

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 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.

$$\begin{aligned} H &= \frac{\partial}{\partial \mathbf{w} \mathbf{w}^\top} \mathbb{E} \left[\frac{1}{2} (\mathbf{w}^\top \mathbf{x} - t)^2 \right] \\ &= \frac{\partial}{\partial \mathbf{w} \mathbf{w}^\top} \mathbb{E} \left[\frac{1}{2} (\mathbf{w}^\top \mathbf{x}) (\mathbf{x}^\top \mathbf{w}) + \text{lin.} + \text{const.} \right] \\ &= \mathbb{E} [\mathbf{x} \mathbf{x}^\top] = \text{Cov}(\mathbf{x}) + \mathbb{E}[\mathbf{x}] \mathbb{E}[\mathbf{x}]^\top = \sigma^2 I + \boldsymbol{\mu} \boldsymbol{\mu}^\top \end{aligned}$$

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

$$\begin{aligned} \lambda_1 &= \max_{\|\mathbf{v}\|=1} \mathbf{v}^\top H \mathbf{v} \\ &= \max_{\|\mathbf{v}\|=1} \mathbf{v}^\top (\sigma^2 I + \boldsymbol{\mu} \boldsymbol{\mu}^\top) \mathbf{v} \\ &= \max_{\|\mathbf{v}\|=1} \sigma^2 + \|\mathbf{v}^\top \boldsymbol{\mu}\|^2 \\ &= \sigma^2 + \left\| \frac{\boldsymbol{\mu}^\top}{\|\boldsymbol{\mu}\|} \boldsymbol{\mu} \right\|^2 \\ &= \sigma^2 + \|\boldsymbol{\mu}\|^2 \\ \lambda_2 &= \max_{\substack{\|\mathbf{v}\|=1 \\ \mathbf{v}^\top \boldsymbol{\mu} = 0}} \mathbf{v}^\top (\sigma^2 I + \boldsymbol{\mu} \boldsymbol{\mu}^\top) \mathbf{v} = \sigma^2 \\ \lambda_3, \dots, \lambda_d &= \sigma^2 \end{aligned}$$

Therefore, $\lambda_1/\lambda_d = (\sigma^2 + \|\boldsymbol{\mu}\|^2)/\sigma^2 = 1 + \|\boldsymbol{\mu}\|^2/\sigma^2$

(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.

Advantage: centering makes λ_1/λ_d lower: $1 + \|\mathbf{0}\|^2/\sigma^2 < 1 + \|\boldsymbol{\mu}\|^2/\sigma^2$, therefore, convergence is faster.
Disadvantage: The set of homogeneous models based on centered data $f(\mathbf{x}) = \mathbf{w}^\top (\mathbf{x} - \mathbb{E}[\mathbf{x}])$ does not contain the ground truth $f(\mathbf{x}) = \mathbf{v}^\top \mathbf{x}$.

Exercise 2: Programming (50 P)

Download the programming files on ISIS and follow the instructions.