

Exercise Sheet 1

Exercise 1: Log-Likelihood (10 + 10 P)

Let μ be the output of a neural network predicting the mean of some Laplacian distribution

$$p(t; \mu) = \frac{1}{2\sigma} \exp\left(-\frac{|t - \mu|}{\sigma}\right).$$

We collect a dataset composed of neural network inputs $\mathbf{x}_1, \dots, \mathbf{x}_N$ and associated targets t_1, \dots, t_N . We denote by μ_1, \dots, μ_N the predictions of the neural network for these points. We would like to optimize the neural network to produce predictions that match the targets. For this, we use the negative log-likelihood loss function given by:

$$\ell(\mu, t) = -\log p(t; \mu).$$

- (a) *Compute* the gradient of the loss function with respect to the output μ of the neural network.
- (b) We now assume that the neural network is made of a single linear layer: $\mu = \mathbf{w}^\top \mathbf{x}$, where \mathbf{w} is the vector of parameters to be learned. *State* the chain rule for transmitting the gradient from the output of the neural network to the model parameters, and use it to compute the gradient of the loss function with respect to \mathbf{w} .

Exercise 2: Shared Parameters (10 + 10 + 10 P)

Let x_1, x_2 be two observed variables. Consider the two-layer neural network that takes these two variables as input and builds the prediction y by computing iteratively:

$$z_3 = x_1 w_{13}, \quad z_4 = x_2 w_{24}, \quad a_3 = 0.5 z_3^2, \quad a_4 = 0.5 z_4^2, \quad y = a_3 + a_4.$$

- (a) *Draw* the neural network graph associated to these computations.

We now consider the loss function $\ell(y, t) = \frac{1}{2}(y - t)^2$ where t is a target variable that the neural network learns to approximate.

- (b) Using the rules for backpropagation, compute the derivatives $\partial\ell/\partial w_{13}$ and $\partial\ell/\partial w_{24}$ required for gradient descent.

Let us now assume that w_{13} and w_{24} cannot be adapted freely, but are a function of the same shared parameter v :

$$w_{13} = \log(1 + \exp(v)) \quad \text{and} \quad w_{24} = -\log(1 + \exp(-v))$$

- (c) Assume you have already computed $\partial\ell/\partial w_{13}$ and $\partial\ell/\partial w_{24}$. *State* the multivariate chain rule that links the gradient $\partial\ell/\partial v$ to these quantities and *apply* it in order to get an analytic expression of $\partial\ell/\partial v$.

Exercise 3: Layered Networks (10 + 10 P)

Consider a deep neural network made of several layers of neurons. Let $\mathbf{a}^{(l)} = (a_1^{(l)}, \dots, a_d^{(l)})$ be the activations of the d neurons at layer (l) , and $\mathbf{a}^{(l+1)} = (a_1^{(l+1)}, \dots, a_h^{(l+1)})$ be the activations of the h neurons at layer $(l+1)$. These activations are computed as:

$$\forall_{j=1}^h : z_j^{(l+1)} = \sum_{i=1}^d a_i^{(l)} w_{ij} + b_j \quad a_j^{(l+1)} = \exp(z_j^{(l+1)}) / (1 + \exp(z_j^{(l+1)}))$$

where w_{ij}, b_j denote the parameters of the mapping between these two layers.

Show by application of the chain rule that the following relations hold:

- (a) $\frac{\partial\ell}{\partial z_j^{(l+1)}} = a_j^{(l+1)} \cdot (1 - a_j^{(l+1)}) \cdot \frac{\partial\ell}{\partial a_j^{(l+1)}}.$
- (b) $\frac{\partial\ell}{\partial \mathbf{a}^{(l)}} = W \cdot (\mathbf{a}^{(l+1)} \odot (1 - \mathbf{a}^{(l+1)})) \odot \frac{\partial\ell}{\partial \mathbf{a}^{(l+1)}}.$

where W is a matrix of size $d \times h$ containing the parameters w_{ij} , and where \odot denotes the element-wise product.

Exercise 4: Programming (30 P)

Download the programming files on ISIS and follow the instructions.