

## Exercise Sheet 3

### Exercise 1: Neural Network Optimization (20 + 20 + 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: Initialization (35 + 10 P)

Consider a deep neural network with  $L$  layers with width  $n$  and a ReLU activation function. Assume the dataset  $X$ , which consists of samples  $x \in \mathbb{R}^n$  which are iid..  $X$  is centered and whitened, i.e.,  $\mathbb{E}[x^{(i)}] = 0$  and  $\text{Var}[x^{(i)}] = 1 \forall i \in \{1, \dots, n\}$ ,  $\text{Cov}(x^{(i)}, x^{(j)}) = 0 \forall i \neq j$  where  $i, j$  indicate the dimensions.

The He-initialization is defined as follows:

$$W_{ij}^{(l)} \sim \mathcal{N}\left(0, \frac{2}{n}\right)$$

$$b_i^{(l)} = 0,$$

where  $W_{(l)}$  is the weight matrix of layer  $l$  and  $b^{(l)}$  is the bias vector of layer  $l$ .

You may use the following assumptions/hints:

- For a random variable  $Y$  centered around 0, i.e.  $\mathbb{E}[Y] = 0$ , we assume  $\mathbb{E}[\text{ReLU}(Y)^2] = \frac{1}{2} \text{Var}(Y)$ .
- For mutually independent random variables  $a, b$ , we have  $\mathbb{E}[ab] = \mathbb{E}[a]\mathbb{E}[b]$ .
- $\mathbb{E}[\sum_i Y^{(i)}] = \sum_i \mathbb{E}[Y^{(i)}] = n\mathbb{E}[Y]$ .
- $a_0$  is the input to the neural network.

(a) *Show* by induction that, when using the initialization scheme of He et al. (2015), the variance of the latent variables  $z_{(l)}^{(i)} = \sum_j W_{(l)}^{(ij)} a_{(l-1)}^{(j)} + b_{(l)}^{(i)}$  for all layers  $l \in \{2, \dots, L\}$  stays constant, i.e.  $\text{Var}(\sum_j W_{(l+1)}^{(ij)} a_{(l)}^{(j)} + b_{(l+1)}^{(i)}) = \text{Var}(\sum_j W_{(l)}^{(ij)} a_{(l-1)}^{(j)} + b_{(l)}^{(i)})$ .

(b) How do you need to choose the initialization for tanh if you assume the following?

- Around 0, we have  $\tanh(x) \approx x$ .

*Show your work.*

## Exercise 1: Neural Network Optimization (20 + 20 + 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  $\mathbf{x} \sim \mathcal{N}(\boldsymbol{\mu}, \sigma^2 \mathbf{I})$  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.

$$\begin{aligned} \frac{\partial J}{\partial w_i} &= \mathbb{E}_{\hat{p}} \left[ \frac{1}{2} (\mathbf{w}^\top \mathbf{x} - t) \cdot x_i \right] & H &= \frac{\partial}{\partial \mathbf{w} \mathbf{w}^\top} \mathbb{E} \left[ \frac{1}{2} (\mathbf{w}^\top \mathbf{x} - t)^2 \right] = \frac{1}{2} \mathbb{E} [\mathbf{w}^\top \mathbf{x} (\mathbf{x}^\top \mathbf{w})] \\ \frac{\partial J}{\partial w_j \partial w_i} &= \mathbb{E} [x_j x_i] & &= \frac{\partial}{\partial \mathbf{w} \mathbf{w}^\top} \mathbb{E} \left[ \frac{1}{2} (\mathbf{w}^\top \mathbf{x}) (\mathbf{w}^\top \mathbf{x}) \right] + \text{lin.} + \text{const.} \\ & & &= \frac{\partial}{\partial \mathbf{w} \mathbf{w}^\top} \mathbb{E} \left[ \frac{1}{2} (\mathbf{w}^\top \mathbf{x})^2 \right] = \mathbb{E} [\mathbf{x} \mathbf{x}^\top] \\ & & &= \text{Cov}(\mathbf{x}, \mathbf{x}^\top) + \mathbb{E}[\mathbf{x}] \mathbb{E}[\mathbf{x}^\top] \\ & & &= \sigma^2 \mathbf{I} + \boldsymbol{\mu} \boldsymbol{\mu}^\top \end{aligned}$$

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

Assume  $\vec{v}$  is the eigenvector of  $\lambda$

$$H \vec{v} = \lambda \vec{v}$$

$$\mathbf{v}^\top H \vec{v} = \lambda \mathbf{v}^\top \vec{v} = \lambda (v_1 \ v_2 \ \dots \ v_d) \begin{pmatrix} v_1 \\ v_2 \\ \vdots \\ v_d \end{pmatrix} = \lambda \sum_{i=1}^d v_i^2 = \lambda \|\vec{v}\|^2$$

$$\begin{aligned} \lambda_1 = \lambda_{\max} &= \frac{\mathbf{v}^\top H \vec{v}}{\|\vec{v}\|^2} = \frac{1}{\|\vec{v}\|^2} \mathbf{v}^\top (\sigma^2 \mathbf{I} + \boldsymbol{\mu} \boldsymbol{\mu}^\top) \mathbf{v} \\ &= \frac{\sigma^2}{\|\vec{v}\|^2} \mathbf{v}^\top \mathbf{v} + \frac{1}{\|\vec{v}\|^2} \mathbf{v}^\top \boldsymbol{\mu} \boldsymbol{\mu}^\top \mathbf{v} \\ &= \sigma^2 + \frac{1}{\|\vec{v}\|^2} \mathbf{v}^\top \boldsymbol{\mu} (\mathbf{v}^\top \boldsymbol{\mu})^\top \\ &= \sigma^2 + \frac{1}{\|\vec{v}\|^2} \cdot \|\mathbf{v}^\top \boldsymbol{\mu}\|^2 \\ &= \sigma^2 + \frac{\sum_i (v_i \mu_i)^2}{\sum_i (v_i)^2} \\ &= \sigma^2 + \sum_i \mu_i^2 \\ &= \sigma^2 + \|\boldsymbol{\mu}\|^2 \end{aligned}$$

$$\lambda_2 = \lambda_{\min} = \sigma^2 \quad (\text{if } \|\boldsymbol{\mu}\| = 0)$$

$$\lambda_1 = \lambda_{\max} = \max_{\|\mathbf{v}\|=1} \mathbf{v}^\top H \mathbf{v}$$

$$\begin{aligned} &= \max_{\|\mathbf{v}\|=1} \mathbf{v}^\top (\sigma^2 \mathbf{I} + \boldsymbol{\mu} \boldsymbol{\mu}^\top) \mathbf{v} \\ &= \max_{\|\mathbf{v}\|=1} \mathbf{v}^\top \mathbf{v} \cdot \sigma^2 + \mathbf{v}^\top \boldsymbol{\mu} \boldsymbol{\mu}^\top \mathbf{v} \\ &= \max_{\|\mathbf{v}\|=1} \|\mathbf{v}\|^2 \sigma^2 + (\boldsymbol{\mu}^\top \mathbf{v})^\top \boldsymbol{\mu}^\top \mathbf{v} \\ &= \max_{\|\mathbf{v}\|=1} \sigma^2 + \|\boldsymbol{\mu}^\top \mathbf{v}\|^2 \end{aligned}$$

$$= \max_{\|\mathbf{v}\|=1} \sigma^2 + \sum_i (\mu_i v_i)^2$$

$$= \sigma^2 + \sum_i \mu_i^2 \quad (\forall v_i, \mu_i = v_i, \text{ i.e. } \boldsymbol{\mu} = \mathbf{v})$$

$$= \sigma^2 + \|\boldsymbol{\mu}\|^2$$

$$\lambda_d = \lambda_{\min}$$

$$= \min_{\|\mathbf{v}\|=1} \sigma^2 + \sum_i (\mu_i v_i)^2$$

$$= \sigma^2 + 0 = \sigma^2$$

RLC Sheet 9

$$\frac{\lambda_1}{\lambda_2} = \frac{\delta^2 + \|M\|^2}{\delta^2} = 1 + \frac{\|M\|^2}{\delta^2}$$

(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.

Pro: Centering makes  $\|M\| = 0$  so that  $\frac{\lambda_1}{\lambda_2} = 1 + \frac{0}{\delta^2} = 1$ . The condition number will be lower and therefore it converge faster. Centering makes condition number  $\frac{\lambda_1}{\lambda_2}$  lower:  $1 + \frac{0}{\delta^2} < 1 + \frac{\|M\|^2}{\delta^2}$  therefore, convergence is faster.

Con: Centering may lead to a higher error cause the error function is influenced by the means and

covariance of the data  $(\xi(w)) = w^T \mathbb{E}[(x - \mu)(x - \mu)^T + \mu \mu^T + \lambda I] w + \text{linear} + \text{const}$

Model performance will likely decrease due to wrong functional form

Model based on centered data  $f(x) = w^T (x - \mathbb{E}[x])$  doesn't contain the ground truth  $f(x) = v^T x$

## Exercise 2: Initialization (35 + 10 P)

Consider a deep neural network with  $L$  layers with width  $n$  and a ReLU activation function. Assume the dataset  $X$ , which consists of samples  $x \in \mathbb{R}^n$  which are iid.  $X$  is centered and whitened, i.e.,  $\mathbb{E}[x^{(i)}] = 0$  and  $\text{Var}[x^{(i)}] = 1 \forall i \in \{1, \dots, n\}$ ,  $\text{Cov}(x^{(i)}, x^{(j)}) = 0 \forall i \neq j$  where  $i, j$  indicate the dimensions. The He-initialization is defined as follows:

$$W_{ij}^{(l)} \sim \mathcal{N}\left(0, \frac{2}{n}\right)$$

$$b_i^{(l)} = 0,$$



where  $W_{(l)}$  is the weight matrix of layer  $l$  and  $b^{(l)}$  is the bias vector of layer  $l$ . You may use the following assumptions/hints:

- For a random variable  $Y$  centered around 0, i.e.  $\mathbb{E}[Y] = 0$ , we assume  $\mathbb{E}[\text{ReLU}(Y)^2] = \frac{1}{2} \text{Var}(Y)$ .
- For mutually independent random variables  $a, b$ , we have  $\mathbb{E}[ab] = \mathbb{E}[a]\mathbb{E}[b]$ .
- $\mathbb{E}[\sum_i Y^{(i)}] = \sum_i \mathbb{E}[Y^{(i)}] = n\mathbb{E}[Y]$ .
- $a_0$  is the input to the neural network.

(a) Show by induction that, when using the initialization scheme of He et al. (2015), the variance of the latent variables  $z_{(l)}^{(i)} = \sum_j W_{(l)}^{(ij)} a_{(l-1)}^{(j)} + b_{(l)}^{(i)}$  for all layers  $l \in \{2, \dots, L\}$  stays constant, i.e.  $\text{Var}(\sum_j W_{(l+1)}^{(ij)} a_{(l)}^{(j)} + b_{(l+1)}^{(i)}) = \text{Var}(\sum_j W_{(l)}^{(ij)} a_{(l-1)}^{(j)} + b_{(l)}^{(i)})$ . 见 solution 需初值  $l=1$  时内容

$$\mathbb{V}[z_i] = \mathbb{V}\left[\sum_j W_{ij}^{(l)} a_{j, l-1}^{(j)} + b_i^{(l)}\right] = \mathbb{V}\left[\sum_j W_{ij}^{(l)} a_{j, l-1}^{(j)}\right]$$

$$= \mathbb{E}\left[\left(\sum_j W_{ij}^{(l)} a_{j, l-1}^{(j)}\right)^2\right] - \mathbb{E}\left[\sum_j W_{ij}^{(l)} a_{j, l-1}^{(j)}\right]^2$$

$$\mathbb{E}\left[\sum_j W_{ij}^{(l)} a_{j, l-1}^{(j)}\right]^2 = \left(n \underbrace{\mathbb{E}[W_{ij}^{(l)}]}_0 \cdot \mathbb{E}[a_{j, l-1}^{(j)}]\right)^2 = 0$$

$$\mathbb{E}\left[\left(\sum_j W_{ij}^{(l)} a_{j, l-1}^{(j)}\right)^2\right] = \mathbb{E}\left[\sum_j \left(W_{ij}^{(l)} a_{j, l-1}^{(j)}\right)^2\right] + \sum_{j \neq m} W_{ij}^{(l)} a_{j, l-1}^{(j)} W_{im}^{(l)} a_{m, l-1}^{(m)}$$

$$\begin{aligned}
&= \mathbb{E} \left[ \sum_j (w_{\ell}^{ij})^2 (a_{\ell-1}^j)^2 \right] + 0 \quad \left( \begin{array}{l} w_{\ell}^{ij}, a_{\ell-1}^j \text{ are independent} \\ \text{so } \sum_{j \neq m} \mathbb{E}[w_{\ell}^{ij} a_{\ell-1}^j w_{\ell}^{im} a_{\ell-1}^m] = 0 \end{array} \right) \\
&= \sum_j \mathbb{E}[(w_{\ell}^{ij})^2] \cdot \mathbb{E}[(a_{\ell-1}^j)^2] \\
&= \sum_j \left( \mathbb{V}[(w_{\ell}^{ij})] + \underbrace{\mathbb{E}^2[w_{\ell}^{ij}]}_0 \right) \cdot \mathbb{E}[(a_{\ell-1}^j)^2] \\
&= \sum_j \left( \frac{2}{n} + 0 \right) \cdot \mathbb{E}[(a_{\ell-1}^j)^2] \\
&= \sum_j \frac{2}{n} \cdot \mathbb{E}[\text{ReLU}(z_{\ell-1}^j)^2] \\
&= \sum_j \frac{2}{n} \cdot \frac{1}{2} \mathbb{V}[z_{\ell-1}^j] \\
&= n \cdot \frac{1}{n} \mathbb{V}[z_{\ell-1}^j] \\
&= \mathbb{V}[z_{\ell-1}^j]
\end{aligned}$$

$$\begin{aligned}
\mathbb{V}[z_{\ell}^i] &= \mathbb{V} \left[ \sum_j w_{\ell}^{ij} a_{\ell-1}^j + b_{\ell}^i \right] = \mathbb{E} \left[ \left( \sum_j w_{\ell}^{ij} a_{\ell-1}^j \right)^2 \right] - \mathbb{E} \left[ \sum_j w_{\ell}^{ij} a_{\ell-1}^j \right]^2 \\
&= \mathbb{V}[z_{\ell-1}^j]
\end{aligned}$$

(b) Now assume instead of ReLU you choose tanh as an activation function. How do you need to choose network parameter initialization if you want to achieve the result from (a) on constant variance of the latent variables? Hint: Around 0, we have  $\tanh(x) \approx x$ . Use this to approximate an expectation value. Show your work.

$$\tanh(z_{\ell}^i) = a_{\ell}^i$$

$$\text{To satisfy: } \mathbb{V}[z_{\ell}^i] = \mathbb{V}[a_{\ell}^i]$$

$$\begin{aligned}
\mathbb{V}[z_{\ell}^i] &= \mathbb{V} \left[ \sum_j w_{\ell}^{ij} a_{\ell-1}^j \right] = \sum_j \left( \mathbb{E}[(w_{\ell}^{ij} a_{\ell-1}^j)^2] - \underbrace{\mathbb{E}[w_{\ell}^{ij} a_{\ell-1}^j]^2}_0 \right) \\
&= \sum_j \left( \mathbb{E}[(w_{\ell}^{ij} a_{\ell-1}^j)^2] \right) \\
&= \mathbb{E} \left[ \sum_j \left( (w_{\ell}^{ij} a_{\ell-1}^j)^2 \right) + \sum_{j \neq m} w_{\ell}^{ij} a_{\ell-1}^j w_{\ell}^{im} a_{\ell-1}^m \right] \\
&= \mathbb{E} \left[ \sum_j (w_{\ell}^{ij})^2 (a_{\ell-1}^j)^2 \right] + 0 \\
&= \sum_j \mathbb{E}[(w_{\ell}^{ij})^2] \cdot \mathbb{E}[(a_{\ell-1}^j)^2] \\
&= \sum_j \left( \mathbb{V}[(w_{\ell}^{ij})] + \underbrace{\mathbb{E}^2[w_{\ell}^{ij}]}_0 \right) \cdot \mathbb{E}[(a_{\ell-1}^j)^2] \\
&= \sum_j \left( \mathbb{V}[w_{\ell}^{ij}] \right) \cdot \mathbb{E}[(a_{\ell-1}^j)^2] \\
&= \sum_j \left( \mathbb{V}[w_{\ell}^{ij}] \right) \cdot \mathbb{E}[\tanh(z_{\ell-1}^j)^2] \\
&\approx \sum_j \left( \mathbb{V}[w_{\ell}^{ij}] \right) \cdot \mathbb{E}[(z_{\ell-1}^j)^2]
\end{aligned}$$

$$\stackrel{!}{=} V[z_{i-1}^j]$$

$$\Leftrightarrow \sum_j (V[w_{i,j}^{ij}] \cdot E[(z_{i-1}^j)^2]) = V[z_{i-1}^j] = E[(z_{i-1}^j)^2] - E[z_{i-1}^j]^2$$

$$n (V[w_{i,j}^{ij}] \cdot E[(z_{i-1}^j)^2]) = E[(z_{i-1}^j)^2] - E[z_{i-1}^j]^2$$

$$V[w_{i,j}^{ij}] = \frac{1}{n} - \frac{E[z_{i-1}^j]^2}{n E[(z_{i-1}^j)^2]} = 0$$

$$= \frac{1}{n}$$

$$\left( \begin{array}{l} \text{Because of } E[z_j^i] = E[\sum_j w_{i,j}^{ij} a_{i,j}^j]^2 \\ = \left( n \underbrace{E[w_{i,j}^{ij}]}_0 \cdot E[a_{i,j}^j] \right)^2 = 0 \end{array} \right)$$

$$\Rightarrow w_{i,j}^{ij} \sim \mathcal{N}(0, \frac{1}{n})$$

$$b_i^i = 0$$

We make approximation  $E[\tanh(t_{i,j}^j)^2] = E[z_{i,j}^j]^2$

from b we have  $V(z_i) = \sum_j \frac{2}{n} E[\tanh(t_{i,j}^j)^2]$

$$= \sum_j \frac{2}{n} E[z_{i,j}^j]^2$$

$$= \sum_j \frac{2}{n} E[(\sum_j w_{i,j}^{ij} a_{i,j}^j)^2]$$

$$= \sum_j \frac{2}{n} (V[\sum_j w_{i,j}^{ij} a_{i,j}^j] + 0)$$