

**Exercise 1: Mixture Density Networks (20 + 10 P)**

In this exercise, we prove some of the results from the paper Mixture Density Networks by Bishop (1994). The mixture density network is given by

$$p(\mathbf{t}|\mathbf{x}) = \sum_{i=1}^m \alpha_i(\mathbf{x}) \phi_i(\mathbf{t}|\mathbf{x})$$

with the mixture elements

$$\phi_i(\mathbf{t}|\mathbf{x}) = \frac{1}{(2\pi)^{c/2} \sigma_i(\mathbf{x})^c} \exp\left(-\frac{\|\mathbf{t} - \boldsymbol{\mu}_i(\mathbf{x})\|^2}{2\sigma_i(\mathbf{x})^2}\right).$$

The contribution to the error function of one data point  $q$  is given by

$$E^q = -\log\left\{\sum_{i=1}^m \alpha_i(\mathbf{x}^q) \phi_i(\mathbf{t}^q|\mathbf{x}^q)\right\}$$

$$\frac{\partial \|\mathbf{x} - \mathbf{y}\|^2}{\partial \mathbf{x}} = 2(\mathbf{x} - \mathbf{y})$$
$$\frac{\partial \|\mathbf{x} - \mathbf{y}\|^2}{\partial \mathbf{y}} = 2(\mathbf{y} - \mathbf{x})$$

We also define the posterior distribution

$$\pi_i(\mathbf{x}, \mathbf{t}) = \frac{\alpha_i \phi_i}{\sum_{j=1}^m \alpha_j \phi_j}$$

which is obtained using the Bayes theorem.

(a) Compute the gradient of the error  $E^q$  w.r.t. the mixture parameters, i.e. show that

(i)  $\frac{\partial E^q}{\partial \alpha_i} = -\frac{\pi_i}{\alpha_i}$

$$\frac{\partial E^q}{\partial \alpha_i} = - \frac{1}{\sum_j \alpha_j \phi_j} \cdot \phi_i \cdot \frac{\alpha_i}{\alpha_i} \pi_i = - \frac{\pi_i}{\alpha_i}$$

(ii)  $\frac{\partial E^q}{\partial \mu_{ik}} = \pi_i \left( \frac{\mu_{ik} - t_k}{\sigma_i^2} \right)$

$$\frac{\partial E}{\partial \mu_{ik}} = \frac{\alpha_i}{\sum_j \alpha_j \phi_j} \cdot \phi_i \cdot \left( \frac{\mu_{ik} - t_k}{\sigma_i^2} \right) = \pi_i \left( \frac{\mu_{ik} - t_k}{\sigma_i^2} \right)$$

(b) We now assume that the neural network produces the mixture coefficients as:

$$\alpha_i = \frac{\exp(z_i^\alpha)}{\sum_{j=1}^M \exp(z_j^\alpha)}$$

$$\left(\frac{f}{g}\right)' = \frac{f'g - fg'}{g^2}$$

where  $z^\alpha$  denotes the outputs of the neural network (after the last linear layer) associated to these mixture coefficients. Compute using the chain rule for derivatives (i.e. by reusing some of the results in the first part of this exercise) the derivative  $\partial E^q / \partial z_i^\alpha$ .

$$\frac{\partial E^q}{\partial z_i^\alpha} = \sum_j \frac{\partial E^q}{\partial \alpha_j} \cdot \frac{\partial \alpha_j}{\partial z_i^\alpha}$$
$$\left[ \frac{\partial \alpha_j}{\partial z_i^\alpha} = \frac{\delta_{ij} \exp(z_i) \cdot (\sum_j \exp(z_j)) - \exp(z_j) \exp(z_i)}{(\sum_j \exp(z_j))^2} \right]$$
$$= \delta_{ij} \alpha_j - \alpha_i \alpha_j$$
$$= \sum_j -\frac{\pi_j}{\alpha_j} (\delta_{ij} \alpha_j - \alpha_i \alpha_j) = \sum_j \left( -\frac{\pi_j}{\alpha_j} \delta_{ij} \alpha_j + \frac{\pi_j}{\alpha_j} \alpha_i \alpha_j \right) = -\pi_i + (\sum_j \pi_j) \alpha_i = -\pi_i + \alpha_i$$

**Exercise 2: Conditional RBM (20 + 10 P)**

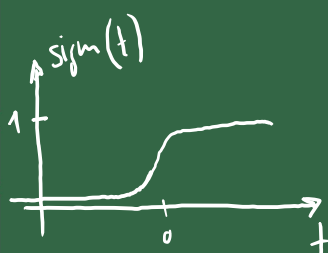
The conditional restricted Boltzmann machine is a system of binary variables comprising inputs  $\mathbf{x} \in \{0, 1\}^d$ , outputs  $\mathbf{y} \in \{0, 1\}^c$ , and hidden units  $\mathbf{h} \in \{0, 1\}^K$ . It associates to each configuration of these binary variables the energy:

$$E(\mathbf{x}, \mathbf{y}, \mathbf{h}) = -\mathbf{x}^\top \mathbf{W} \mathbf{h} - \mathbf{y}^\top \mathbf{U} \mathbf{h} = \underbrace{\sum_k -x^\top W_{:,k} h_k - y^\top U_{:,k} h_k}_{E(\mathbf{x}, \mathbf{y}, \mathbf{h}_k)}$$

and the probability associated to each configuration is then given as:

$$p(\mathbf{x}, \mathbf{y}, \mathbf{h}) = \frac{1}{Z} \exp(-E(\mathbf{x}, \mathbf{y}, \mathbf{h}))$$

where  $Z$  is a normalization constant that makes probabilities sum to one.



(a) Let  $\text{sigm}(t) = \exp(t)/(1 + \exp(t))$  be the sigmoid function. Show that

(i)  $p(h_k = 1 | \mathbf{x}, \mathbf{y}) = \text{sigm}(\mathbf{x}^\top \mathbf{W}_{:,k} + \mathbf{y}^\top \mathbf{U}_{:,k})$

$$p(h_k = 1 | \mathbf{x}, \mathbf{y}) = \frac{p(h_k = 1, \mathbf{x}, \mathbf{y})}{p(\mathbf{x}, \mathbf{y})}$$

product rule

$$= \frac{\sum_{h_{-k}} p(h_k = 1, \mathbf{x}, \mathbf{y}, h_{-k})}{\sum_{q \in \{0,1\}} \sum_{h_{-k}} p(h_k = q, \mathbf{x}, \mathbf{y}, h_{-k})}$$

marginalization

$$= \frac{\sum_{h_{-k}} \exp(x^\top W_{:,k} 1 + y^\top U_{:,k} 1 - E(\mathbf{x}, \mathbf{y}, h_{-k}))}{\sum_{q \in \{0,1\}} \sum_{h_{-k}} \exp(x^\top W_{:,k} q + y^\top U_{:,k} q - E(\mathbf{x}, \mathbf{y}, h_{-k}))}$$
$$= \frac{\exp(x^\top U_{:,k} 1 + y^\top U_{:,k} 1) \cdot \sum_{h_{-k}} \exp(-E(\mathbf{x}, \mathbf{y}, h_{-k}))}{\sum_{q \in \{0,1\}} \sum_{h_{-k}} \exp(x^\top W_{:,k} q + y^\top U_{:,k} q) \cdot \exp(-E(\mathbf{x}, \mathbf{y}, h_{-k}))}$$
$$= \frac{\exp(x^\top W_{:,k} 1 + y^\top U_{:,k} 1)}{1 + \exp(x^\top W_{:,k} 1 + y^\top U_{:,k} 1)} = \text{sigm}(x^\top W_{:,k} + y^\top U_{:,k})$$

all hidden units except  $h_k$

$$\sum_{h_{-k}} : \sum_{q_1 \in \{0,1\}} \sum_{q_2 \in \{0,1\}} \sum_{q_3 \in \{0,1\}} \dots$$

(ii)  $p(y_j = 1 | \mathbf{h}, \mathbf{x}) = \text{sigm}(\mathbf{U}_{j,:}^\top \mathbf{h})$

(b) Show that

$$p(\mathbf{x}, \mathbf{y}) = \frac{1}{Z} \exp(-F(\mathbf{x}, \mathbf{y}))$$

where

$$F(\mathbf{x}, \mathbf{y}) = -\sum_{k=1}^K \log(1 + \exp(\mathbf{x}^\top \mathbf{W}_{:,k} + \mathbf{y}^\top \mathbf{U}_{:,k}))$$

is the free energy and where  $Z$  is again a normalization constant.

$$Z \cdot p(\mathbf{x}, \mathbf{y}) = \sum_{\mathbf{h}} p(\mathbf{x}, \mathbf{y}, \mathbf{h}) \cdot Z$$
$$= \sum_{\mathbf{h}} \exp(-E(\mathbf{x}, \mathbf{y}, \mathbf{h})) = \sum_{\mathbf{h}} \exp(-\sum_k E(\mathbf{x}, \mathbf{y}, h_k))$$
$$= \sum_{\mathbf{h}} \prod_k \exp(-E(\mathbf{x}, \mathbf{y}, h_k)) \quad \sum_i \pi_j \alpha_{ij} = \pi_i \sum_j \alpha_{ij}$$
$$= \prod_k \sum_{q_k \in \{0,1\}} \exp(-E(\mathbf{x}, \mathbf{y}, h_k = q_k))$$
$$= \prod_k (1 + \exp(-E(\mathbf{x}, \mathbf{y}, h_k = 1)))$$
$$= \exp\left(\log\left(\prod_k (1 + \exp(-E(\mathbf{x}, \mathbf{y}, h_k = 1)))\right)\right)$$
$$= \exp\left(\sum_k \log(1 + \exp(-E(\mathbf{x}, \mathbf{y}, h_k = 1)))\right)$$
$$= \exp(-F(\mathbf{x}, \mathbf{y}))$$