)] - \mathbb{E}_q[\log q_\nu(Z)]
\end{aligned} \tag{12}$$

Another way to reach this same equation is by using Jensen's inequality. Consider the log marginal likelihood:

$$\log p(X | \Phi) = \log \sum_{z \in \mathbf{Z}} p(X, z | \Phi) \tag{13}$$

The sum marginalizes out the hidden variables  $z$  in the joint probability distribution. Picking any variational distribution  $q(z)$  we can multiply them by  $\frac{q(z)}{q(z)}$ :

$$\log \sum_{\forall z \in \mathbf{Z}} (p(x, z | \Phi) * \frac{q(z)}{q(z)}) = \log \sum_{\forall z \in \mathbf{Z}} q(z) \frac{p(x, z | \Phi)}{q(z)} \tag{14}$$

Jensen's inequality implies

$$\log \sum_{\forall z \in \mathbf{Z}} q(z) \frac{p(x, z | \Phi)}{q(z)} \geq \sum_{\forall z \in \mathbf{Z}} q(z) \log \frac{p(x, z | \Phi)}{q(z)} \tag{15}$$

This equation can be broken into:

$$\begin{aligned}
\sum_{\forall z \in \mathbf{Z}} q(z) \log \frac{p(x, z | \Phi)}{q(z)} &= \sum_{\forall z \in \mathbf{Z}} q(z) (\log p(x, z | \Phi) - \log q(z)) = \\
&\quad \sum_{\forall z \in \mathbf{Z}} q(z) \log p(x, z | \Phi) - \sum_{\forall z \in \mathbf{Z}} q(z) \log q(z) = \\
&\quad \sum_{\forall z \in \mathbf{Z}} q(z) \log p(x, z | \Phi) + \mathcal{H}(q)
\end{aligned} \tag{16}$$

$$\tag{17}$$

where

$$\mathcal{H}(q) = - \sum_{\forall z \in \mathbf{Z}} q(z) \log q(z) \tag{18}$$

(Blei et al., 2017) This first term is of the form of our expected value definition, so our equation becomes:

$$\log p(x | \Phi) \geq \mathbb{E}_q[\log p(x, z | \Phi)] + \mathcal{H}(q) \tag{19}$$

This derivation yields an important fact:

$$\log p(X | \Phi) - KL(q(Z) || p(Z | X, \Phi)) = \mathbb{E}_q[\log p(z, x | \Phi)] + H(q) \quad (20)$$

From this equation, we can see why minimizing KL divergence gives us the best possible value for our marginal likelihood.

## B Deriving Variational Updates

### B.1 Mean Field Approximation

Recall our mean-field assumption to treat each variational distribution as conditionally independent:  $q(Z) = \prod_i q_i(z_i)$ ; also recall that our bound on the log marginal likelihood was:

$$L(q) \geq \sum_{z_i \in Z} q(Z) \log p(X, Z | \Phi) + H(q) \quad (21)$$

For the sake of concise notation, let  $q_j = q_{\nu_j}(Z_j)$  and disregard the conditioner  $\Phi$  for now. Replace  $q(Z)$  with this approximating product:

$$\begin{aligned} L(q) &\geq \sum_{z_i \in Z} \left( \prod_i q_i(z_i) \right) \log p(X, Z | \Phi) + H(q) \\ L(q) &\geq \mathbb{E}_{\prod_i q_i(z_i)} \log p(X, Z | \Phi) [\log p(X, Z | \Phi)] + H(q) \end{aligned} \quad (22)$$

Using the chain rule and by expanding the entropy term, we can rewrite this expression as

$$\log p(X | \Phi) + \sum_{i=1}^{|Z|} \mathbb{E}_q[\log p(z_i | X, z_1, \dots, z_{i-1}, \Phi)] - \sum_{i=1}^{|Z|} \mathbb{E}_q[\log q_{\nu_i}(z_i)] \quad (23)$$

Since  $p(X | \Phi)$  does not depend on the variational parameter  $\nu_i$  it factors out as a constant (recall that this is a lower bound, not an exact equality). We can reorder the elements of  $Z$  in any way we wish. If we reorder them each time so that  $z_i$  comes last, we can say:

$$\uparrow_i = \mathbb{E}_q[\log p(z_i | Z_{-i}, X, \Phi)] - \mathbb{E}_q[\log q_{\nu_i}(z_i)] \quad (24)$$

(Blei and Jordan, 2006)

Note that for any exponential family distribution  $q_{\nu_i}$ ,

$$q_{\nu_i}(z_i) = h(z_i) \exp \{ \nu_i^T z_i - a(\nu_i) \} \quad (25)$$

where  $a(\nu_i)$  is the cumulant function, which for the first three derivatives is equivalent to the corresponding derivatives of the moment generating function. We can rewrite our equation using this form for  $q_{\nu_i}(z_i)$ :

$$\begin{aligned}
\Downarrow_i &= \mathbb{E}_q[\log p(z_i|Z_{-i}, X, \Phi)] - \mathbb{E}_q\left[\log \left(h(z_i) \exp \{\nu_i^T z_i - a(\nu_i)\}\right)\right] \\
&= \mathbb{E}_q[\log p(z_i|Z_{-i}, X, \Phi)] - \mathbb{E}_q\left[\log (h(z_i)) + \nu_i^T z_i - a(\nu_i)\right] \\
&= \mathbb{E}_q[\log p(z_i|Z_{-i}, X, \Phi)] - \mathbb{E}_q[\log (h(z_i))] - \mathbb{E}_q[\nu_i^T z_i] + \mathbb{E}_q[a(\nu_i)] \\
&= \mathbb{E}_q[\log p(z_i|Z_{-i}, X, \Phi)] - \mathbb{E}_q[\log (h(z_i))] - \nu_i^T a'(\nu_i) + a(\nu_i)
\end{aligned} \tag{26}$$

Note that  $\mathbb{E}_q[\nu_i^T z_i] = \nu_i^T a'(\nu_i)$  since  $E_q(z_i) = a'(\nu_i)$  and  $\nu_i^T$  factors out as a constant when taking the expectation with respect to  $q$ . The main premise of variational inference is to cast the intractable calculation of the posterior as an optimization problem. In most optimization problems, there are two general steps:

1. computing an objective function which will allow us to
2. optimize the function by adjusting the parameters.

Recall that to avoid expensive computations, we use exponential family distributions which allow us to simplify the problem. We are trying to optimize the function by adjusting the variational parameters, so we take the partial derivative of our function with respect to  $\nu_i$ :

$$\begin{aligned}
\frac{\delta}{\delta \nu_i} \Downarrow_i &= \frac{\delta}{\delta \nu_i} (\mathbb{E}_q[\log p(z_i|Z_{-i}, X, \Phi)] - \mathbb{E}_q[\log (h(z_i))] - \nu_i^T a'(\nu_i) + a(\nu_i)) \\
&= \frac{\delta}{\delta \nu_i} (\mathbb{E}_q[\log p(z_i|Z_{-i}, X, \Phi)] - \mathbb{E}_q[\log h(z_i)]) - (\nu_i^T a''(\nu_i) + a''(\nu_i)) + a''(\nu_i) \\
&= \frac{\delta}{\delta \nu_i} (\mathbb{E}_q[\log p(z_i|Z_{-i}, X, \Phi)] - \mathbb{E}_q[\log h(z_i)]) - \nu_i^T a''(\nu_i)
\end{aligned} \tag{27}$$

Setting this to 0 we get:

$$\begin{aligned}
0 &= \frac{\delta}{\delta \nu_i} (\mathbb{E}_q[\log p(z_i|Z_{-i}, X, \Phi)] - \mathbb{E}_q[\log h(z_i)]) - \nu_i^T a''(\nu_i) \\
\nu_i^T a''(\nu_i) &= \frac{\delta}{\delta \nu_i} (\mathbb{E}_q[\log p(z_i|Z_{-i}, X, \Phi)] - \mathbb{E}_q[\log h(z_i)]) \\
\nu_i &= \left( \frac{\delta}{\delta \nu_i} \mathbb{E}_q[\log p(z_i|Z_{-i}, X, \Phi)] - \frac{\delta}{\delta \nu_i} \mathbb{E}_q[\log h(z_i)] \right) (a''(\nu_i))^{-1}
\end{aligned} \tag{28}$$

If  $p(z_i|Z_{-i}, X, \Phi)$  is also a member of the exponential family, it can be rewritten:

$$p(z_i|Z_{-i}, X, \Phi) = h(z_i) \exp \{g_i(Z_{-i}, X, \Phi)^T z_i - a(g_i(Z_{-i}, X, \Phi))\} \tag{29}$$

where  $g_i(Z_{-i}, X, \Phi)$  is the natural parameter of distribution  $p$ . Replacing  $p(z_i|Z_{-i}, X, \Phi)$  (first in the expected values for the sake of readability) and taking the derivative gives us

$$\mathbb{E}_q[\log p(z_i|Z_{-i}, X, \Phi)] = \mathbb{E}_q[\log h(z_i)] + \mathbb{E}_q[g_i(Z_{-i}, X, \Phi)]^T a'(\nu_i) - \mathbb{E}_q[a(g_i(Z_{-i}, X, \Phi))] \quad (30)$$

$$\begin{aligned} \frac{\delta}{\delta \nu_i} \mathbb{E}_q[p(z_i|Z_{-i}, X, \Phi)] &= \frac{\delta}{\delta \nu_i} \mathbb{E}_q[\log h(z_i)] \\ &+ \left( \frac{\delta}{\delta \nu_i} (\mathbb{E}_q[g_i(Z_{-i}, X, \Phi)]^T) a'(\nu_i) + \mathbb{E}_q[g_i(Z_{-i}, X, \Phi)]^T a''(\nu_i) \right) - \frac{\delta}{\delta \nu_i} \mathbb{E}_q[a(g_i(Z_{-i}, X, \Phi))] \\ &= \frac{\delta}{\delta \nu_i} \mathbb{E}_q[\log h(z_i)] + \mathbb{E}_q[g_i(Z_{-i}, X, \Phi)]^T a''(\nu_i) \end{aligned} \quad (31)$$

Notice that many of the expectations drop out. Substituting this for  $\frac{\delta}{\delta \nu_i} (\mathbb{E}_q[\log p(z_i|Z_{-i}, X, \Phi)])$  in our first differentiation, we get:

$$\begin{aligned} \nu_i &= \left( \frac{\delta}{\delta \nu_i} \mathbb{E}_q[\log h(z_i)] + \mathbb{E}_q[g_i(Z_{-i}, X, \Phi)]^T a''(\nu_i) - \frac{\delta}{\delta \nu_i} \mathbb{E}_q[\log h(z_i)] \right) (a''(\nu_i))^{-1} \\ &= (\mathbb{E}_q[g_i(Z_{-i}, X, \Phi)]^T a''(\nu_i)) (a''(\nu_i))^{-1} \\ &= \mathbb{E}_q[g_i(Z_{-i}, X, \Phi)] \end{aligned} \quad (32)$$

So the optimal value (when the derivative is 0) of  $\nu_i$  is  $\nu_i = \mathbb{E}_q[g_i(Z_{-i}, X, \Phi)]$ .

Now we need to obtain a closed-form expression for  $\mathbb{E}_q[g_i(Z_{-i}, X, \Phi)]$ . Firstly, let  $\Phi$  be composed of 2 parts  $\phi_1$  and  $\phi_2$ , where  $\phi_1$  is the number of observations contributed by the prior, and  $\phi_2$  corresponds to the total effect of the observations on the sufficient statistic. Because of the factorization and exponential family assumptions we made earlier, we can say:

$$\begin{aligned} P_\pi(z_i|\phi_1, \phi_2) &= f(\phi_1, \phi_2) \exp \{ \eta^T \phi_1 - \phi_2 a(\eta) \} \\ &= f(\phi_1, \phi_2) g(\eta)^{\phi_2} \exp \{ \eta^T \phi_1 \} \\ &\propto g(\eta)^{\phi_2} \exp \{ \eta^T \phi_1 \} \end{aligned} \quad (33)$$

where  $\eta$  are the natural parameters for the distribution,  $a(\eta)$  is the cumulant function, and  $f(\phi_1, \phi_2)$  is a normalizing function.

Assuming the posterior over data and local hidden variables  $P(X, Z_{-i} | z_i)$  is also in the exponential family and factorizes, we can say that for one data point  $x_n$

$$\begin{aligned} P(x_n, z_n | z_i) &= h(x_n, z_n) g(z_i) \exp \{ z_i^T t(x_n, z_n) \} \\ \Rightarrow P(X, Z_{-i} | z_i) &= \prod_{z_n \in Z_{-i}} h(x_n, z_n) g(z_i)^N \exp \{ z_i^T t(x_n, z_n) \} \end{aligned} \quad (34)$$

where  $t(x_n, z_n)$  is the sufficient statistic, which in most cases is simply the count of occurrences of  $x_n$  or  $z_n$ . By Bayes rule, the distribution over the selected hidden variable rewrites as

$$\begin{aligned}
P(z_i|X, Z_{-i}, \Phi) &\propto P(X, Z_{-i}|z_i)P(z_i|\Phi) \\
&= \prod_{n=0}^N h(x_n, z_n)g(z_i) \exp \{z_i^T t(x_n, z_n)\} g(\eta)^{\phi_2} \exp \{\eta^T \phi_1\} \\
&\propto g(z_i)^N \exp \{z_i^T t(x_n, z_n)\} g(\eta)^{\phi_2} \exp \{\eta^T \phi_1\} \\
&\propto g(z_i)^{N+\phi_2} \exp \left\{ z_i^T \left( \phi_1 + \sum_{z_n \in Z_{-i}} t(x_n, z_n) \right) \right\}
\end{aligned} \tag{35}$$

Because all exponential family distributions have conjugate priors, this result implies that the posterior  $P(z_i | X, Z_{-i}, \Phi)$  is the same type of distribution as the prior with parameters:

$$P(z_i|X, Z_{-i}, \Phi) = P_\pi(z_i|\phi_2 + \sum_{z_n \in Z_{-i}} t(x_n, z_n), \phi_2 + N) \tag{36}$$

This means that the natural parameters of the posterior distribution on global hidden variables has the natural parameters  $\phi_2 + \sum_{z_n \in Z_{-i}} t(x_n, z_n)$  and  $\phi_2 + N$ , giving us a closed form for our expectation in equation TODO: add final equation number here:

$$\mathbb{E}_q[g_i(Z_{-i}, X, \Phi)] = \mathbb{E}_q \left[ \frac{\phi_2 + \sum_{z_n \in Z_{-i}} t(x_n, z_n)}{\phi_2 + N} \right] \tag{37}$$

(Hoffman et al., 2013)



## References

- Blei, D. M. and Jordan, M. I. (2006). Variational inference for Dirichlet process mixtures. *Bayesian Analysis*, 1(1):121–143.
- Blei, D. M., Kucukelbir, A., and McAuliffe, J. D. (2017). Variational inference: A review for statisticians. *Journal of the American Statistical Association*, (just-accepted).
- Dempster, A. P., Laird, N. M., and Rubin, D. B. (1977). Maximum likelihood from incomplete data via the EM algorithm. *Journal of the royal statistical society. Series B (methodological)*, pages 1–38.
- Hill III, I. (1979). Introduction to automata theory, languages, and computation.
- Hoffman, M. D., Blei, D. M., Wang, C., and Paisley, J. (2013). Stochastic variational inference. *The Journal of Machine Learning Research*, 14(1):1303–1347.
- Johnson, M. (2008). Using Adaptor Grammars to Identify Synergies in the Unsupervised Acquisition of Linguistic Structure. In *ACL*, pages 398–406.
- Johnson, M., Griffiths, T. L., and Goldwater, S. (2007). Adaptor grammars: A framework for specifying compositional nonparametric Bayesian models. In *Advances in neural information processing systems*, pages 641–648.
- Lee, C.-Y. and Glass, J. (2012). A nonparametric Bayesian approach to acoustic model discovery. In *Proceedings of the 50th Annual Meeting of the Association for Computational Linguistics: Long Papers-Volume 1*, pages 40–49. Association for Computational Linguistics.
- Lee, C.-Y., O’Donnell, T. J., and Glass, J. (2015). Unsupervised lexicon discovery from acoustic input. *Transactions of the Association for Computational Linguistics*, 3:389–403.
- Neal, R. M. and Hinton, G. E. (1998). A view of the EM algorithm that justifies incremental, sparse, and other variants. In *Learning in graphical models*, pages 355–368. Springer.
- Pitman, J. and Yor, M. (1997). The two-parameter Poisson-Dirichlet distribution derived from a stable subordinator. *The Annals of Probability*, pages 855–900.
- Zhai, K., Boyd-Graber, J., Asadi, N., and Alkhrouja, M. L. (2012). Mr. LDA: A flexible large scale topic modeling package using variational inference in MapReduce. In *Proceedings of the 21st international conference on World Wide Web*, pages 879–888. ACM.