

## Exercise Sheet 3

### Exercise 1: Neural Network Optimization (20 + 15 + 15 P)

Consider the one-layer neural network

$$f(\mathbf{x}) = \mathbf{w}^\top \mathbf{x}$$

applied to data points  $\mathbf{x} \in \mathbb{R}^d$ , and where  $\mathbf{w} \in \mathbb{R}^d$  is the parameter of the model. We would like to optimize the mean square error objective:

$$J(\mathbf{w}) = \mathbb{E}_{\hat{p}} \left[ \frac{1}{2} (\mathbf{w}^\top \mathbf{x} - t)^2 \right],$$

where the expectation is computed over an empirical approximation  $\hat{p}$  of the true joint distribution  $p(\mathbf{x}, t)$ . The ground truth is known to be of type:  $t|\mathbf{x} = \mathbf{v}^\top \mathbf{x} + \varepsilon$ , with the parameter  $\mathbf{v}$  unknown, and where  $\varepsilon$  is some small i.i.d. Gaussian noise. The input data follows the distribution  $\mathbf{x} \sim \mathcal{N}(\boldsymbol{\mu}, \sigma^2 I)$  where  $\boldsymbol{\mu}$  and  $\sigma^2$  are the mean and variance.

(a) *Compute* the Hessian of the objective function  $J$  at the current location  $\mathbf{w}$  in the parameter space, and as a function of the parameters  $\boldsymbol{\mu}$  and  $\sigma$  of the data.

(b) *Show* that the condition number of the Hessian is given by:  $\frac{\lambda_1}{\lambda_d} = 1 + \frac{\|\boldsymbol{\mu}\|^2}{\sigma^2}$ .

(c) *Explain* for this particular problem what would be the advantages and disadvantages of centering the data before training. Your answer could include the following aspects: (1) condition number and speed of convergence, (2) ability to reach a low prediction error.

### Exercise 2: Programming (50 P)

Download the programming files on ISIS and follow the instructions.