

# Homework 6.1

Deep Learning 2024 Spring

Due on 2024/4/29

## 1 Q&A

**Problem 1.** (Noise Contrastive Estimation) Suppose we are using NCE to train a language model. Denote the context as  $\mathbf{h}$ , the target word as  $\mathbf{w}$ , and noise word samples as  $\bar{\mathbf{w}}$ .  $\tilde{p}_{\mathbf{w}|\mathbf{h}}(w|h)$  and  $p_{\mathbf{w}|\mathbf{h}}^\theta(w|h)$  are the target word distribution under context  $h$  of the corpus and the learning model respectively.  $q_{\bar{\mathbf{w}}}(\bar{w})$  is the noise distribution introduced by NCE. Assume the learning distribution is self-normalized, i.e.,

$$p_{\mathbf{w}|\mathbf{h}}^\theta(w|h) = \frac{u^\theta(w, h)}{Z_h^\theta} \approx u^\theta(w, h).$$

If we choose 1 positive sample and  $k$  negative samples, the loss of NCE is given by

$$L_{\text{NCE}}^k(\theta; h) = \sum_w \tilde{p}_{\mathbf{w}|\mathbf{h}}(w|h) \log \left( \frac{u^\theta(w, h)}{u^\theta(w, h) + k q_{\bar{\mathbf{w}}}(w)} \right) + \sum_{1 \leq i \leq k, \bar{w}} q_{\bar{\mathbf{w}}}(\bar{w}) \log \left( \frac{k q_{\bar{\mathbf{w}}}(\bar{w})}{u^\theta(\bar{w}, h) + k q_{\bar{\mathbf{w}}}(\bar{w})} \right)$$

Prove that as  $k \rightarrow \infty$ ,  $\nabla L_{\text{NCE}}^k(\theta; h) \rightarrow \nabla L_{\text{MLE}}(\theta; h)$ , where

$$\nabla L_{\text{MLE}}(\theta; h) = \sum_w \left( \tilde{p}_{\mathbf{w}|\mathbf{h}}(w|h) - p_{\mathbf{w}|\mathbf{h}}^\theta(w|h) \right) \nabla \log u^\theta(w, h).$$