

## Machine Learning Worksheet 07

### Neural Networks 1

**Problem 1:** Consider a two-layer neural network in which the hidden-unit nonlinear sigmoid activation functions  $\phi(\cdot)$  are given by

$$\sigma(x) = \frac{1}{1 + \exp(-x)}. \quad (1)$$

Show that there exists an equivalent network, which computes exactly the same function, but with hidden unit activation functions given by  $\tanh(x)$ .

**Problem 2:** Show that the derivative of the sigmoid activation function given in Eq. (1) can be expressed in terms of the function value itself, thus allowing for efficient implementation. Also derive the corresponding result for the tanh activation function.

**Problem 3:** If we have multiple target variables  $\mathbf{z}$ , and we assume that they are independent conditional on  $\mathbf{x}$  and  $\mathbf{w}$  with shared noise precision  $\beta$  then the conditional distribution of the target values is given by

$$p(\mathbf{z} \mid \mathbf{x}, \mathbf{w}) = \mathcal{N}(\mathbf{z} \mid \mathbf{y}(\mathbf{x}, \mathbf{w}), \beta^{-1} \mathbf{I})$$

Show that maximising the resulting likelihood function under the above conditional distribution for a multi-output neural network is equivalent to minimising a sum-of-squares error function.

**Problem 4:** You know that the sum of squared errors is related to the Gaussian distribution—differently put, if you assume a normal distribution of the data around their expectation, the maximum likelihood estimate (MLE) is reached when the summed squared errors is minimised.

The same is true for a Laplace distribution and the sum of absolute errors. In particular, if the data observes a Laplacian distribution

$$p(\mathbf{z} \mid \mathbf{x}, \mathbf{w}, \beta) = \prod_{n=1}^N p(z_n \mid x_n, \mathbf{w}, \beta) = \prod_{n=1}^N \frac{1}{2\beta} \exp\left(-\frac{|z_n - y(x_n, \mathbf{w})|}{\beta}\right)$$

then minimising the summed absolute errors

$$\sum_{n=1}^N |z_n - y(x_n, \mathbf{w})|$$

leads to MLE. In these equations,  $\mathbf{x}$  is the vector of all inputs,  $x_n$  is the input of sample  $n$ , while  $y(x_n, \mathbf{w})$  is the neural network prediction on  $x_n$ . Then,  $z_n$  is the desired output for  $x_n$ .

Show that the MLE of the Laplace distribution minimises the sum of absolute errors.

**Problem 5:** Experiment with the uploaded Jupyter notebook `neuralnetworks1.ipynb`.

Plot several (at least 6 total) training curves for different values of the learning rate, when classifying XOR with `mlp_xor.NeuralNetwork(X,y,0.1)` as well as the  $\sin(x)$  data with `mlp_sin.NeuralNetwork(X,y,0.1)`.

**Problem 6:** Look at the uploaded Jupyter notebook `neuralnetworks1.ipynb`. In the definition of the input pattern,

```
X = np.array([ [0,0,1],
               [0,1,1],
               [1,0,1],
               [1,1,1] ])
```

why is the input in the last column always 1?

**Problem 7:** Look at the uploaded Jupyter notebook `neuralnetworks1.ipynb`. What is wrong with the use of this notebook for the XOR problem?