



1. **Theory Question:** The last lectures discussed the Laplace approximation for deep neural networks  $f(x, \theta) : \mathbb{X} \times \mathbb{R}^D \rightarrow \mathbb{R}^K$  (for simplicity, we will here assume  $K = 1$ ) with parameters  $\theta$ , trained by regularized empirical risk minimization to find the minimum of the loss (negative log posterior)

$$\theta_* = \arg \min_{\theta \in \mathbb{R}^D} \mathcal{L}(\theta) = \arg \min_{\theta \in \mathbb{R}^D} \sum_{i=1}^N \ell(y_i, f(x_i, \theta)) + r(\theta).$$

and approximating the associated posterior as  $p(\theta|X, \mathbf{y}) \approx \mathcal{N}(\theta|\theta_*, -\Psi^{-1})$ , where  $\Psi = \nabla \nabla^\top \mathcal{L}(\theta_*)$  is the Hessian of the loss at the optimum. An obvious challenge for this formalism is that the Hessian  $\Psi$  is of size  $D \times D$ , and thus inverting it is computationally expensive. In this exercise, we will explore a popular way to address this issue, known as the *Generalized Gauss-Newton (GGN)* approximation.<sup>1</sup>

- (a) Using the Chain rule, show that the Hessian  $\Psi$  can be written as

$$\Psi = \sum_{i=1}^N (J(x_i, \theta) H(y_i, f(x_i, \theta)) J(x_i, \theta)^\top + L(y_i, f(x_i, \theta)) F(x_i, \theta)) + \nabla \nabla^\top r(\theta)$$

where  $H \in \mathbb{R}$ ,  $J \in \mathbb{R}^D$ ,  $F \in \mathbb{R}^{D \times D}$  and  $L \in \mathbb{R}$  and

$$\begin{aligned} [J(x_i, \theta)]_j &= \frac{\partial f(x_i, \theta)}{\partial \theta_j} && \text{is the Jacobian of } f \text{ wrt. the parameters} \\ H(y_i, f(x_i, \theta)) &= \frac{\partial^2 \ell(y_i, f(x_i, \theta))}{\partial f^2} && \text{is the Hessian of } \ell \text{ wrt. the network} \\ [F(x_i, \theta)]_{nm} &= \frac{\partial^2 f(x_i, \theta)}{\partial \theta_n \partial \theta_m} && \text{is the Hessian of } f \\ L(y_i, f(x_i, \theta)) &= \frac{\partial \ell(y_i, f(x_i, \theta))}{\partial f(x_i, \theta)} && \text{is the Jacobian of } \ell \text{ wrt. the network} \end{aligned}$$

Note that  $J$  is the vector that we also used in the lecture to construct the Laplace tangent kernel (for multi-output models, it is a rectangular matrix).

- (b) Consider the two standard choices  $\ell_2(y_i, f(x_i)) = \frac{1}{2}(y_i - f(x_i))^2$  (square loss) and  $\ell_{\text{ce}}(y_i, f(x_i)) = -\log(1 + \exp(-y_i f(x_i)))$  (cross-entropy loss). Compute the Jacobian  $L$  and Hessian  $H$  of the loss wrt. the network output  $f(x_i)$  for both choices of loss. Use the results to show that, in each case, a “well-trained network” (i.e. a choice  $\theta_*$  that minimizes the loss) satisfies  $L(y_i, f(x_i, \theta_*)) \approx 0$ , and  $H \geq 0$ . Thus, the Hessian can be computed efficiently (in  $\mathcal{O}(DN)$ ) using only the first term in  $\Psi$  above, and the resulting approximation is positive semi-definite.<sup>2</sup>
- (c) If we consider  $r(\theta) = \alpha I$ , then the above results suggest that we can approximately write the Hessian as a scalar plus rank- $N$  matrix :

$$\Psi = \alpha I + \sum_{i=1}^N g_i g_i^\top \quad (\text{with } g_i = J(x_i, \theta_*) \sqrt{H(y_i, f(x_i, \theta_*))})$$

<sup>1</sup>Nicol N Schraudolph, *Fast Curvature Matrix-Vector Products for Second-Order Gradient Descent*. Neural computation, **14**(7), 2002.

<sup>2</sup>Specifically for these two choices, the resulting form of  $\Psi$  is also equal to the Fisher information matrix of the model. But this is not important to understand this question. Cf. J. Martens *New insights and perspectives on the natural gradient method*. Journal of Machine Learning Research, **21**(146):1–76, 2020.

Use the matrix inversion Lemma to provide an explicit expression of the inverse  $\Psi^{-1}$  used in the Laplace approximation. What is the computational complexity of computing  $\Psi^{-1}$  in this case, in terms of  $N$  and  $D$ ?

2. **Practical Question:** can be found in `Ex10.ipynb`