

Supplementary Material to “Momentum Batch Normalization for Deep Learning with Small Batch Size”

Hongwei Yong^{1,2}, Jianqiang Huang², Deyu Meng^{3,4}, Xiansheng Hua², and Lei Zhang^{1,2}

¹ Department of Computing, The Hong Kong Polytechnic University
{cshyong, cslzhang}@comp.polyu.edu.hk

² DAMO Academy, Alibaba Group

{jianqiang.jqh, huaxiansheng}@gmail.com

³ Macau University of Science and Technology

⁴ School of Mathematics and Statistics, Xi'an Jiaotong University
dymeng@mail.xjtu.edu.cn

In this supplementary file, we provide proofs of the theoretical results in the main paper, including Theorem 1 and Theorem 2.

A1. Proof of Theorem 1

Theorem 1 *Suppose samples x_i for $i = 1, 2, \dots, m$ are i.i.d. with $E[x] = \mu$ and $\text{Var}[x] = \sigma^2$, ξ_μ and ξ_σ are defined in Eq.(2), we have:*

$$\lim_{m \rightarrow \infty} p(\xi_\mu) \rightarrow \mathcal{N}(0, \frac{1}{m}), \quad \lim_{m \rightarrow \infty} p(\xi_\sigma) \rightarrow \frac{1}{m} \chi^2(m-1).$$

Proof. From the classical central limit theorem, we have

$$\lim_{m \rightarrow \infty} p\left(\sum_{i=1}^m x_i\right) \rightarrow \mathcal{N}(m\mu, m\sigma^2),$$

that is

$$\lim_{m \rightarrow \infty} p(m\mu_B) \rightarrow \mathcal{N}(m\mu, m\sigma^2).$$

Therefore $\lim_{m \rightarrow \infty} p(\mu_B) \rightarrow \mathcal{N}(\mu, \frac{\sigma^2}{m})$, and $\xi_\mu = \frac{\mu - \mu_B}{\sigma}$ is a linear function of μ_B . Then according to the property of Gaussian distribution, we can obtain that

$$\lim_{m \rightarrow \infty} p(\xi_\mu) \rightarrow \mathcal{N}(0, \frac{1}{m}).$$

For χ^2 distribution, it has the following property:

$$\lim_{m \rightarrow \infty} \frac{\chi^2(m-1)}{m} = \lim_{m \rightarrow \infty} \frac{\chi^2(m-1)}{m-1} \rightarrow \mathcal{N}(1, \frac{2}{m}).$$

And we have $\lim_{m \rightarrow \infty} \mu_B = \mu$. Then for ξ_σ we can also use the central limit theorem to obtain:

$$\lim_{m \rightarrow \infty} p(\xi_\sigma) = \lim_{m \rightarrow \infty} p\left(\frac{1}{m\sigma^2} \sum_{i=1}^m (x_i - \mu_B)^2\right) = \lim_{m \rightarrow \infty} p\left(\sum_{i=1}^m \left(\frac{x_i - \mu}{\sigma}\right)^2\right) \rightarrow \mathcal{N}\left(1, \frac{\kappa}{m}\right),$$

where κ is the kurtosis of x . When m is a very large number, both $\frac{\kappa}{m}$ and $\frac{2}{m}$ are close to zeros so that $\mathcal{N}\left(1, \frac{\kappa}{m}\right) \simeq \mathcal{N}\left(1, \frac{2}{m}\right)$. Therefore, in this case the distribution of ξ_σ can be viewed as $\frac{1}{m}\chi^2(m-1)$.

The proof is completed. ■

A2. Proof of Theorem 2

Theorem 2: *If the infinite derivative of $l(x)$ exists for any x , given two random variables ξ_μ and ξ_σ (> 0), then we have the Taylor expansion for $l\left(\frac{x+\xi_\mu}{\sqrt{\xi_\sigma}}\right)$:*

$$E_{\xi_\mu, \xi_\sigma} \left[l\left(\frac{x+\xi_\mu}{\sqrt{\xi_\sigma}}\right) \right] = l(x) + R^{add}(x) + R^{mul}(x) + R(x), \quad R(x) = \sum_{n=1}^{\infty} \frac{E[\xi_\mu^n]}{n!} \frac{d^n R^{mul}(x)}{dx^n} \quad (1)$$

where $R^{add}(x)$ and $R^{mul}(x)$ are defined in Eq.(5) and (6), respectively.

Proof.

$$\begin{aligned} E_{\xi_\mu, \xi_\sigma} \left[l\left(\frac{x+\xi_\mu}{\sqrt{\xi_\sigma}}\right) \right] &= E_{\xi_\mu, \xi_\sigma} \left[l\left(\frac{x}{\sqrt{\xi_\sigma}} + \frac{\xi_\mu}{\sqrt{\xi_\sigma}}\right) \right] \\ &= E_{\xi_\sigma} \left[l\left(\frac{x}{\sqrt{\xi_\sigma}}\right) \right] + E_{\xi_\mu, \xi_\sigma} \left[\sum_{n=1}^{\infty} \frac{\left(\frac{\xi_\mu}{\sqrt{\xi_\sigma}}\right)^n}{n!} \frac{d^n l\left(\frac{x}{\sqrt{\xi_\sigma}}\right)}{d\left(\frac{x}{\sqrt{\xi_\sigma}}\right)^n} \right] \\ &= l(x) + R^{mul}(x) + E_{\xi_\mu, \xi_\sigma} \left[\sum_{n=1}^{\infty} \frac{\xi_\mu^n}{n!} \frac{d^n l\left(\frac{x}{\sqrt{\xi_\sigma}}\right)}{dx^n} \right] \\ &= l(x) + R^{mul}(x) + E_{\xi_\mu} \left[\sum_{n=1}^{\infty} \frac{\xi_\mu^n}{n!} \frac{d^n E_{\xi_\sigma} \left[l\left(\frac{x}{\sqrt{\xi_\sigma}}\right) \right]}{dx^n} \right] \\ &= l(x) + R^{mul}(x) + E_{\xi_\mu} \left[\sum_{n=1}^{\infty} \frac{\xi_\mu^n}{n!} \frac{d^n (l(x) + R^{mul}(x))}{dx^n} \right] \\ &= l(x) + R^{mul}(x) + E_{\xi_\mu} \left[\sum_{n=1}^{\infty} \frac{\xi_\mu^n}{n!} \frac{d^n (l(x))}{dx^n} \right] + E_{\xi_\mu} \left[\sum_{n=1}^{\infty} \frac{\xi_\mu^n}{n!} \frac{d^n (R^{mul}(x))}{dx^n} \right] \\ &= l(x) + R^{mul}(x) + R^{add}(x) + R(x). \end{aligned} \quad (2)$$

The proof is completed. ■