

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 ε 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 σ^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 σ 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 . Also, the $W_{ij}^{(l)}$ are mutually independent.

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 $z_{(l)}^{(j)} = \sum_i W_{(l)}^{(ij)} a_{(l-1)}^{(i)} + b_{(l)}^{(j)}$ for all layers $l \in \{1, \dots, L\}$ stays constant, i.e. $\text{Var}(\sum_i W_{(l+1)}^{(ij)} a_{(l)}^{(i)} + b_{(l+1)}^{(j)}) = \text{Var}(\sum_i W_{(l)}^{(ij)} a_{(l-1)}^{(i)} + b_{(l)}^{(j)})$.

(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.