# EQ2415 – Machine Learning and Data Science HT22

Tutorial 4                                                                  A. Honoré, A. Ghosh

## 1   Neural networks

### 1.1   Loss function for regression and classification

**Question 1.** Suppose that $(\mathbf{t}, \mathbf{x}) \in \mathbb{R}^q \times \mathbb{R}^d$ is a target/input pair for a regression problem, $\mathbf{f}$ is a neural network (NN) model parameterized with a vector $\mathbf{w}$.

**Question 1a.** Write the likelihood function for the distribution of a targets conditional on $\mathbf{x}$ and $\mathbf{w}$. Assume that the noise of the model is Gaussian, that the targets component are independent and that all components share the same noise precision $\beta$.

**Question 1b.** Suppose that you are given a dataset consisting of $n$ independent target/input pairs. Show that maximizing the likelihood of the dataset wrt $\mathbf{w}$ is equivalent to minimizing the MSE wrt $\mathbf{w}$.

**Question 2.** (Bishop 5.4) Suppose you are given a binary classification task. You are given a set of $n$ independent training data points $\mathcal{D} = \{(\mathbf{x}_{(i)}, y_{(i)})\}_{i=1}^n$, where $\mathbf{x}_{(i)} \in \mathbb{R}^d$ and $y_{(i)} \in \{0, 1\}$. In general, for a binary classifier, the probability that an input $\mathbf{x}$ is classified with label $y = 1$ is expressed

$$p(y = 1|\mathbf{x}) = y_W(\mathbf{x}), \tag{1}$$

where $y_W : \mathbb{R}^d \to [0, 1]$ is a function of the input $\mathbf{x}$ parameterized with $W$. Importantly, you are told that the data is mislabeled with probability $\epsilon$. To model this situation, we can introduce a binary random variable modeling the true and unobserved label $y_r$ of an input $\mathbf{x}$. We can also introduce an unobserved binary random variable $m$ associated with each label, indicating whether the label is true or false.

**Question 2a.** How would you introduce the probability of mislabeling in the output of your classifier ?

**Question 2b.** Write the distribution of the Bernoulli variable $y|\mathbf{x}$.

**Question 2c.** Show that the negative likelihood function on the dataset corresponds to the cross entropy function when $\epsilon = 0$.

### 1.2   Standard results on activation functions

Let us consider a real valued functions $h : \mathbb{R} \to \mathbb{R}$. We calculate the derivative of $h$ wrt to its argument $x \in \mathbb{R}$ for different values of $h$.

**Question 1.**

$$h(x) = \sigma(x) = \frac{1}{1 + e^{-x}} \tag{2}$$

**Question 2.**

$$h(x) = \tanh(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}} \tag{3}$$

**Question 3.**

$$h(x) = \max(0, x) \tag{4}$$

## 1.3 Multi-layer perceptron

Suppose you are given $n$ data points for a regression or classification task in the form of two matrices: $X \in \mathbb{R}^{n \times d}$ and $Y \in \mathbb{R}^{n \times q}$. Note that this time the data and target vectors are row vectors. This is the standard notation in ML.

In what follows, we will derive the back-propagation update rules for an artificial neural network composed of 1 hidden layer with $m$ hidden neurons and a component wise sigmoid activation function $\sigma$. We denote the input weight matrix by $W^{(1)} \in \mathbb{R}^{d \times m}$, and the output weight matrix by $W^{(2)} \in \mathbb{R}^{m \times q}$.

**Question 0.** Draw the network. Specify the meaning of the edge and nodes in terms of the parameters, inputs and outputs of the network.

**Question 1.** Express the output of the network $\hat{Y} \in \mathbb{R}^{n \times q}$ in terms of the network parameters and activation function.

We train the NN to minimize the MSE loss wrt both $W^{(1)}$ and $W^{(2)}$. The loss can be written:

$$
\begin{aligned}
E(\hat{Y}) &= \frac{1}{2n} ||\hat{Y} - Y||_F^2 \\
&= \frac{1}{n} \sum_{k=1}^{n} \sum_{j=1}^{q} \frac{1}{2} (\hat{y}_{k,j} - y_{k,j})^2.
\end{aligned}
\tag{5}
$$

**Question 2.** (Back-propagation update rules.)

**Question 2a.** Calculate the Jacobian matrix of $E(\hat{Y})$ wrt $W^{(2)} \in \mathbb{R}^{m \times q}$:

$$\frac{\partial E(\hat{Y})}{\partial W^{(2)}}. \tag{6}$$

**Hint:** You can first derive the value of every index, and then find the Jacobian in matrix form.

**Question 2b.** What is the derivative of the composition of scalar functions: $l \circ f \circ h \circ g$ ?

**Question 2c.** Calculate the Jacobian matrix of $E(\hat{Y})$ wrt $W^{(1)} \in \mathbb{R}^{d \times m}$:

$$\frac{\partial E(\hat{Y})}{\partial W^{(1)}}. \tag{7}$$

**Question 2d.** Write the update rule for $W^{(1)}$ and $W^{(2)}$, assuming that the network is trained using batch gradient descent and with learning rate $\eta > 0$. Suppose that we are computing the update rule for step $k + 1$, i.e. we can denote $\hat{Y}_k$, the output of the network when the weights are $W_k^{(1)}$ and $W_k^{(2)}$.

**Question 3.** Implement backprop for this 1 hidden layer neural network example !