

## Exercise Sheet 12

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

Consider the one-layer neural network

$$y = \mathbf{w}^\top \mathbf{x} + b$$

applied to data points  $\mathbf{x} \in \mathbb{R}^d$ , and where  $\mathbf{w} \in \mathbb{R}^d$  and  $b \in \mathbb{R}$  are the parameters of the model. We consider the optimization of the objective:

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

where the expectation is computed over an empirical approximation  $\hat{p}$  of the true joint distribution  $p(\mathbf{x}, t)$  and  $t \in \{-1, 1\}$ . 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:  $\frac{\lambda_1}{\lambda_d} = 1 + \frac{\|\boldsymbol{\mu}\|^2}{\sigma^2}$ .

### Exercise 2: Neural Network Regularization (10 + 10 + 10 P)

For a neural network to generalize from limited data, it is desirable to make it sufficiently invariant to small local variations. This can be done by limiting the gradient norm  $\|\partial f / \partial \mathbf{x}\|$  for all  $\mathbf{x}$  in the input domain. As the input domain can be high-dimensional, it is impractical to minimize the gradient norm directly. Instead, we can minimize an upper-bound of it that depends only on the model parameters.

We consider a two-layer neural network with  $d$  input neurons,  $h$  hidden neurons, and one output neuron. Let  $W$  be a weight matrix of size  $d \times h$ , and  $(b_j)_{j=1}^h$  a collection of biases. We denote by  $W_{i,:}$  the  $i$ th row of the weight matrix and by  $W_{:,j}$  its  $j$ th column. The neural network computes:

$$\begin{aligned} a_j &= \max(0, W_{:,j}^\top \mathbf{x} + b_j) && \text{(layer 1)} \\ f(\mathbf{x}) &= \sum_j s_j a_j && \text{(layer 2)} \end{aligned}$$

where  $s_j \in \{-1, 1\}$  are fixed parameters. The first layer detects patterns of the input data, and the second layer computes a fixed linear combination of these detected patterns.

- (a) *Show* that the gradient norm of the network can be upper-bounded as:

$$\left\| \frac{\partial f}{\partial \mathbf{x}} \right\| \leq \sqrt{h} \cdot \|W\|_F$$

- (b) Let  $\|W\|_{\text{Mix}} = \sqrt{\sum_i \|W_{i,:}\|_1^2}$  be a  $\ell_1/\ell_2$  mixed matrix norm. *Show* that the gradient norm of the network can be upper-bounded by it as:

$$\left\| \frac{\partial f}{\partial \mathbf{x}} \right\| \leq \|W\|_{\text{Mix}}$$

- (c) *Show* that the mixed norm provides a bound that is tighter than the one based on the Frobenius norm, i.e. show that:

$$\|W\|_{\text{Mix}} \leq \sqrt{h} \cdot \|W\|_F$$

.

### Exercise 3: Programming (40 P)

Download the programming files on ISIS and follow the instructions.

## ML1 Sheet 12

1) a) The Hessian of  $J$  can be obtained by taking the second ~~derivative~~ partial derivative:

$$\begin{aligned} H &= \frac{\partial}{\partial w} \left( \frac{\partial}{\partial w} J(w) \right) = \frac{\partial}{\partial w} \left( \frac{\partial}{\partial w} E_p \left[ \frac{1}{2} (1 - y + t)^2 \right] \right) \\ &= \frac{\partial}{\partial w} \left( \frac{\partial}{\partial w} E_p \left[ \frac{1}{2} (1 - (w^T x + b) + t)^2 \right] \right) \\ &= \frac{\partial}{\partial w w^T} E_p \left[ \frac{1}{2} (1^2 - 2(w^T x + b)t + (w^T x + b)^2 t^2) \right] \\ &= \frac{\partial}{\partial w w^T} E_p \left[ \frac{1}{2} - w^T x t + b t + t^2 (w^T x x^T w + 2w^T x b + b^2) \right] \\ &= \frac{\partial}{\partial w w^T} E_p \left[ \frac{1}{2} - w^T x t + b t + t^2 w^T x x^T w + t^2 w^T x b + t^2 b^2 \frac{1}{2} \right] \\ &= E_p \left[ 0 - 0 + 0 + t^2 x x^T + 0 + 0 \right] \\ &= \underbrace{t^2}_{=1} E_p [x x^T] = E_p [x x^T] \\ &= \text{cov}(x) + E[x] E[x]^T = \sigma^2 \mathbb{I} + \vec{\mu} \vec{\mu}^T \end{aligned}$$

b) Show  $\frac{\lambda_1}{\lambda_d} = 1 + \frac{\|\mu\|^2}{\sigma^2}$

$$\begin{aligned} \lambda_1 &= \max_{\|v\|=1} v^T H v = \max_{\|v\|=1} v^T (\sigma^2 \mathbb{I} + \vec{\mu} \vec{\mu}^T) v = \max_{\|v\|=1} \sigma^2 + (v^T \vec{\mu})^2 \\ &= \sigma^2 + \left\| \frac{\vec{\mu}^T}{\|\vec{\mu}\|} \vec{\mu} \right\|^2 = \sigma^2 + \|\vec{\mu}\|^2 \end{aligned}$$

$$\lambda_2 = \max_{\substack{\|v\|=1 \\ v \perp \mu}} v^T H v = \max_{\substack{\|v\|=1 \\ v^T \mu = 0}} v^T (\sigma^2 \mathbb{I} + \mu \mu^T) v = \sigma^2 = \lambda_3 = \dots = \lambda_d$$

$$\rightarrow \frac{\lambda_1}{\lambda_d} = \frac{\sigma^2 + \|\mu\|^2}{\sigma^2} = 1 + \frac{\|\mu\|^2}{\sigma^2}$$



2

a)

$$\begin{aligned}\left\| \frac{\partial f}{\partial x} \right\|^2 &= \sum_{i=1}^d \left( \frac{\partial f}{\partial x_i} \right)^2 = \sum_{i=1}^d \left( \sum_{j=1}^h s_j \cdot 1_{q_j > 0} \cdot w_{ij} \right)^2 \\ &= \sum_{i=1}^d \left( \sum_{j=1}^h (s_j \cdot 1_{q_j > 0})^2 \sum_{j=1}^h w_{ij}^2 \right) \\ &\leq \sum_{i=1}^d h \sum_{j=1}^h w_{ij}^2 = h \|W\|_F^2\end{aligned}$$

$$\Rightarrow \left\| \frac{\partial f}{\partial x} \right\| \leq \sqrt{h} \|W\|_F$$

$$\begin{aligned}b) \left\| \frac{\partial f}{\partial x} \right\|^2 &= \sum_{i=1}^d \left( \frac{\partial f}{\partial x_i} \right)^2 = \sum_{i=1}^d \left( \sum_{j=1}^h |s_j| 1_{q_j > 0} |w_{ij}| \right)^2 \\ &\leq \sum_{i=1}^d \left( \sum_{j=1}^h |w_{ij}| \right)^2 = \sum_{i=1}^d \|w_{i:}\|_1^2\end{aligned}$$

$$\Rightarrow \left\| \frac{\partial f}{\partial x} \right\| \leq \|W\|_{\max}$$

$$\begin{aligned}c) \|W\|_{\text{Mix}}^2 &= \sum_{i=1}^d \left( \sum_{j=1}^h 1_j \right)^2 \leq \sum_{i=1}^d \underbrace{\left( \sum_{j=1}^h 1_j^2 \right)}_h \left( \sum_{j=1}^h w_{ij}^2 \right) \\ &= h \|W\|_F^2\end{aligned}$$

Cauchy-Schwarz:  
 $\langle a, b \rangle^2 \leq \|a\|^2 \|b\|^2$