

RDO-Q: Extremely Fine-Grained Channel-Wise Quantization via Rate-Distortion Optimization —Supplementary Material

Zhe Wang¹, Jie Lin¹, Xue Geng¹, Mohamed M. Sabry Aly², and
Vijay Chandrasekhar^{1,2}

¹ Institute for Infocomm Research, A*STAR, Singapore 138632
wangz@i2r.a-star.edu.sg, jie.dellinger@gmail.com,
geng_xue@i2r.a-star.edu.sg, vijay.cmu@gmail.com

² Nanyang Technological University, 50 Nanyang Ave, Singapore 639798
msabry@ntu.edu.sg

Abstract. In this supplementary materials, we provide the analysis for the additivity property of output distortion. We show a mathematical derivation for the additivity property by linearizing the output distortion using Taylor series expansion.

1 Additivity of Output Distortion

The output distortion δ , caused by quantizing all weight channels and activation layers, equals the sum of all output distortion due to the quantization of each individual weight channel and activation layer

$$\delta = \sum_{i=1}^l \sum_{j=1}^{n_i} \delta_{i,j}^w + \sum_{i=1}^l \delta_i^a \quad (1)$$

if the neural network is continuously differentiable in every layer, the quantization errors can be considered as small deviations distributed with zero mean.

Proof. We first define the main notations. Let

$$\mathcal{F}(W_{1,1}, \dots, W_{1,n_1}, \dots, W_{l,1}, \dots, W_{l,n_l})$$

denote a neural network and

$$\tilde{\mathcal{F}}(W_{1,1}, \dots, W_{1,n_1}, \dots, W_{l,1}, \dots, W_{l,n_l}, s_1, \dots, s_l)$$

denote a modified neural network of \mathcal{F} where an element-wise add layer with parameter s_i is followed for each activation a_i (see Fig. 1). Based on this definition, we have

$$\mathcal{F}(W_{1,1}, \dots, W_{l,n_l}) = \tilde{\mathcal{F}}(W_{1,1}, \dots, W_{l,n_l}, 0, \dots, 0) \quad (2)$$

Define two variables X_0 and ΔX , where

$$X_0 = (W_{1,1}, \dots, W_{l,n_l}, 0, \dots, 0)$$

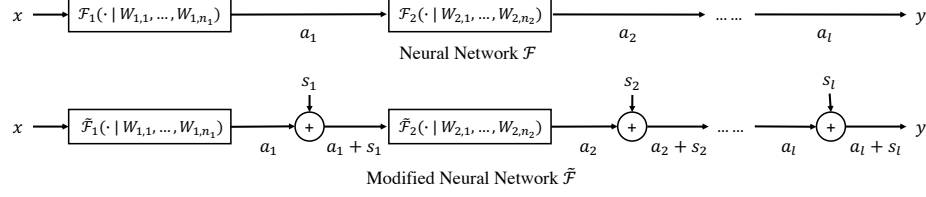


Fig. 1. Examples of a neural network \mathcal{F} and a modified neural network $\tilde{\mathcal{F}}$.

$$\Delta X = (\Delta W_{1,1}, \dots, \Delta W_{l,n_l}, \Delta s_1, \dots, \Delta s_l)$$

Assume that the quantization error can be considered as small deviation. We apply the Taylor series expansion up to first order term on $\tilde{\mathcal{F}}$ at X_0 ,

$$\begin{aligned} \tilde{\mathcal{F}}(X_0 + \Delta X) - \tilde{\mathcal{F}}(X_0) &= \sum_{i,j} \frac{\partial \tilde{\mathcal{F}}}{\partial W_{i,j}} \cdot \Delta W_{i,j} \\ &+ \sum_i \frac{\partial \tilde{\mathcal{F}}}{\partial s_i} \cdot \Delta s_i. \end{aligned} \quad (3)$$

Then $\|\tilde{\mathcal{F}}(X_0 + \Delta X) - \tilde{\mathcal{F}}(X_0)\|^2$ can be written as

$$\begin{aligned} &\left(\sum_{i,j} \Delta W_{i,j}^\top \cdot \frac{\partial \tilde{\mathcal{F}}}{\partial W_{i,j}}^\top + \sum_i \Delta s_i^\top \cdot \frac{\partial \tilde{\mathcal{F}}}{\partial s_i}^\top \right) \\ &\cdot \left(\sum_{i,j} \frac{\partial \tilde{\mathcal{F}}}{\partial W_{i,j}} \cdot \Delta W_{i,j} + \sum_i \frac{\partial \tilde{\mathcal{F}}}{\partial s_i} \cdot \Delta s_i \right) \end{aligned} \quad (4)$$

Because quantization errors in different layers are independently distributed with zero mean, the cross terms of (4) disappear when taking the expectation. That is:

$$\begin{aligned} &E(\Delta W_{i,j}^\top \cdot \frac{\partial \tilde{\mathcal{F}}}{\partial W_{i,j}}^\top \cdot \frac{\partial \tilde{\mathcal{F}}}{\partial W_{i,j}} \cdot \Delta W_{i,j}) \\ &= E(\Delta W_{i,j}^\top) \cdot \frac{\partial \tilde{\mathcal{F}}}{\partial W_{i,j}}^\top \cdot \frac{\partial \tilde{\mathcal{F}}}{\partial W_{i,j}} \cdot E(\Delta W_{i,j}) = 0 \end{aligned} \quad (5)$$

as is the case also for the cross products between $W_{i,j}$ and s_i (all i, j), and s_i and s_j ($i \neq j$). Then, we can obtain

$$\begin{aligned} &E(\|\tilde{\mathcal{F}}(X_0 + \Delta X) - \tilde{\mathcal{F}}(X_0)\|^2) \\ &= \sum_{i,j} E\left(\left\|\frac{\partial \tilde{\mathcal{F}}}{\partial W_{i,j}} \cdot \Delta W_{i,j}\right\|^2\right) + \sum_i E\left(\left\|\frac{\partial \tilde{\mathcal{F}}}{\partial s_i} \cdot \Delta s_i\right\|^2\right) \end{aligned} \quad (6)$$

Eq. (6) is the result we want because, again, according to the Taylor series expansion up to first order term,

$$\begin{aligned} \frac{\partial \tilde{\mathcal{F}}}{\partial W_{i,j}} \cdot \Delta W_{i,j} &= \tilde{\mathcal{F}}(\dots, W_{i,j} + \Delta W_{i,j}, \dots, W_{l,n_l}, 0, \dots) \\ &\quad - \tilde{\mathcal{F}}(\dots, W_{i,j}, \dots, W_{l,n_l}, 0, \dots) \end{aligned} \quad (7)$$

$$\begin{aligned} \frac{\partial \tilde{\mathcal{F}}}{\partial s_i} \cdot \Delta s_i &= \tilde{\mathcal{F}}(W_{1,1}, \dots, W_{l,n_l}, 0, \dots, \Delta s_i, \dots) \\ &\quad - \tilde{\mathcal{F}}(W_{1,1}, \dots, W_{l,n_l}, 0, \dots, 0, \dots) \end{aligned} \quad (8)$$

After dividing both sides of (6) by the dimensionality of the output vector of the neural network, the left side becomes δ and the right side becomes the sum of all output distortion due to the quantization of each individual weight channel and activation layer.