

1 Derivation RBM negative phase

The 'cost function' of the RBM is given by the Kullback-Leibler divergence,

$$C_{\lambda} = D_{KL}(q||p_{\lambda}) = \sum_{\mathbf{v}} q(\mathbf{v}) \log \left(\frac{q(\mathbf{v})}{p_{\lambda}(\mathbf{v})} \right) = \sum_{\mathbf{v}} q(\mathbf{v}) \log(q(\mathbf{v})) - \sum_{\mathbf{v}} q(\mathbf{v}) \log(p_{\lambda}(\mathbf{v})) \quad (1)$$

The first term of the last equation is the Shannon entropy of the data distribution and is not dependent on the model parameters λ and will therefore not matter if we take the derivative with respect to these parameters. The second term of the last equation can be rewritten as,

$$- \sum_{\mathbf{v}} q(\mathbf{v}) \log(p_{\lambda}(\mathbf{v})) = -\langle \log(p_{\lambda}(\mathbf{v})) \rangle_{q(\mathbf{v})}, \quad (2)$$

which is the expectation value of the log-likelihood of $p_{\lambda}(\mathbf{v})$ w.r.t. the data distribution $q(\mathbf{v})$. The calculation can be simplified, by just maximizing the negative log-likelihood. Which means nothing else than making the training data the most likeli. The gradient of the likelihood (resp. KL-divergence) of a single training example \mathbf{v} reads:

$$-\nabla_{\lambda} \log(p_{\lambda}(\mathbf{v})) = \nabla_{\lambda} \mathcal{E}_{\lambda}(\mathbf{v}) + \nabla_{\lambda} \log(Z_{\lambda}) = \nabla_{\lambda} \mathcal{E}_{\lambda}(\mathbf{v}) + \frac{1}{Z_{\lambda}} \nabla_{\lambda} \sum_{\mathbf{v}', \mathbf{h}} e^{-E(\mathbf{v}', \mathbf{h})} = \quad (3)$$

$$\nabla_{\lambda} \mathcal{E}_{\lambda}(\mathbf{v}) - \frac{1}{Z_{\lambda}} \sum_{\mathbf{v}', \mathbf{h}} e^{-E(\mathbf{v}', \mathbf{h})} \nabla_{\lambda} E(\mathbf{v}', \mathbf{h}) = \nabla_{\lambda} \mathcal{E}_{\lambda}(\mathbf{v}) - \sum_{\mathbf{v}', \mathbf{h}} p_{\lambda}(\mathbf{v}', \mathbf{h}) \nabla_{\lambda} E(\mathbf{v}', \mathbf{h}) \quad (4)$$

Like in normal neural networks we would like to calculate the derivative for all data \mathbf{v} . But this is generally not feasible and we calculate it for batches \mathcal{D} of the data distribution $q(\mathbf{v})$. It is very important to note, that the second term of the last equation is already summed over all \mathbf{v}' . Therefore this part is already averaged over all possible configurations \mathbf{v} and $\frac{1}{|\mathcal{D}|} \sum_{\mathbf{v} \in \mathcal{D}}$ cancels for this part.

$$-\nabla_{\lambda} \langle \log(p_{\lambda}(\mathbf{v})) \rangle_{q(\mathbf{v})} \approx \frac{1}{|\mathcal{D}|} \sum_{\mathbf{v} \in \mathcal{D}} \nabla_{\lambda} \mathcal{E}_{\lambda}(\mathbf{v}) - \sum_{\mathbf{v}', \mathbf{h}} p_{\lambda}(\mathbf{v}', \mathbf{h}) \nabla_{\lambda} E(\mathbf{v}', \mathbf{h}) \quad (5)$$

In the second term we can again contract the variable \mathbf{h} and we get:

$$-\nabla_{\lambda} \langle \log(p_{\lambda}(\mathbf{v})) \rangle_{q(\mathbf{v})} \approx \frac{1}{|\mathcal{D}|} \sum_{\mathbf{v} \in \mathcal{D}} \nabla_{\lambda} \mathcal{E}_{\lambda}(\mathbf{v}) - \sum_{\mathbf{v}} p_{\lambda}(\mathbf{v}) \nabla_{\lambda} \mathcal{E}(\mathbf{v}) = \quad (6)$$

$$\langle \nabla_{\lambda} \mathcal{E}_{\lambda}(\mathbf{v}) \rangle_{\mathcal{D}(\mathbf{v})} - \langle \nabla_{\lambda} \mathcal{E}_{\lambda}(\mathbf{v}) \rangle_{p_{\lambda}(\mathbf{v})} \quad (7)$$