

# Notes on Mean Field Theory for Sigmoid Belief Networks

Caio Corro

September 18, 2021

This are my personal notes on the paper "Mean Field Theory for Sigmoid Belief Networks" (Saul, Jaakkola and Jordan).

## 1 Definition

Let  $Z \in \{0, 1\}^m$  and  $X \in \{0, 1\}^n$  be the sets of latent and observed random variables, respectively. We assume that the data is generated as follows:

$$\begin{aligned} (1) \quad \mathbf{z} &\sim p(Z; \boldsymbol{\mu}) && \text{where } \boldsymbol{\mu} = \sigma(\mathbf{a}) \\ (2) \quad \mathbf{x} &\sim p(X|Z = \mathbf{z}; \boldsymbol{\phi}) && \text{where } \boldsymbol{\phi} = \sigma(\mathbf{B}\mathbf{z} + \mathbf{c}) \end{aligned}$$

where  $\sigma$  is the element-wise sigmoid function and  $\mathbf{a} \in \mathbb{R}^m$ ,  $\mathbf{B} \in \mathbb{R}^{n \times m}$  and  $\mathbf{c} \in \mathbb{R}^n$  are the parameters of the model. The distributions  $p(Z)$  and  $p(X|Z)$  are distributions over vector of binary values where we assume independence between coordinates, i.e.  $p(Z) = \prod_i p(Z_i)$  and  $p(X|Z) = \prod_i p(X_i|Z)$ . In order to highlight the dependency between parameters and distribution, we will write  $p(Z; \mathbf{a})$  and  $p(X|Z; \mathbf{B}, \mathbf{c})$ . We will omit the parameters when clear from context.

During training, the goal is to maximize the log-likelihood of a dataset of observed values  $\mathbf{x}$  by tuning the parameters of the model. In theory, this can be done via gradient descent by marginalizing over latent variables, i.e. computing  $\log p(\mathbf{x}) = \log \sum_{\mathbf{z}} p(\mathbf{x})p(\mathbf{x}|\mathbf{z})$ . However, as we use a neural parameterization of the conditional distribution, it cannot be computed efficiently with a closed form solution. It is therefore required to sum over the  $2^m$  possible latent variable assignments, which is prohibitive in practice. A second solution would be to rely on the Expectation-Maximization algorithm, but similar issues appear when computing the posterior distribution in the expectation step.

Using mean field theory, we show how we can efficiently derive a lower bound on the true likelihood which can be optimized via gradient ascent.

## 2 Mean Field Theory

We assume a parameterized proposal distribution where each coordinate is independent:

$$q(\mathbf{z}|\mathbf{x}; \bar{\boldsymbol{\mu}}) = \prod_i q(z_i|\mathbf{x}; \bar{\mu}_i)$$

where  $\bar{\boldsymbol{\mu}} \in [0, 1]^m$  are the parameters of the proposal distribution, i.e. each  $q(z_i|\mathbf{x}; \bar{\mu}_i)$  is a Bernoulli parameterized by  $\bar{\mu}_i$ . We omit to explicitly note the dependency between  $\bar{\boldsymbol{\mu}}$  and  $\mathbf{x}$  for simplicity, but one proposal distribution per datapoint is computed in practice (i.e. the proposal is not amortized on the dataset). In this section, we highlight properties of statistical independency that will be useful in the following.

First, we show that the entropy of the joint distribution can be computed efficiently. We have:

$$H[q(Z_i|\mathbf{x})] = - \sum_{z_i} p(z_i|\mathbf{x}) \log p(z_i|\mathbf{x}) = -p(0|\mathbf{x}) \log p(0|\mathbf{x}) - p(1|\mathbf{x}) \log p(1|\mathbf{x}) = -\bar{\mu}_i \log \bar{\mu}_i - (1 - \bar{\mu}_i) \log(1 - \bar{\mu}_i)$$

As each coordinate is independent, we can show that the entropy of the joint distribution decomposes as the sum of entropies:

$$H[q(Z|\mathbf{x})] = \sum_i H[q(Z_i|\mathbf{x})] = \sum_i -\bar{\mu}_i \log \bar{\mu}_i - (1 - \bar{\mu}_i) \log(1 - \bar{\mu}_i)$$

Note that the partial derivative wrt any given  $\bar{\mu}_j$  as the form:

$$\begin{aligned}
\frac{\partial}{\partial \bar{\mu}_j} H[q(Z|\mathbf{x})] &= \frac{\partial}{\partial \bar{\mu}_j} \sum_i -\bar{\mu}_i \log \bar{\mu}_i - (1 - \bar{\mu}_i) \log(1 - \bar{\mu}_i) \\
&= \frac{\partial}{\partial \bar{\mu}_j} -\bar{\mu}_j \log \bar{\mu}_j - (1 - \bar{\mu}_j) \log(1 - \bar{\mu}_j) \\
&= -\log \bar{\mu}_j - 1 + \log(1 - \bar{\mu}_j) + 1 \\
&= -\log \bar{\mu}_j + \log(1 - \bar{\mu}_j) \\
&= -\log \frac{\bar{\mu}_j}{1 - \bar{\mu}_j}
\end{aligned}$$

### 3 Efficient approximation of the log-likelihood

In this section we show that under the mean field assumption we can derive tractable lower bound of the true likelihood that can be optimized via gradient ascent. Let  $q(\mathbf{z}|\mathbf{x}; \bar{\boldsymbol{\mu}}) = \prod_i q(z_i|\mathbf{x}; \bar{\mu}_i)$  be the proposal mean field distribution. The usual evidence lower bound is derived as follows:

$$\begin{aligned}
\log p(\mathbf{x}) &= \log \sum_{\mathbf{z}} p(\mathbf{z})p(\mathbf{x}|\mathbf{z}) \\
&= \log \sum_{\mathbf{z}} \frac{q(\mathbf{z}|\mathbf{x})}{q(\mathbf{z}|\mathbf{x})} p(\mathbf{z})p(\mathbf{x}|\mathbf{z})
\end{aligned}$$

By Jensen's inequality:

$$\begin{aligned}
&\geq \sum_{\mathbf{z}} q(\mathbf{z}|\mathbf{x}) \log \frac{p(\mathbf{z})p(\mathbf{x}|\mathbf{z})}{q(\mathbf{z}|\mathbf{x})} \\
&= \mathbb{E}_{q(\mathbf{z}|\mathbf{x})} [\log p(\mathbf{z})] + \mathbb{E}_{q(\mathbf{z}|\mathbf{x})} [p(\mathbf{x}|\mathbf{z})] + H[q]
\end{aligned}$$

Note that the last term is constant wrt to the network parameters, so its gradient is not required. We are left with the two expectations. We show how to compute them efficiently, meaning that gradient wrt the parameters can be deduced via the backpropagation algorithm.

First, note that we have:

$$\log p(\mathbf{z}; \mathbf{a}) = \log \prod_i p(z_i; \mathbf{a}) = \sum_i \log p(z_i; a_i) \quad \text{where} \quad \log p(z_i; a_i) = \begin{cases} \sigma(a_i) & \text{if } z_i = 1 \\ 1 - \sigma(a_i) & \text{otherwise} \end{cases}$$

We can rewrite the probabilities as follow:

$$\begin{aligned}
\sigma(a_i) &= \frac{\exp(a_i)}{1 + \exp(a_i)} = \frac{\exp(a_i z_i)}{1 + \exp(a_i)} \quad \text{when } z_i = 1 \\
1 - \sigma(a_i) &= \frac{1 + \exp(a_i)}{1 + \exp(a_i)} - \frac{\exp(a_i)}{1 + \exp(a_i)} = \frac{1}{1 + \exp(a_i)} = \frac{\exp(a_i z_i)}{1 + \exp(a_i)} \quad \text{when } z_i = 0
\end{aligned}$$

We can therefore write:

$$\log p(\mathbf{z}; \mathbf{a}) = \sum_i \log \frac{\exp(a_i z_i)}{1 + \exp(a_i)} = \sum_i (a_i z_i - \log(1 + \exp(a_i)))$$

Taking the expectation under the proposal distribution, we have:

$$\begin{aligned}
\mathbb{E}_{q(\mathbf{z}|\mathbf{x})} [\log p(\mathbf{z})] &= \mathbb{E}_{q(\mathbf{z}|\mathbf{x})} \left[ \sum_i (a_i z_i - \log(1 + \exp(a_i))) \right] \\
&= \sum_i (\mathbb{E}_{q(\mathbf{z}|\mathbf{x})} [a_i z_i] - \log(1 + \exp(a_i))) \\
&= \sum_i (\mathbb{E}_{q(z_i|\mathbf{x})} [(a_i z_i)] - \log(1 + \exp(a_i))) \\
&= \sum_i (a_i \bar{\mu}_i - \log(1 + \exp(a_i)))
\end{aligned}$$

Note that this term is easy to compute.

For the second term in the ELBO, note that we have:

$$\begin{aligned}\log p(\mathbf{x}|\mathbf{z}; \mathbf{B}, \mathbf{c}) &= \sum_i \log p(x_i|\mathbf{z}; \mathbf{B}, \mathbf{c}) \\ &= \sum_i \log \frac{\exp(x_i(c_i + \sum_j B_{i,j}z_j))}{1 + \exp(c_i + \sum_j B_{i,j}z_j)} \\ &= \sum_i \left( c_i x_i + \sum_j B_{i,j} z_j x_i - \log(1 + \exp(c_i + \sum_j B_{i,j}z_j)) \right)\end{aligned}$$

Taking the expectation wrt the proposal distribution, we have:

$$\begin{aligned}\mathbb{E}_{q(\mathbf{z}|\mathbf{x})} [\log p(\mathbf{x}|\mathbf{z})] &= \mathbb{E}_{q(\mathbf{z}|\mathbf{x})} \left[ \sum_i \left( c_i x_i + \sum_j B_{i,j} z_j x_i - \log(1 + \exp(c_i + \sum_j B_{i,j}z_j)) \right) \right] \\ &= \sum_i \mathbb{E}_{q(\mathbf{z}|\mathbf{x})} \left[ c_i x_i + \sum_j B_{i,j} z_j x_i - \log(1 + \exp(c_i + \sum_j B_{i,j}z_j)) \right] \\ &= \sum_i c_i x_i + \sum_j B_{i,j} \bar{\mu}_j x_i - \mathbb{E}_{q(\mathbf{z}|\mathbf{x})} \left[ \log(1 + \exp(c_i + \sum_j B_{i,j}z_j)) \right]\end{aligned}$$

Unfortunately the last term is intractable as we cannot break down the expectation over the joint distribution in simpler computation. However, we can derive a lower bound using Jensen's inequality:

$$\begin{aligned}&\geq \sum_i c_i x_i + \sum_j B_{i,j} \bar{\mu}_j x_i - \log \left( \mathbb{E}_{q(\mathbf{z}|\mathbf{x})} \left[ 1 + \exp(c_i) \prod_j \exp(B_{i,j}z_j) \right] \right) \\ &= \sum_i c_i x_i + \sum_j B_{i,j} \bar{\mu}_j x_i - \log \left( 1 + \exp(c_i) \prod_j (\bar{\mu}_j \exp(B_{i,j}) + 1 - \bar{\mu}_j) \right)\end{aligned}$$

which can be efficiently computed.

As such, we showed that if we are given the proposal distribution parameters  $\bar{\boldsymbol{\mu}}$  we can efficiently learn the parameters  $\mathbf{a}, \mathbf{B}$  and  $\mathbf{c}$  of the SBN via gradient ascent. In the next section, we show how we can compute the optimal proposal mean field distribution.

## 4 Proposal distribution estimation

To compute the parameters  $\bar{\boldsymbol{\mu}}$ , i.e. the expectation step, we can maximize the ELBO wrt to  $\bar{\boldsymbol{\mu}}$ :

$$\begin{aligned}\arg \max_{\bar{\boldsymbol{\mu}}} \quad & \mathbb{E}_{q(\mathbf{z}|\mathbf{x})} [\log p(\mathbf{z})p(\mathbf{x}|\mathbf{z})] + H[q] \\ \text{s.t.} \quad & \bar{\mu}_i \geq 0 \quad \forall i \\ & \bar{\mu}_i \leq 1 \quad \forall i\end{aligned}$$

We will ignore the inequality constraints for simplicity. A more formal derivation should include them, but their contribution to the stationarity conditions will vanish as the associated dual variables are fixed to zero by complementary slackness and the fact the each  $\bar{\mu}_i$  will have a sigmoidal shape.

Hence, if we focus on the unconstrained variant of the problem that is concave and differentiable, we can derive optimal parameters  $\bar{\boldsymbol{\mu}}$  it by solving the following system of equations:

$$\nabla_{\bar{\boldsymbol{\mu}}} (\mathbb{E}_{q(\mathbf{z}|\mathbf{x})} [\log p(\mathbf{z})p(\mathbf{x}|\mathbf{z})] + H[q]) = \mathbf{0}$$

However, as we saw in the previous section, the expectation step is intractable. We replace this objective by

the lower bound from previous section and obtain the following partial derivative constraint:

$$\begin{aligned}
\log \frac{\bar{\mu}_k}{1 - \bar{\mu}_k} &= \frac{\partial}{\partial \bar{\mu}_k} \left( \sum_i c_i x_i + \sum_j B_{i,j} \bar{\mu}_j x_i - \log \left( 1 + \exp(c_i) \prod_j (\bar{\mu}_j \exp(B_{i,j}) + 1 - \bar{\mu}_j) \right) \right) \\
&= B_{i,k} x_i - \frac{1}{1 + \exp(c_i) \prod_j (\bar{\mu}_j \exp(B_{i,j}) + 1 - \bar{\mu}_j)} \frac{\partial}{\partial \bar{\mu}_k} \left( 1 + \exp(c_i) \prod_j (\bar{\mu}_j \exp(B_{i,j}) + 1 - \bar{\mu}_j) \right) \\
&= B_{i,k} x_i - \frac{(\exp(B_{i,k}) - 1) \exp(c_i) \prod_{j \neq k} (\bar{\mu}_j \exp(B_{i,j}) + 1 - \bar{\mu}_j)}{1 + \exp(c_i) \prod_j (\bar{\mu}_j \exp(B_{i,j}) + 1 - \bar{\mu}_j)} \\
\bar{\mu}_k &= \sigma \left( B_{i,k} x_i - \frac{(\exp(B_{i,k}) - 1) \exp(c_i) \prod_{j \neq k} (\bar{\mu}_j \exp(B_{i,j}) + 1 - \bar{\mu}_j)}{1 + \exp(c_i) \prod_j (\bar{\mu}_j \exp(B_{i,j}) + 1 - \bar{\mu}_j)} \right)
\end{aligned}$$

Note that this does not result in a closed form expression for  $\bar{\mu}$ . However, we can compute iteratively a solution to this set of equations as follows:

$$\bar{\mu}_k^{(t+1)} = \sigma \left( B_{i,k} x_i - \frac{(\exp(B_{i,k}) - 1) \exp(c_i) \prod_{j \neq k} (\bar{\mu}_j^{(t)} \exp(B_{i,j}) + 1 - \bar{\mu}_j^{(t)})}{1 + \exp(c_i) \prod_j (\bar{\mu}_j^{(t)} \exp(B_{i,j}) + 1 - \bar{\mu}_j^{(t)})} \right)$$