

# Supplementary Material for NeurIPS Rebuttal

## Supplementary Material for NeurIPS Rebuttal

Correpondance between the quantized activation and the quantized objective function

## Correpondance between the quantized activation and the quantized objective function

For the answer of the question 1 and 2, we insist the quantization of an objective function and the activation function is almost equivalent.

First, we consider the following quantization for an objective function  $f : \mathbf{R}^d \rightarrow \mathbf{R}$  with respect to an activation function as follows:

$$f = \sum_{k=0}^{\infty} f_k b^{n-k} = \sum_{k=0}^{m-1} f_k b^{n-k} + \sum_{k=m}^{\infty} f_k b^{n-k} = f^Q + O(b^m), \quad \because f^Q = \sum_{k=0}^m f_k b^{n-k}, \quad (1)$$

where  $b \in \mathbf{Z}^+$  denotes the base of the number system for quantization, and  $f_k \in \mathbf{Z}^+[0, b)$  denotes a coefficient for  $b^{n-k}$ . Hence considering a binary number system such that  $b = 2$ , we note that  $f_i \in \{0, 1\} \forall i \in \mathbf{Z}^+$ .

Assume that there exist a neural network contains  $l$  layers which contain the map  $\mathbf{h}^l : \mathbf{R}^m \rightarrow \mathbf{R}^d$  consisted with the activation function  $h_i^l : \mathbf{R} \rightarrow \mathbf{R}^d$  at the  $i$ -th node in the  $l$ -th Layer such that  $h_i^l(y) \triangleq h_i^l(\mathbf{w}_i \mathbf{h}^{l-1})$ ,  $\mathbf{h} = [h_i^l]_{i=1}^d$  where  $\mathbf{w}_i \in \mathbf{R}^m$ ,  $[w_i^j] \in \mathbf{R}^{d \times m}$  denotes the weight vector for the  $i$  th node.

Additionally, we let a quantized activation function  $\mathbf{h}_{s_q}^{lQ}$  with the quantization step defined as the reciprocal of the quantization parameter  $\mathbf{Q}_p^{-1} \in \mathbf{Q}^d$  such that  $\mathbf{h}_{s_q}^{lQ} = \mathbf{h}_0^{lQ} + s_q \mathbf{Q}_p^{-1}$ ,  $s_q \in \mathbf{Z}$ , where each component of  $\mathbf{Q}_p^{-1}$ , i.e.  $Q_{p,i}^{-1}$  represents one of the elements to the set  $\{-Q_p^{-1}, 0, Q_p^{-1}\}$ .

Consider the second-order Taylor series for the objective function  $f$ . Particularly, we set  $s_q = 1$  for convenience, then

$$\begin{aligned} f(\mathbf{h}_1^{lQ}) &= f(\mathbf{h}_0^{lQ}) + \nabla_{\mathbf{h}} f(\mathbf{h}_0^{lQ}) \cdot \mathbf{Q}_p^{-1} + \frac{1}{2} \mathbf{Q}_p^{-1} \cdot \nabla_{\mathbf{h}}^2 f(\mathbf{h}_0^{lQ}) \cdot \mathbf{Q}_p^{-1} + O(\|\mathbf{Q}_p^{-1}\|^3) \\ &\approx \bar{f}(\mathbf{h}_1^{lQ}) + O(Q_p^{-3}), \end{aligned} \quad (3)$$

where  $\bar{f}(\mathbf{h}_1^{lQ}) \triangleq f(\mathbf{h}_0^{lQ}) + \nabla_{\mathbf{h}} f(\mathbf{h}_0^{lQ}) \cdot \mathbf{Q}_p^{-1} + \frac{1}{2} \mathbf{Q}_p^{-1} \cdot \nabla_{\mathbf{h}}^2 f(\mathbf{h}_0^{lQ}) \cdot \mathbf{Q}_p^{-1}$ .

Assume that the quantization step  $\mathbf{Q}_p^{-1}$  is sufficiently small such that  $[O(\|\mathbf{Q}_p^{-1}\|^3)]^Q = 0$ . Then, we calculate the Taylor expansion of the quantization of  $f$  such that

$$\begin{aligned}
f^Q(\mathbf{h}_1^{l^Q}) &= f^Q(\mathbf{h}_0^{l^Q}) + \left[ \nabla_{\mathbf{h}} f(\mathbf{h}_0^{l^Q}) \right]^Q \cdot \mathbf{Q}_p^{-1} + \frac{1}{2} \mathbf{Q}_p^{-1} \cdot \left[ \nabla_{\mathbf{h}}^2 f(\mathbf{h}_0^{l^Q}) \right]^Q \cdot \mathbf{Q}_p^{-1} + [O(\|\mathbf{Q}_p^{-1}\|^3)]^Q \\
&= f(\mathbf{h}_0^{l^Q}) + \varepsilon_q \mathbf{Q}_p^{-1} + \left( \nabla_{\mathbf{h}} f(\mathbf{h}_0^{l^Q}) + \varepsilon_q \mathbf{Q}_p^{-1} \right) \cdot \mathbf{Q}_p^{-1} + \frac{1}{2} \mathbf{Q}_p^{-1} \cdot \left( \nabla_{\mathbf{h}}^2 f(\mathbf{h}_0^{l^Q}) + \varepsilon_q \mathbf{Q}_p^{-1} \right) \cdot \mathbf{Q}_p^{-1} + [O(\|\mathbf{Q}_p^{-1}\|^3)]^Q \\
&= f(\mathbf{h}_0^{l^Q}) + \nabla_{\mathbf{h}} f(\mathbf{h}_0^{l^Q}) \cdot \mathbf{Q}_p^{-1} + \frac{1}{2} \mathbf{Q}_p^{-1} \cdot \nabla_{\mathbf{h}}^2 f(\mathbf{h}_0^{l^Q}) \cdot \mathbf{Q}_p^{-1} + \varepsilon_q \mathbf{Q}_p^{-1} + O(|\varepsilon_q \cdot \mathbf{Q}_p^{-1}|^2) \\
&= \bar{f}(\mathbf{h}_1^{l^Q}) + \varepsilon_q \mathbf{Q}_p^{-1} + O(|\varepsilon_q \cdot \mathbf{Q}_p^{-1}|^2) \\
&\approx \bar{f}^Q(\mathbf{h}_1^{l^Q}) + O(Q_p^{-2}).
\end{aligned} \tag{4}$$

As shown in (3) and (4), we get

$$|f(\mathbf{h}_1^{l^Q}) - \bar{f}(\mathbf{h}_1^{l^Q})| \approx |f^Q(\mathbf{h}_1^{l^Q}) - \bar{f}^Q(\mathbf{h}_1^{l^Q})| + O(Q_p^{-2}). \tag{2}$$

Consequently, if  $Q_p^{-1}$  is sufficiently small, the objective function calculated from the quantized activation is almost equivalent to the quantization of the objective function.

This result demonstrates that we can develop a learning equation based on the quantized objective function that is equivalent to the learning equation based on the quantized activation.