

# Lab Exercise 2: Restricted Boltzmann Machine

Caio Corro

## 1 Introduction

In this lab exercise, you will build a Restricted Boltzmann Machine (RBM) with continuous observed variables and binary latent variables. Hence, we have the following sample spaces:

- $\mathbf{x} \in \mathbb{R}^2$  for the observed variables,
- $\mathbf{y} \in \{0, 1\}^n$  for the latent variables.

The joint probability distribution over all random variables is defined as follows:

$$p(\mathbf{x}, \mathbf{y}; \theta) = \frac{\exp(w(\mathbf{x}, \mathbf{y}; \theta))}{\int_{\mathbf{x}'} \sum_{\mathbf{y}'} \exp(w(\mathbf{x}', \mathbf{y}'; \theta)) d\mathbf{x}'} = \frac{\exp(w(\mathbf{x}, \mathbf{y}; \theta))}{Z(\theta)} = Z(\theta)^{-1} \exp(w(\mathbf{x}, \mathbf{y}; \theta)) = \exp(w(\mathbf{x}, \mathbf{y}; \theta) - c(\theta))$$

where  $w : \mathbb{R}^2 \times \{0, 1\}^n \rightarrow \mathbb{R}$  is the energy function parameterized by  $\theta$ ,  $Z(\theta)$  is the partition function and  $c(\theta) = \log Z(\theta)$  is the log-partition (or cumulant) function. Note that the integral and the sum are over all possible values for the observed and latent variables, respectively. As mentioned in the course, in a RBM we have the following conditional independence by construction (i.e. by how we define the function  $w$ ):

- $p(\mathbf{x}|\mathbf{y}; \theta) = p(x_1|\mathbf{y}; \theta)p(x_2|\mathbf{y}; \theta)$
- $p(\mathbf{y}|\mathbf{x}; \theta) = \prod_{i=1}^n p(y_i|\mathbf{x}; \theta)$

Moreover, both  $p(\mathbf{x}|\mathbf{y}; \theta)$  and  $p(\mathbf{y}|\mathbf{x}; \theta)$  are well known probability distributions. In the example of the course, they were both independent Bernoullis. In this lab exercise,  $p(\mathbf{x}|\mathbf{y}; \theta)$  are independent Gaussians and  $p(\mathbf{y}|\mathbf{x}; \theta)$  are independent Bernoullis (see below for a derivation). This fact allows us to rely on efficient Gibbs sampling in order to sample from the joint distribution  $p(\mathbf{x}, \mathbf{y}; \theta)$ : we start from a random assignation of either  $\mathbf{y}$  or  $\mathbf{x}$ , and then iteratively sample from  $p(X|\mathbf{y}; \theta)$  and  $p(Y|\mathbf{x}; \theta)$ . For example:

1.  $\forall i : \mathbf{y}_i^{(1)} \sim \mathbb{B}(0.5)$  where  $\mathbb{B}(0.5)$  is an unbiased Bernoulli distribution,
2.  $\mathbf{x}^{(2)} \sim p(X|\mathbf{y}^{(1)}; \theta)$ ,
3.  $\mathbf{y}^{(2)} \sim p(Y|\mathbf{x}^{(2)}; \theta)$ ,
4.  $\mathbf{x}^{(3)} \sim p(X|\mathbf{y}^{(2)}; \theta)$ ,
5.  $\mathbf{y}^{(3)} \sim p(Y|\mathbf{x}^{(3)}; \theta)$ ,
6. ...

And we keep one or several couple of values  $\langle \mathbf{x}^{(t)}, \mathbf{y}^{(t)} \rangle$  from this Markov chain as samples.

## 2 Weighting function

In this lab exercise, the RBM weighting function is defined as follows:

$$w(\mathbf{x}, \mathbf{y}; \theta) = -\sum_{i=1}^2 \frac{(x_i - b_i)^2}{2\sigma_i^2} + \sum_{i=1}^2 \sum_{j=1}^n \frac{x_i}{\sigma_i^2} W_{ij} y_j + \sum_{j=1}^n d_j y_j$$

where  $\theta = \{\mathbf{b}, \sigma, \mathbf{W}, \mathbf{d}\}$  is the set of parameters that must be learned from data:

- $\mathbf{b} \in \mathbb{R}^2$ ,
- $\sigma \in \mathbb{R}_{++}^2$  (that is  $\sigma_1 > 0$  and  $\sigma_2 > 0$ ),
- $\mathbf{W} \in \mathbb{R}^{2 \times n}$ ,
- $\mathbf{d} \in \mathbb{R}^n$ .

We can compute analytically the conditional probability of  $x_1$ :

$$\begin{aligned}
p(x_1|x_2, \mathbf{y}; \theta) &= \frac{p(x_1, x_2, \mathbf{y}; \theta)}{p(x_2, \mathbf{y}; \theta)} \\
&= \frac{p(x_1, x_2, \mathbf{y}; \theta)}{\int_{x'_1} p(x'_1, x_2, \mathbf{y}; \theta) dx'_1} \\
&= \frac{Z(\theta)^{-1} \exp\left(-\frac{(x_1-b_1)^2}{2\sigma_1^2} - \frac{(x_2-b_2)^2}{2\sigma_2^2} + \sum_{j=1}^n \frac{x_1}{\sigma_1^2} W_{1j} y_j + \sum_{j=1}^n \frac{x_2}{\sigma_2^2} W_{2j} y_j + \sum_{j=1}^n d_j y_j\right)}{\int_{x'_1} Z(\theta)^{-1} \exp\left(-\frac{(x'_1-b_1)^2}{2\sigma_1^2} - \frac{(x_2-b_2)^2}{2\sigma_2^2} + \sum_{j=1}^n \frac{x'_1}{\sigma_1^2} W_{1j} y_j + \sum_{j=1}^n \frac{x_2}{\sigma_2^2} W_{2j} y_j + \sum_{j=1}^n d_j y_j\right) dx'_1}
\end{aligned}$$

Remember the exponential rule  $\exp(a + b) = \exp(a)\exp(b)$ . In the denominator, we can get out of the integral the partition function term  $Z(\theta)^{-1}$  and all terms in the exponential that do not depend on  $x'_1$  thanks to the exponential rule. Then, these terms will cancel with their counterparts in the nominator, so we end up with:

$$= \frac{\exp\left(-\frac{(x_1-b_1)^2}{2\sigma_1^2} + \sum_{j=1}^n \frac{x_1}{\sigma_1^2} W_{1j} y_j\right)}{\int_{x'_1} \exp\left(-\frac{(x'_1-b_1)^2}{2\sigma_1^2} + \sum_{j=1}^n \frac{x'_1}{\sigma_1^2} W_{1j} y_j\right) dx'_1}$$

which shows the conditional independence of  $x_1$  given  $x_2$  as the latter value does not appear in the equation anymore. Let  $a = \sum_{j=1}^n W_{1j} y_j$ . Then we can rewrite the argument of the exponential in the numerator argument as:

$$\frac{ax_1}{\sigma_1^2} - \frac{(x_1 - b_1)^2}{2\sigma_1^2} = \frac{2ax_1 - x_1^2 - b_1^2 + 2b_1x_1}{2\sigma_1^2} = \frac{2ax_1 - x_1^2 - b_1^2 + 2b_1x_1 + a^2 - a^2 + 2ab_1 - 2ab_1}{2\sigma_1^2} = \frac{-(x_1 - a - b_1)^2 + a^2 + 2ab_1}{2\sigma_1^2}$$

By using this relation and similarly in the denominator, we can rewrite the conditional probability as:

$$= \frac{\exp\left(\frac{-(x_1 - \sum_{j=1}^n y_j W_{1j} - b_1)^2 + (\sum_{j=1}^n W_{1j} y_j)^2 + 2 \sum_{j=1}^n W_{1j} y_j b_1}{2\sigma_1^2}\right)}{\int_{x'_1} \exp\left(\frac{-(x'_1 - \sum_{j=1}^n W_{1j} y_j - b_1)^2 + (\sum_{j=1}^n W_{1j} y_j)^2 + 2 \sum_{j=1}^n W_{1j} y_j b_1}{2\sigma_1^2}\right) dx'_1}$$

Again, we can simplify by removing terms that do not depend on  $x'_1$  or  $x_1$  by using the exponential rule and taking them out of the integral in the denominator:

$$\begin{aligned}
&= \frac{\exp\left(\frac{-(x_1 - \sum_{j=1}^n W_{1j} y_j - b_1)^2}{2\sigma_1^2}\right)}{\int_{x'_1} \exp\left(\frac{-(x'_1 - \sum_{j=1}^n W_{1j} y_j - b_1)^2}{2\sigma_1^2}\right) dx'_1} \\
&= \frac{\exp\left(-\frac{1}{2} \left(\frac{x_1 - \sum_{j=1}^n W_{1j} y_j - b_1}{\sigma_1}\right)^2\right)}{\int_{x'_1} \exp\left(-\frac{1}{2} \left(\frac{x'_1 - \sum_{j=1}^n W_{1j} y_j - b_1}{\sigma_1}\right)^2\right) dx'_1}
\end{aligned}$$

Now, this value looks like a PDF where the denominator act as a normalizer that ensures that the integral of the PDF is equal to 1. You can observe that the numerator is an unnormalized Gaussian, therefore the normalization term is equal to  $\sigma_1 \sqrt{2\pi}$ . This results can also be shown using closed form formula for Gaussian function integration.<sup>1</sup>By rewriting the denominator as  $\sigma_1 \sqrt{2\pi}$  we have:

$$\begin{aligned}
&= \frac{1}{\sigma_1 \sqrt{2\pi}} \exp\left(-\frac{1}{2} \left(\frac{x_1 - (\sum_{j=1}^n W_{1j} y_j + b_1)}{\sigma_1}\right)^2\right) \\
&= \mathcal{N}\left(x_1 \middle| b_1 + \sum_{j=1}^n W_{1j} y_j, \sigma_1^2\right)
\end{aligned}$$

that is  $p(x_1|x_2, \mathbf{y}; \theta) = p(x_1|\mathbf{y}; \theta)$  is a Gaussian with mean  $b_1 + \sum_{j=1}^n W_{1j} y_j$  and standard deviation  $\sigma_1^2$ . We can do a similar computation for the other observed variable:

$$p(x_2|x_1, \mathbf{y}; \theta) = p(x_2|\mathbf{y}; \theta) = \mathcal{N}\left(x_2 \middle| b_2 + \sum_{j=1}^n W_{2j} y_j, \sigma_2^2\right)$$

<sup>1</sup>See [https://en.wikipedia.org/wiki/Gaussian\\_integral](https://en.wikipedia.org/wiki/Gaussian_integral)

Similarly, we can compute the conditional probability of an observation. If we take the specific case  $p(y_1|x, y_2\dots y_n; \theta)$ :

$$\begin{aligned} p(y_1|x, y_2\dots y_n; \theta) &= \frac{\exp(w(x, [y_1, y_2\dots y_n]; \theta))}{\sum_{y'_1} \exp(w(x, [y'_1, y_2\dots y_n]; \theta))} \\ &= \frac{\exp\left(-\sum_{i=1}^2 \frac{(x_i - b_i)^2}{2\sigma_i^2} + \sum_{i=1}^2 \sum_{j=1}^n \frac{x_i}{\sigma_i^2} W_{ij} y_j + \sum_{j=1}^n d_j y_j\right)}{\sum_{y'_1} \exp\left(-\sum_{i=1}^2 \frac{(x_i - b_i)^2}{2\sigma_i^2} + \sum_{i=1}^2 \frac{x_i}{\sigma_i^2} W_{i1} y'_1 + d_1 y'_1 + \sum_{i=1}^2 \sum_{j=2}^n \frac{x_i}{\sigma_i^2} W_{ij} y_j + \sum_{j=2}^n d_j y_j\right)} \end{aligned}$$

As previously, there are many terms that cancels:

$$= \frac{\exp\left(\sum_{i=1}^2 \frac{x_i}{\sigma_i^2} W_{i1} y_1 + d_1 y_1\right)}{\sum_{y'_1} \exp\left(\sum_{i=1}^2 \frac{x_i}{\sigma_i^2} W_{i1} y'_1 + d_1 y'_1\right)}$$

The variable  $y'_1$  can either be 0 or 1 in the denominator's sum. In the case it is 0, the exponential simplifies as  $\exp(0) = 1$ . Moreover, we can factorize by  $y_1$  in the numerator's exponential, therefore:

$$= \frac{\exp\left(y_1 \left(\sum_{i=1}^2 \frac{x_i}{\sigma_i^2} W_{i1} + d_1\right)\right)}{1 + \exp\left(\sum_{i=1}^2 \frac{x_i}{\sigma_i^2} W_{i1} + d_1\right)}$$

Hence, we can see that  $y_1$  is conditionally independent of other latent variables  $p(y_1|x, y_2\dots y_n; \theta) = p(y_1|x; \theta)$ . Moreover, from the formula we deduce that  $p(y_1|x; \theta)$  is a Bernoulli distribution with parameter  $\mu = \text{sigmoid}\left(\sum_{i=1}^2 \frac{x_i}{\sigma_i^2} W_{i1} + d_1\right)$  which we write as:

$$= \mathcal{B}\left(y_1 \left| \text{sigmoid}\left(d_1 + \sum_{i=1}^2 \frac{x_i}{\sigma_i^2} W_{i1}\right)\right.\right)$$

A similar computation can be done for other latent variables. Note that for both condition distributions over observed and latent variables we can use tensor operations instead of computing parameters for each variable independently. Try to visualize how does it works!

### 3 Training

The parameter  $\theta$  must be set so that the probability of the training data is maximized. To do this, we reformulate the probability of observation using the log-partition function of  $p(\mathbf{y}|\mathbf{x})$ , denoted  $c(\mathbf{x}; \theta)$ , which can be computed in closed form:

$$\begin{aligned} p(\mathbf{x}; \theta) &= \sum_{\mathbf{y}} p(\mathbf{x}, \mathbf{y}; \theta) \\ &= \sum_{\mathbf{y}} \frac{\exp(w(\mathbf{x}, \mathbf{y}; \theta))}{\int_{\mathbf{x}'} \sum_{\mathbf{y}'} \exp(w(\mathbf{x}', \mathbf{y}'; \theta)) d\mathbf{x}'} \\ &= \frac{\sum_{\mathbf{y}} \exp\left(-\sum_{i=1}^2 \frac{(x_i - b_i)^2}{2\sigma_i^2} + \sum_{i=1}^2 \sum_{j=1}^n \frac{x_i}{\sigma_i^2} W_{ij} y_j + \sum_{j=1}^n d_j y_j\right)}{\int_{\mathbf{x}'} \sum_{\mathbf{y}'} \exp\left(-\sum_{i=1}^2 \frac{(x'_i - b_i)^2}{2\sigma_i^2} + \sum_{i=1}^2 \sum_{j=1}^n \frac{x'_i}{\sigma_i^2} W_{ij} y'_j + \sum_{j=1}^n d_j y'_j\right) d\mathbf{x}'} \end{aligned}$$

We can use the fact that  $\exp(a + b) = \exp(a) \exp(b)$  to rewrite both exponentials:

$$= \frac{\sum_{\mathbf{y}} \left[ \exp\left(-\sum_{i=1}^2 \frac{(x_i - b_i)^2}{2\sigma_i^2}\right) \prod_{j=1}^n \exp\left(d_j y_j + \sum_{i=1}^2 \frac{x_i}{\sigma_i^2} W_{ij} y_j\right) \right]}{\int_{\mathbf{x}'} \left[ \sum_{\mathbf{y}'} \exp\left(-\sum_{i=1}^2 \frac{(x'_i - b_i)^2}{2\sigma_i^2}\right) \prod_{j=1}^n \left(d_j y'_j + \sum_{i=1}^2 \frac{x'_i}{\sigma_i^2} W_{ij} y'_j\right) \right] d\mathbf{x}'}$$

Note that the first exponential term in the nominator does not depend on  $\mathbf{y}$ , so we can take it out of the sum, and similarly in the denominator:

$$= \frac{\exp\left(-\sum_{i=1}^2 \frac{(x_i - b_i)^2}{2\sigma_i^2}\right) \sum_{\mathbf{y}} \left[ \prod_{j=1}^n \exp\left(d_j y_j + \sum_{i=1}^2 \frac{x_i}{\sigma_i^2} W_{ij} y_j\right) \right]}{\int_{\mathbf{x}'} \exp\left(-\sum_{i=1}^2 \frac{(x'_i - b_i)^2}{2\sigma_i^2}\right) \sum_{\mathbf{y}'} \left[ \prod_{j=1}^n \left(d_j y'_j + \sum_{i=1}^2 \frac{x'_i}{\sigma_i^2} W_{ij} y'_j\right) \right] d\mathbf{x}'}$$

Now we can rewrite using the distributivity of multiplication. The sum will simplify as a sum over a single latent variable, see the course for a more detailed explanation:

$$= \frac{\exp\left(-\sum_{i=1}^2 \frac{(x_i - b_i)^2}{2\sigma_i^2}\right) \prod_{j=1}^n \left[ \sum_{y_j} \exp\left(d_j y_j + \sum_{i=1}^2 \frac{x_i}{\sigma_i^2} W_{ij} y_j\right) \right]}{\int_{\mathbf{x}'} \exp\left(-\sum_{i=1}^2 \frac{(x'_i - b_i)^2}{2\sigma_i^2}\right) \prod_{j=1}^n \left[ \sum_{y'_j} \left(d_j y'_j + \sum_{i=1}^2 \frac{x'_i}{\sigma_i^2} W_{ij} y'_j\right) \right] d\mathbf{x}'}$$

The latent variable  $y_j$  can only take two values (0 and 1) in the sum in the numerator. When  $y_j = 0$ , we have  $\exp(0) = 1$ . By applying the same trick in the nominator and denominator, we have:

$$= \frac{\exp\left(-\sum_{i=1}^2 \frac{(x_i - b_i)^2}{2\sigma_i^2}\right) \prod_{j=1}^n \left[1 + \exp\left(d_j + \sum_{i=1}^2 \frac{x_i}{\sigma_i^2} W_{ij}\right)\right]}{\int_{\mathbf{x}'} \exp\left(-\sum_{i=1}^2 \frac{(x'_i - b_i)^2}{2\sigma_i^2}\right) \prod_{j=1}^n \left[1 + \exp\left(d_j + \sum_{i=1}^2 \frac{x'_i}{\sigma_i^2} W_{ij}\right)\right] d\mathbf{x}'}$$

Now observe that if we can take the exponential-log of the numerator and then use the fact that  $\log(ab) = \log a + \log b$  to change the product by a sum (moreover the log-exp in the first term then cancel - again, check the course for an explanation of this trick) we have:

$$\begin{aligned} &= \frac{\exp\left(-\sum_{i=1}^2 \frac{(x_i - b_i)^2}{2\sigma_i^2} + \sum_{j=1}^n \log\left[1 + \exp\left(d_j + \sum_{i=1}^2 \frac{x_i}{\sigma_i^2} W_{ij}\right)\right]\right)}{\int_{\mathbf{x}'} \exp\left(-\sum_{i=1}^2 \frac{(x'_i - b_i)^2}{2\sigma_i^2}\right) + \sum_{j=1}^n \log\left[1 + \exp\left(d_j + \sum_{i=1}^2 \frac{x'_i}{\sigma_i^2} W_{ij}\right)\right] d\mathbf{x}'} \\ &= \frac{\exp(c(\mathbf{x}; \theta))}{\int_{\mathbf{x}'} \exp(c(\mathbf{x}'; \theta)) d\mathbf{x}'} \end{aligned}$$

where  $c(\mathbf{x}; \theta) = \sum_{i=1}^2 \frac{(x_i - b_i)^2}{2\sigma_i^2} - \sum_{j=1}^n \log\left[1 + \exp\left(d_j + \sum_{i=1}^2 \frac{x_i}{\sigma_i^2} W_{ij}\right)\right]$  is the log-partition function of  $p(\mathbf{y}|\mathbf{x})$ . Note that  $c(\mathbf{x}; \theta)$  can be computed efficiently because we transformed the original sum over an exponential number of values  $\sum_{\mathbf{y}}$  (sum over  $2^n$  possible assignation for  $\mathbf{y}$ !) as  $n$  sums. The marginalization in the denominator is still intractable, but during training we will approximate its contribution to the gradient via Monte Carlo estimation.

Let  $D = \{\mathbf{x}^{(1)} \dots \mathbf{x}^{(d)}\}$  be the training set. During training, we seek to maximize the log probability of the dataset:

$$\begin{aligned} \hat{\theta} &= \arg \max_{\theta} \sum_{i=1}^d \log p(\mathbf{x}^{(i)}; \theta) \\ &= \arg \min_{\theta} \sum_{i=1}^d \mathcal{L}(\theta, \mathbf{x}^{(i)}) \end{aligned}$$

where  $\mathcal{L}(\theta, \mathbf{x}) = -\log p(\mathbf{x}; \theta)$  is the loss function. The gradient of the loss function is defined as:

$$\begin{aligned} \nabla \mathcal{L}(\theta, \mathbf{x}) &= \nabla -\log p(\mathbf{x}; \theta) \\ &= \nabla -\log \frac{\exp(c(\mathbf{x}; \theta))}{\int_{\mathbf{x}'} \exp(c(\mathbf{x}'; \theta)) d\mathbf{x}'} \\ &= -\nabla c(\mathbf{x}; \theta) + \nabla \log \left( \int_{\mathbf{x}'} \exp(c(\mathbf{x}'; \theta)) d\mathbf{x}' \right) \\ &= -\nabla c(\mathbf{x}; \theta) + \frac{1}{\int_{\mathbf{x}'} \exp(c(\mathbf{x}'; \theta)) d\mathbf{x}'} \nabla \left( \int_{\mathbf{x}'} \exp(c(\mathbf{x}'; \theta)) d\mathbf{x}' \right) \end{aligned}$$

Similarly to sums, the gradient of an integral is the integral of the gradient:

$$\begin{aligned} &= -\nabla c(\mathbf{x}; \theta) + \frac{1}{\int_{\mathbf{x}'} \exp(c(\mathbf{x}'; \theta)) d\mathbf{x}'} \int_{\mathbf{x}'} \nabla \exp(c(\mathbf{x}'; \theta)) d\mathbf{x}' \\ &= -\nabla c(\mathbf{x}; \theta) + \frac{1}{\int_{\mathbf{x}'} \exp(c(\mathbf{x}'; \theta)) d\mathbf{x}'} \int_{\mathbf{x}'} \exp(c(\mathbf{x}'; \theta)) \nabla c(\mathbf{x}'; \theta) d\mathbf{x}' \\ &= -\nabla c(\mathbf{x}; \theta) + \int_{\mathbf{x}'} \frac{\exp(c(\mathbf{x}'; \theta))}{\int_{\mathbf{x}''} \exp(c(\mathbf{x}''; \theta)) d\mathbf{x}''} \nabla c(\mathbf{x}'; \theta) d\mathbf{x}' \\ &= -\nabla c(\mathbf{x}; \theta) + \int_{\mathbf{x}'} p(\mathbf{x}'; \theta) \nabla c(\mathbf{x}'; \theta) d\mathbf{x}' \\ &= -\nabla c(\mathbf{x}; \theta) + \mathbb{E}_{\mathbf{x}' \sim p(\mathbf{x}'; \theta)} [\nabla c(\mathbf{x}'; \theta)] \end{aligned}$$

Computing the expectation exactly is intractable but we can rely on Monte Carlo estimation:

$$\simeq -\nabla c(\mathbf{x}; \theta) + \frac{1}{K} \sum_{k=1}^K \nabla c(\tilde{\mathbf{x}}^{(k)}; \theta)$$

where  $\tilde{\mathbf{x}}^{(k)} \sim p(\tilde{\mathbf{x}}^{(k)}; \theta)$  are  $K$  samples from the model distribution. To sample from the distribution, we can rely on Gibbs sampling Markov Chain Monte Carlo as it is easy to sample from both conditional distributions. In practice, we will rely on the contrastive divergence objective, i.e. we take a single sample using a single step of a Markov chain initialized with the datapoint  $\mathbf{x}$ .