

PGM - Variational methods for sigmoid belief networks

Caio Corro

Contents

1	Introduction	1
2	Two layers sigmoid belief networks (SBN)	1
2.1	Lower bound on the evidence	2
3	Upper bound on the evidence	4
4	Expectation-Maximization algorithm	5

1 Introduction

Computing the probability of an observation requires to marginalize out latent (i.e. non observed) variables (if any). For, many models, this marginalization is too expensive to be computed exactly. There exists two standard approximation methods: stochastic approximations and variational approximations. These notes present variational approximation techniques for two layers sigmoid belief networks that were introduced in [Saul et al., 1996] and [Jaakkola and Jordan, 1996].

The main idea is as follows. Assume a joint probability distribution $p(\mathcal{X}, \mathcal{Y})$ where \mathcal{X} and \mathcal{Y} are the observed and latent random variables, respectively. We wish to estimate the evidence (log-probability) of an observation \mathbf{x} :

$$\log(\mathcal{X} = \mathbf{x}) = \log \mathbb{E}_{p(\mathcal{Y})}[p(\mathcal{X} = \mathbf{x}|\mathcal{Y})]$$

This quantity can be useful if we want to estimate the probability that a given point comes from the distribution or to learn the parameters of the distribution by maximizing the log-likelihood of a dataset. In many cases, the expectation (i.e. the marginalization) leads to an intractable evidence. Therefore, we may want to approximate it using parameterized functions:

$$\begin{aligned} \log(\mathcal{X} = \mathbf{x}) &\geq f(\mathbf{x}, \phi) \\ \text{or } \log(\mathcal{X} = \mathbf{x}) &\leq h(\mathbf{x}, \epsilon), \end{aligned}$$

where we assume that computing $f(\mathbf{x}, \phi)$ and $h(\mathbf{x}, \epsilon)$ is easy. The parameters ϕ and ϵ are variational parameters that control the quality of the bounds, i.e. we assume that these bounds are true for any ϕ and ϵ . That is, we replaced an intractable expectation with an easy to compute parameterized function. Note that, unfortunately, finding the best possible variational parameters can be a difficult problem.

2 Two layers sigmoid belief networks (SBN)

Let \mathcal{X} be the observed random variable, taking values in $[0, 1]^d$, and \mathcal{Y} the latent random variable, taking values in $[0, 1]^k$. A two layers sigmoid belief networks (SBN) defines a joint probability distribution

on \mathcal{X} and \mathcal{Y} as follows:

$$\begin{aligned} p_\theta(\mathbf{x}, \mathbf{y}) &= p_\theta(\mathbf{y})p_\theta(\mathbf{x}|\mathbf{y}) \\ &= \prod_{i=1}^k p_\theta(y_i) \prod_{i=1}^d p_\theta(x_i|\mathbf{y}) \\ &= \prod_{i=1}^k \frac{\exp(y_i a_i)}{1 + \exp(a_i)} \prod_{i=1}^d \frac{\exp(x_i(\mathbf{B}_i \mathbf{y} + c_i))}{1 + \exp(\mathbf{B}_i \mathbf{y} + c_i)} \end{aligned}$$

where $\theta = \{\mathbf{a}, \mathbf{B}, \mathbf{c}\}$ are the parameters of the model with $\mathbf{a} \in \mathbb{R}^k$, $\mathbf{B} \in \mathbb{R}^{d \times k}$ and $\mathbf{c} \in \mathbb{R}^d$. In this model, the evidence is intractable as it requires to sum of 2^k possible assignments for the latent variable \mathcal{Y} .

Generating data from a trained model is easy via ancestral sampling: first sample from the latent variable distribution (independent Bernoullis) and then sample from the conditional distribution of observed variables (independent Bernoullis again). The generative story can be described as follows:

1. $\mathbf{y} \sim p_\theta(\mathcal{Y})$
2. $\mathbf{x} \sim p_\theta(\mathcal{X}|\mathcal{Y} = \mathbf{y})$

2.1 Lower bound on the evidence

In this section, we derive a lower bound on the evidence based on the introduction of a proposal distribution and Jensen’s inequality. The resulting lower bound is not equivalent to the “standard” ELBO, as one of the term of the ELBO is intractable for SBNs. Contrary to Gaussian Mixture Models, the posterior distribution cannot be computed exactly (as it requires summing over 2^k terms in the denominator of the Bayes rule). Therefore, we cannot just introduce a proposal distribution that is made equal to the true posterior as in the Expectation step of EM for GMM. Instead, we introduce a proposal distribution that factorizes across latent variable, i.e. each latent variable is independent wrt to other latent variable. In other words, the proposal distribution is defined as follows:

$$q_\phi(\mathbf{z}|\mathbf{x}) = \prod_{i=1}^k (\phi_i^{(\mathbf{x})})^{z_i} (1 - \phi_i^{(\mathbf{x})})^{1-z_i}$$

where $\phi^{(\mathbf{x})} \in [0, 1]^k$ are the parameters of the proposal distribution associated with observation \mathbf{x} . We can observe that they are independent Bernoullis. This technique of simplifying the distribution of the proposal is called Mean Field Theory in statistical physics.

The bound we present here is a simplified version of the one proposed by [Saul et al., 1996]. We first derive the evidence lower bound via Jensen’s inequality:

$$\begin{aligned} \log p_\theta(\mathbf{x}) &= \log \sum_{\mathbf{y} \in \mathcal{Y}} p_\theta(\mathbf{y})p_\theta(\mathbf{x}|\mathbf{y}) \\ &\geq \underbrace{\mathbb{E}_{q_\phi(\mathcal{Y}|\mathbf{x})}[\log p_\theta(\mathcal{Y})]}_{\text{(a)}} + \underbrace{\mathbb{E}_{q_\phi(\mathcal{Y}|\mathbf{x})}[\log p_\theta(\mathbf{x}|\mathcal{Y})]}_{\text{(b)}} + \underbrace{H^S[q_\phi(\mathcal{Y}|\mathbf{x})]}_{\text{(c)}} \end{aligned}$$

where H^S is the Shannon entropy. Although this equation doesn’t seem simpler (we still need to sum over 2^k values in the tree expectations), we show in the following that the mean field assumption on q allows to compute these expectations efficiently (except (b) which requires to derive another bound on this term).

Term **(a)** can be rewritten as:

$$\begin{aligned}\mathbb{E}_{q_\phi(\mathcal{Y}|\mathbf{x})}[\log p_\theta(\mathcal{Y})] &= \mathbb{E}_{q_\phi(\mathcal{Y}|\mathbf{x})} \left[\log \prod_{i=1}^k \frac{\exp(\mathcal{Y}_i a_i)}{1 + \exp(a_i)} \right] \\ &= \mathbb{E}_{q_\phi(\mathcal{Y}|\mathbf{x})} \left[\sum_{i=1}^k (\mathcal{Y}_i a_i - \log(1 + \exp(a_i))) \right]\end{aligned}$$

Note that by linearity of the expectation, we can separate the two part of the subtraction. Moreover, the term $\log(1 + \exp(a_i))$ does not depends on Y , therefore:

$$\begin{aligned}&= \sum_{i=1}^k \mathbb{E}_{q_\phi(\mathcal{Y}|\mathbf{x})} [\mathcal{Y}_i a_i] - \sum_{i=1}^k \log(1 + \exp(a_i)) \\ &= \sum_{i=1}^k a_i \phi_i^{(\mathbf{x})} - \sum_{i=1}^k \log(1 + \exp(a_i)) \\ &= \langle \mathbf{a}, \boldsymbol{\phi}^{(\mathbf{x})} \rangle - \sum_{i=1}^k \log(1 + \exp(a_i))\end{aligned}$$

which is easy to compute, we don't need to sum over 2^k terms in this form.

Term **(b)** can be rewritten as:

$$\begin{aligned}\mathbb{E}_{q_\phi(\mathcal{Y}|\mathbf{x})}[\log p_\theta(\mathbf{x}|\mathcal{Y})] &= \mathbb{E}_{q_\phi(\mathcal{Y}|\mathbf{x})} \left[\log \prod_{i=1}^d \frac{\exp(x_i(\mathbf{B}_i \mathcal{Y} + c_i))}{1 + \exp(\mathbf{B}_i \mathcal{Y} + c_i)} \right] \\ &= \mathbb{E}_{q_\phi(\mathcal{Y}|\mathbf{x})} \left[\sum_{i=1}^d (x_i(\mathbf{B}_i \mathcal{Y} + c_i) - \log(1 + \exp(\mathbf{B}_i \mathcal{Y} + c_i))) \right] \\ &= \sum_{i=1}^d \mathbb{E}_{q_\phi(\mathcal{Y}|\mathbf{x})} [(x_i(\mathbf{B}_i \mathcal{Y} + c_i))] - \sum_{i=1}^d \mathbb{E}_{q_\phi(\mathcal{Y}|\mathbf{x})} [\log(1 + \exp(\mathbf{B}_i \mathcal{Y} + c_i))]\end{aligned}$$

As for term **(a)**, the first argument of the subtraction can be simplified. However, the second argument cannot be simplified directly and requires to sum over 2^k values. We derive a lower bound on this term using Jensen's inequality again (therefore, we actually compute a lower bound to the standard ELBO):

$$\begin{aligned}&\geq \sum_{i=1}^d x_i \mathbb{E}_{q_\phi(\mathcal{Y}|\mathbf{x})} \left[\sum_{j=1}^k \mathbf{B}_{i,k} \mathcal{Y}_k \right] + \sum_{i=1}^d x_i c_i \\ &\quad - \sum_{i=1}^d \log \left(\mathbb{E}_{q_\phi(\mathcal{Y}|\mathbf{x})} \left[1 + \exp(c_i) \prod_{j=1}^k \exp(B_{i,j} \mathcal{Y}_j) \right] \right) \\ &= \langle \mathbf{x}, \mathbf{B} \boldsymbol{\phi}^{(\mathbf{x})} \rangle + \langle \mathbf{x}, \mathbf{c} \rangle - \sum_{i=1}^d \log \left(1 + \exp(c_i) \prod_{j=1}^k (1 - \phi_j^{(\mathbf{x})} + \phi_j^{(\mathbf{x})} \exp(B_{i,j})) \right),\end{aligned}$$

which is easy to compute.

Lastly, the entropy term **(c)** can be rewritten as:

$$\begin{aligned}
H^S[q_\phi(\mathcal{Y}|\mathbf{x})] &= \mathbb{E}_{q_\phi(\mathcal{Y}|\mathbf{x})}[\log q_\phi(\mathcal{Y}|\mathbf{x})] \\
&= \mathbb{E}_{q_\phi(\mathcal{Y}|\mathbf{x})} \left[\log \left(\prod_{i=1}^k q_\phi(\mathcal{Y}_i|\mathbf{x}) \right) \right] \\
&= \sum_{i=1}^k \mathbb{E}_{q_\phi(\mathcal{Y}|\mathbf{x})} [\log q_\phi(\mathcal{Y}_i|\mathbf{x})]
\end{aligned}$$

Note that the value inside the expectation depends only on a single random variable element \mathcal{Y}_i , therefore it simplifies to:

$$\begin{aligned}
&= \sum_{i=1}^k \mathbb{E}_{q_\phi(\mathcal{Y}_i|\mathbf{x})} [\log q_\phi(\mathcal{Y}_i|\mathbf{x})] \\
&= \sum_{i=1}^k \left(\phi_i^{(\mathbf{x})} \log \phi_i^{(\mathbf{x})} + (1 - \phi_i^{(\mathbf{x})}) \log(1 - \phi_i^{(\mathbf{x})}) \right)
\end{aligned}$$

which is tractable.

3 Upper bound on the evidence

We now show how to derive an upper bound on the evidence. Although it is not a good idea to rely on this upper bound for training a SBN (i.e. at training time we want to maximize a lower bound, hopping that maximizing a lower bound will also push up the actual evidence of the training data), it can be useful for example to estimate the quality of the lower bound derived in the previous section by computing the gap between the two bounds.

The method described in this section is a simplification of the approach proposed in [Jaakkola and Jordan, 1996]. The probability of an observation is defined as:

$$p_\theta(\mathbf{x}) = \sum_{\mathbf{y} \in Y} p_\theta(\mathbf{y}) p_\theta(\mathbf{x}|\mathbf{y})$$

As $1 - \sigma(u) = \sigma(-u)$, we can write:

$$= \sum_{\mathbf{y} \in Y} p_\theta(\mathbf{y}) \prod_{i=1}^d \sigma((2x_i - 1)(\mathbf{B}_i \mathbf{y} + c_i))$$

We can replace the sigmoid by its variational formulation, i.e. $\sigma(u) = \inf_{\epsilon \in [0,1]} \exp(\epsilon u - H^{\text{FD}}[\epsilon])$ where $H^{\text{FD}}[\epsilon] = -\epsilon \log \epsilon - (1 - \epsilon) \log(1 - \epsilon)$ if the Fermi-Dirac entropy. Note that we need one variational parameter ϵ_i per observed variable dimension. We obtain:

$$= \sum_{\mathbf{y} \in Y} p_\theta(\mathbf{y}) \prod_{i=1}^d \inf_{\epsilon_i \in [0,1]} \exp(\epsilon_i (2x_i - 1)(\mathbf{B}_i \mathbf{y} + c_i) - H^{\text{FD}}[\epsilon_i])$$

The objectives in infinimizations are strictly positive, therefore we can take out the infinimizations from the product:

$$\begin{aligned}
&= \sum_{\mathbf{y} \in Y} p_\theta(\mathbf{y}) \inf_{\substack{\epsilon_i \in [0,1], \\ \forall i \in \{1 \dots d\}}} \prod_{i=1}^d \exp(\epsilon_i(2x_i - 1)(\mathbf{B}_i \mathbf{y} + c_i) - H^{\text{FD}}[\epsilon_1]) \\
&= \sum_{\mathbf{y} \in Y} p_\theta(\mathbf{y}) \inf_{\substack{\epsilon_i \in [0,1], \\ \forall i \in \{1 \dots d\}}} \exp(-\sum_{i=1}^d H^{\text{FD}}[\epsilon_1]) \exp(\sum_{i=1}^d \epsilon_i(2x_i - 1)(\mathbf{B}_i \mathbf{y} + c_i)) \\
&= \sum_{\mathbf{y} \in Y} p_\theta(\mathbf{y}) \inf_{\substack{\epsilon_i \in [0,1], \\ \forall i \in \{1 \dots d\}}} \exp(-\sum_{i=1}^d H^{\text{FD}}[\epsilon_1]) \exp(\sum_{i=1}^d \epsilon_i c_i(2x_i - 1)) \prod_{j=1}^k \exp(\sum_{i=1}^d \epsilon_i \mathbf{B}_{i,j}(2x_i - 1))^{y_j}
\end{aligned}$$

Note that the outside sum is an expectation:

$$= \mathbb{E}_{p_\theta(\mathcal{Y})} \left[\inf_{\substack{\epsilon_i \in [0,1], \\ \forall i \in \{1 \dots d\}}} \exp(-\sum_{i=1}^d H^{\text{FD}}[\epsilon_1]) \exp(\sum_{i=1}^d \epsilon_i c_i(2x_i - 1)) \prod_{j=1}^k \exp(\sum_{i=1}^d \epsilon_i \mathbf{B}_{i,j}(2x_i - 1))^{y_j} \right]$$

We obtain an upper bound by moving the expectation inside the infinimizations. To understand why this gives us an upper bound, instead of finding the optimal variational parameters for each latent variable assignments \mathbf{y} , now the assignment is shared across all possible $\mathbf{y} \in Y$. As it is an infinimization, which is then “less expressive”, we obtain an upperbound.

$$\begin{aligned}
&\leq \inf_{\substack{\epsilon_i \in [0,1], \\ \forall i \in \{1 \dots d\}}} \exp(-\sum_{i=1}^d H^{\text{FD}}[\epsilon_1]) \exp(\sum_{i=1}^d \epsilon_i c_i(2x_i - 1)) \prod_{j=1}^k \mathbb{E}_{p_\theta(\mathcal{Y}_i | \mathbf{x})} \left[\exp(\sum_{i=1}^d \epsilon_i \mathbf{B}_{i,j}(2x_i - 1))^{y_j} \right] \\
&= \inf_{\substack{\epsilon_i \in [0,1], \\ \forall i \in \{1 \dots d\}}} \exp(-\sum_{i=1}^d H^{\text{FD}}[\epsilon_1]) \exp(\sum_{i=1}^d \epsilon_i c_i(2x_i - 1)) \prod_{j=1}^k \left(1 - \sigma(a_i) + \sigma(a_i) \exp(\sum_{i=1}^d \epsilon_i \mathbf{B}_{i,j}(2x_i - 1)) \right)
\end{aligned}$$

Computing the objective in the infinimizations in the last expression are all tractable. Therefore, we have the following upper bound on the evidence:

$$\log p_\theta(\mathbf{x}) \leq \left(-\sum_{i=1}^d H^{\text{FD}}[\epsilon_1] + \sum_{i=1}^d \epsilon_i c_i(2x_i - 1) + \sum_{j=1}^k \left(1 - \sigma(a_i) + \sigma(a_i) \exp(\sum_{i=1}^d \epsilon_i \mathbf{B}_{i,j}(2x_i - 1)) \right) \right)$$

where $\epsilon \in [0, 1]^d$ are the variational parameters of the bound. Although this bound is easy to compute, this does not mean that computing the parameters ϵ that leads to the tightness bound possible is easy. Methods to compute the variational parameters are described in [Jaakkola and Jordan, 1996].

4 Expectation-Maximization algorithm

In this section, we show how the variational lower bound can be used to learn the parameters θ of a sigmoid belief network. The Expectation-Maximization algorithm is a learning algorithm the relies on block-coordinate ascent to optimize the evidence lower bound (ELBO). Given a dataset $\mathcal{D} = \{\mathbf{x}^{(i)}\}_{i=1}^n$, the log-likelihood of the dataset and its ELBO are defined as follows:

$$\max_{\theta \in \Theta} \sum_{\mathbf{x} \in \mathcal{D}} \log p_\theta(\mathbf{x}) \geq \max_{\theta \in \Theta, \phi \in \Phi} \sum_{\mathbf{x} \in \mathcal{D}} \text{ELBO}(\mathbf{x}, \theta, \phi)$$

The two step of the EM algorithms are:

1. Expectation step: maximize over variational parameters ϕ ;
2. Maximization step: maximize over model parameters θ .

The iterative optimization process simply interleaves between E and M steps. In the case of GMMs, both steps enjoy simple closed-form expression. This is not the case of SBNs. We show below how the E step for SBN can be approximated efficiently. For the M step, one can simply rely on a single step of gradient ascent.

To simplify notation, for each training point $\mathbf{x} \in \mathcal{D}$ we introduce vector $\mathbf{w}^{(\mathbf{x})} \in \mathbb{R}^d$ and matrix $\mathbf{U}^{(\mathbf{x})} \in \mathbb{R}^{d,k}$ defined as follows:

$$w_j^{(\mathbf{x})} = \exp(c_j) \prod_{l=1}^k (1 - \phi_l^{(\mathbf{x})} + \phi_l^{(\mathbf{x})} \exp(B_{j,l}))$$

$$U_{j,l}^{(\mathbf{x})} = \frac{\partial}{\partial \phi_l^{(\mathbf{x})}} \log(1 + w_j^{(\mathbf{x})}) = \frac{w_j^{(\mathbf{x})}}{1 + w_j^{(\mathbf{x})}} \times \frac{\exp(B_{j,l}) - 1}{1 - \phi_l^{(\mathbf{x})} + \phi_l^{(\mathbf{x})} \exp(B_{j,l})}$$

The objective of the E step is defined as follows:

$$\max_{\phi \in \Phi} \sum_{\mathbf{x} \in \mathcal{D}} \text{ELBO}(\mathbf{x}, \theta, \phi) = \max_{\substack{\phi^{(\mathbf{x})} \in [0,1]^k \\ i \in \{1 \dots n\}}} \sum_{\mathbf{x} \in \mathcal{D}} \left(\begin{array}{l} \langle \mathbf{a}, \phi^{(\mathbf{x})} - \sum_{j=1}^k \log(1 + \exp(a_j)) \\ \langle \mathbf{x}, \mathbf{B} \phi^{(\mathbf{x})} \rangle + \langle \mathbf{x}, \mathbf{c} \rangle - \sum_{j=1}^d \log(1 + w_j^{(\mathbf{x})}) \\ + \sum_{j=1}^k (\phi_j^{(\mathbf{x})} \log \phi_j^{(\mathbf{x})} + (1 - \phi_j^{(\mathbf{x})}) \log(1 - \phi_j^{(\mathbf{x})})) \end{array} \right)$$

Ignoring the constraints on ϕ , by first order optimality conditions we have:

$$\log \frac{\phi_m^{(\mathbf{x})}}{1 - \phi_m^{(\mathbf{x})}} = a_m + \langle \mathbf{x}^{(i)}, \mathbf{B}_{-,m} \rangle - \sum_{j=1}^d U_{j,m}^{(\mathbf{x})}$$

$$\phi_m^{(\mathbf{x})} = \sigma \left(a_m + \langle \mathbf{x}^{(i)}, \mathbf{B}_{-,m} \rangle - \sum_{j=1}^d U_{j,m}^{(\mathbf{x})} \right)$$

where $\mathbf{B}_{-,m}$ denotes column m of matrix \mathbf{B} . Note that in the last equation, each $\phi_m^{(\mathbf{x})}$ depends on the values in $\phi^{(\mathbf{x})}$ via the matrix $\mathbf{U}^{(\mathbf{x})}$. However, each vector $\phi^{(\mathbf{x})}$ is independent of the proposal parameters of other datapoints. We can solve the equations defined by first order optimality conditions via standard iterative methods, i.e. we initialize $\phi^{(\mathbf{x})}$ randomly and update its value with the right-hand side of the last equation for a prefixed number of iterations. The resulting solution will always satisfies constraints on $\phi^{(\mathbf{x})}$.

References

- [Jaakkola and Jordan, 1996] Jaakkola, T. S. and Jordan, M. I. (1996). Computing upper and lower bounds on likelihoods in intractable networks. In *Proceedings of the Twelfth Conference on Uncertainty in Artificial Intelligence (UAI1996)*.
- [Saul et al., 1996] Saul, L. K., Jaakkola, T., and Jordan, M. I. (1996). Mean field theory for sigmoid belief networks. *Journal of artificial intelligence research*, 4:61–76.