

RBM

Oliver K. Ernst

September 30, 2020

1 Theory

The objective function is the KL divergence:

$$\mathcal{D}_{\mathcal{KL}}(p||\tilde{p}) = \int dx p(x) \ln \frac{p(x)}{\tilde{p}(x, \nu)} \quad (1)$$

where $p(x)$ is the true data distribution and $\tilde{p}(x, \nu)$ is the model distribution:

$$\begin{aligned} \tilde{p}(x, \nu) &= \frac{1}{Z(\nu)} \exp[-E(x, \nu)] \\ Z(\nu) &= \int dy \exp[-E(y, \nu)] \end{aligned} \quad (2)$$

for some energy function $E(x, \nu)$ with interactions ν . The gradients are

$$\begin{aligned} \frac{\partial \mathcal{D}_{\mathcal{KL}}}{\partial \nu} &= - \int dx p(x) \left(\frac{\tilde{p}(x, \nu)}{p(x)} \right) p(x) (\tilde{p}(x, \nu))^{-2} \frac{\partial \tilde{p}(x, \nu)}{\partial \nu} = - \int dx \frac{p(x)}{\tilde{p}(x, \nu)} \frac{\partial \tilde{p}(x, \nu)}{\partial \nu} \\ \frac{\partial \tilde{p}(x, \nu)}{\partial \nu} &= - \frac{\partial E(x, \nu)}{\partial \nu} \tilde{p}(x, \nu) - \frac{1}{Z(\nu)^2} \exp[-E(x, \nu)] \frac{\partial Z(\nu)}{\partial \nu} \\ &= - \frac{\partial E(x, \nu)}{\partial \nu} \tilde{p}(x, \nu) + \frac{1}{Z(\nu)} \tilde{p}(x, \nu) \int dy \frac{\partial E(y, \nu)}{\partial \nu} \exp[-E(y, \nu)] \\ &= - \frac{\partial E(x, \nu)}{\partial \nu} \tilde{p}(x, \nu) + \tilde{p}(x, \nu) \int dy \frac{\partial E(y, \nu)}{\partial \nu} \tilde{p}(y, \nu) \\ &= - \frac{\partial E(x, \nu)}{\partial \nu} \tilde{p}(x, \nu) + \tilde{p}(x, \nu) \left\langle \frac{\partial E}{\partial \nu} \right\rangle_{\tilde{p}} \\ \frac{\partial \mathcal{D}_{\mathcal{KL}}}{\partial \nu} &= \int dx \frac{p(x)}{\tilde{p}(x, \nu)} \frac{\partial E(x, \nu)}{\partial \nu} \tilde{p}(x, \nu) - \int dx \frac{p(x)}{\tilde{p}(x, \nu)} \tilde{p}(x, \nu) \left\langle \frac{\partial E}{\partial \nu} \right\rangle_{\tilde{p}} \\ &= \int dx p(x) \frac{\partial E(x, \nu)}{\partial \nu} - \int dx p(x) \left\langle \frac{\partial E}{\partial \nu} \right\rangle_{\tilde{p}} \\ &= \left\langle \frac{\partial E}{\partial \nu} \right\rangle_p - \left\langle \frac{\partial E}{\partial \nu} \right\rangle_{\tilde{p}} \end{aligned} \quad (3)$$

where the second integral is unity because p is normalized by definition.

$$\boxed{\frac{\partial \mathcal{D}_{\mathcal{KL}}}{\partial \nu} = \left\langle \frac{\partial E}{\partial \nu} \right\rangle_p - \left\langle \frac{\partial E}{\partial \nu} \right\rangle_{\tilde{p}}} \quad (4)$$

The first term is often called the *awake phase* moment, or moment under the *data distribution* p ; the second term is often called the *asleep phase* moment, or moment under the *model distribution* \tilde{p} .

For the discrete case on a lattice with visible units \mathbf{v} and hidden units \mathbf{h} :

$$\mathcal{D}_{\mathcal{KL}}(p||\tilde{p}) = \sum_{\mathbf{v}} \sum_{\mathbf{h}} p(\mathbf{v}, \mathbf{h}) \ln \frac{p(\mathbf{v}, \mathbf{h})}{\tilde{p}(\mathbf{v}, \mathbf{h})} \quad (5)$$

and a common energy function is:

$$E(\mathbf{v}, \mathbf{h}) = -\mathbf{a}^\top \mathbf{v} - \mathbf{b}^\top \mathbf{h} - \mathbf{v}^\top W \mathbf{h} \quad (6)$$

for biases \mathbf{a}, \mathbf{b} and weight matrix W , which play the role of ν . The gradients are:

$$\begin{aligned} \frac{\partial E}{\partial \mathbf{a}} &= -\mathbf{v} \\ \frac{\partial E}{\partial \mathbf{b}} &= -\mathbf{h} \\ \frac{\partial E}{\partial W} &= -\mathbf{v} \otimes \mathbf{h} \end{aligned} \quad (7)$$

leading to the gradients:

$$\boxed{\begin{aligned} \frac{\partial \mathcal{D}_{\mathcal{KL}}}{\partial \mathbf{a}} &= \langle \mathbf{v} \rangle_{\tilde{p}} - \langle \mathbf{v} \rangle_p \\ \frac{\partial \mathcal{D}_{\mathcal{KL}}}{\partial \mathbf{b}} &= \langle \mathbf{h} \rangle_{\tilde{p}} - \langle \mathbf{h} \rangle_p \\ \frac{\partial \mathcal{D}_{\mathcal{KL}}}{\partial W} &= \langle \mathbf{v} \otimes \mathbf{h} \rangle_{\tilde{p}} - \langle \mathbf{v} \otimes \mathbf{h} \rangle_p \end{aligned}} \quad (8)$$

2 Implementation

2.1 Gradients

In practice, concerning the moments $\langle \dots \rangle$:

- In the continuous case, we usually cannot analytically perform the integral $\int dx$.
- In the discrete case, we usually cannot enumerate all possible states appearing in the sum \sum_x .

Therefore, these moments are estimated using batches. In the continuous case, let the batch be $X_{\tilde{p}}$ of size N , typically small ($N \sim 5 - 10$). For some observable $\chi(x)$:

$$\langle \chi(x) \rangle_{\tilde{p}} = \int dx \tilde{p}(x, \nu) \chi(x) \sim \frac{1}{N} \sum_{i=1}^N \chi(X_{\tilde{p},i}) \quad (9)$$

and similarly for moments with respect to p . In the discrete case, with the batch represented as $\mathbf{V}_{\tilde{p}}, \mathbf{H}_{\tilde{p}}$:

$$\langle \chi(\mathbf{v}, \mathbf{h}) \rangle_{\tilde{p}} = \sum_{\mathbf{v}} \sum_{\mathbf{h}} \tilde{p}(\mathbf{v}, \mathbf{h}) \chi(\mathbf{v}, \mathbf{h}) \sim \frac{1}{N} \sum_{i=1}^N \chi(\mathbf{V}_{\tilde{p},i}, \mathbf{H}_{\tilde{p},i}) \quad (10)$$

The gradients are therefore estimated as:

$$\boxed{\begin{aligned} \frac{\partial \mathcal{D}_{\mathcal{KL}}}{\partial \mathbf{a}} &= \frac{1}{N} \sum_{i=1}^N \mathbf{V}_{\tilde{p},i} - \frac{1}{N} \sum_{i=1}^N \mathbf{V}_{p,i} \\ \frac{\partial \mathcal{D}_{\mathcal{KL}}}{\partial \mathbf{b}} &= \frac{1}{N} \sum_{i=1}^N \mathbf{H}_{\tilde{p},i} - \frac{1}{N} \sum_{i=1}^N \mathbf{H}_{p,i} \\ \frac{\partial \mathcal{D}_{\mathcal{KL}}}{\partial W} &= \frac{1}{N} \sum_{i=1}^N \mathbf{V}_{\tilde{p},i} \otimes \mathbf{H}_{\tilde{p},i} - \frac{1}{N} \sum_{i=1}^N \mathbf{V}_{p,i} \otimes \mathbf{H}_{p,i} \end{aligned}} \quad (11)$$

2.2 Sampling

Where do such batches come from? They should be samples of the distributions p or \tilde{p} as appropriate.

- For sampling the model distribution \tilde{p} , we have several options:
 - In the continuous case, without knowing further about $\tilde{p}(x)$, we can always perform Markov Chain Monte Carlo (MCMC) to sample the distributions (other sampling methods *may* be possible).
 - In the discrete case, we can use *Gibbs sampling*, which works by iteratively sampling $\tilde{p}(v_i|\mathbf{h})$ and $\tilde{p}(h_i|\mathbf{v})$. These are derived as follows:

$$\begin{aligned}\tilde{p}(h_i|\mathbf{v}) &= \frac{\tilde{p}(h_i, \mathbf{v})}{\tilde{p}(\mathbf{v})} \\ \tilde{p}(\mathbf{v}) &= \sum_{\mathbf{h}} \tilde{p}(\mathbf{v}, \mathbf{h}) \propto \exp[\mathbf{a}^\top \mathbf{v}] \\ \tilde{p}(\mathbf{h}, \mathbf{v}) &= \prod_i \tilde{p}(h_i, \mathbf{v})\end{aligned}\tag{12}$$

where the last line follows because in an RBM, the hidden variables are conditionally independent of all other hidden variables, and similarly for visible-variables, then:

$$\tilde{p}(h_i, \mathbf{v}) \propto \exp[\mathbf{a}^\top \mathbf{v} + b_i h_i + \mathbf{v}^\top \text{col}_i(W) h_i]\tag{13}$$

then it follows:

$$\boxed{\tilde{p}(h_i|\mathbf{v}) \propto \exp[b_i h_i + \mathbf{v}^\top \text{col}_i(W) h_i]}\tag{14}$$

and similarly

$$\boxed{\tilde{p}(v_i|\mathbf{h}) \propto \exp[a_i v_i + v_i \text{row}_i(W) \mathbf{h}]}\tag{15}$$

In *contrastive divergence*, this procedure of iteratively sampling $\tilde{p}(v_i|\mathbf{h})$ and $\tilde{p}(h_i|\mathbf{v})$ is performed only a few times (or even only once!), starting from an initial data vector \mathbf{v} . This greatly improves computational efficiency; alternatively, you can run this sampler for a long time to let the *chain converge*.

After sampling, we obtain the desired samples for the batch $\mathbf{V}_{\tilde{p},i}, \mathbf{H}_{\tilde{p},i}$. The sampling can be performed in *parallel* for higher efficiency, evaluating all N items in the batch $\mathbf{V}_{\tilde{p}}, \mathbf{H}_{\tilde{p}}$.

In *persistent contrastive divergence*, we do not throw out the hidden states $\mathbf{V}_{\tilde{p}}, \mathbf{H}_{\tilde{p}}$ after one gradient step and restart from new data vectors \mathbf{v} . Instead, we keep these states, and use them in the next gradient step again after sampling $\tilde{p}(v_i|\mathbf{h})$ and $\tilde{p}(h_i|\mathbf{v})$ for a few more steps.

- The samples \mathbf{V}_p from the data distribution p are obvious; they are provided as training data. The samples \mathbf{H}_p are obtained by *clamping the visible units to the data vectors* \mathbf{V}_p and sampling *only* $\tilde{p}(h_i|\mathbf{v})$. Note that this is similarly possible with MCMC for the general continuous case; assuming we partition x into observed and latent variables, we can sample only the latent variables, keeping the observed variables clamped.

2.3 Objective function

It would be great if we could code up the objective function $\mathcal{D}_{\mathcal{KL}}$ onto a computer, but this is **not trivial**. Instead, we can make an important restriction:

Assume that our **optimizer only uses first-order gradients**.

If we make this restriction, we can consider the following objective function

$$\begin{aligned}
S = & \mathbf{a}^\top \left(\frac{1}{N} \sum_{i=1}^N \mathbf{v}_{\tilde{p},i} - \frac{1}{N} \sum_{i=1}^N \mathbf{v}_{p,i} \right) \\
& + \mathbf{b}^\top \left(\frac{1}{N} \sum_{i=1}^N \mathbf{H}_{\tilde{p},i} - \frac{1}{N} \sum_{i=1}^N \mathbf{H}_{p,i} \right) \\
& + \left(\frac{1}{N} \sum_{i=1}^N \mathbf{v}_{\tilde{p},i}^\top \mathbf{W} \mathbf{H}_{\tilde{p},i} - \frac{1}{N} \sum_{i=1}^N \mathbf{v}_{p,i}^\top \mathbf{W} \mathbf{H}_{p,i} \right)
\end{aligned} \tag{16}$$

which has the same **first-order** gradients as (11).

Note that the second-order gradients will obviously be incorrect!

This trick allows us to easily implement an RBM in TensorFlow.