

Learning Sparse Neural Networks Through Mixture-Distributed Regularization

Chang-Ti Huang¹, Jun-Cheng Chen², Ja-Ling Wu¹

1. National Taiwan University, Taipei, Taiwan
2. Academia Sinica, Taipei, Taiwan

cthuang@cmlab.csie.ntu.edu.tw, pullpull@citi.sinica.edu.tw, wjl@cmlab.csie.ntu.edu.tw

A. The gradient variance of the MDR

In the following subsections, we show both the qualitative and quantitative analysis for the proposed MDR method.

A.1. Qualitative visualization

As shown in the third column of Fig. (3), we can qualitatively verify that the gradient variances of both the power-law function distribution and the exponential-uniform distribution are lower than the gradient variance of the exponential distribution.

A.2. Quantitative analysis

As for the three MDRs compared with the Concrete distribution w.r.t. the gradient variance, it is hard to distinguish the pros and cons from Fig. (3). Hence, we quantitatively validate on the gradient variances of both the MDR and the Concrete distribution using a toy experiment. We generate 10^6 samples drawn from the PDF of the MDRs, i.e.:

$$Q(\zeta) = (1 - q) r(\zeta|z = 0) + qr(\zeta|z = 1), \quad (1)$$

where q is the mixture weight associated with 1-component, setting $q = 0.5$. Besides, we also generate 10^6 samples from the PDF of the Concrete distribution, setting $\log \alpha = 0$. The experiment consists of two measurements, one of which is the gradient variance of drawn samples ζ and the other is the average L_1 distance $\|\zeta - z\|_1$, where

$$z_i = \begin{cases} 0 & \text{if } \zeta_i \leq 0.5 \\ 1 & \text{if } \zeta_i > 0.5 \end{cases}. \quad (2)$$

Each curve in Fig. (4) is plotted with 8 dots, which indicate different temperatures $\beta \in \{40, 34, 28, 22, 16, 10, 7, 4\}$ for the Concrete distribution, $\beta \in \{8, 9, 10, \dots, 15\}$ for the exponential distribution, $\beta \in \{9, 10, 11, \dots, 16\}$ for the exponential-uniform distribution, and $\beta \in \{10, 20, 30, \dots, 80\}$ for the power-law function distribution, respectively.

There is a trade-off between the gradient variance and the estimated binary gate ζ , as shown in Fig. (4). As the gradient variance becomes smaller (i.e. decreasing β for the MDRs), ζ tends to be biased toward the original binary gate z . Hence, the choice of β is still an empirical study. However, from Fig. (4) we observe that for a given gradient variance, the power-law function distribution provides more accurate estimates to the binary gate z .

B. Deriving the reparameterization method for the mixture distribution

Let $r(\zeta|z)$ be the PDF of the mixture component and q be the mixture weight where $q \equiv q(z = 1|\theta)$. Applying the definition of a univariate CDF to the components $z = \{0, 1\}$ yields

$$\text{CDF}^{\text{mix}}(\zeta) = (1 - q) \underbrace{\int_{t=-\infty}^{\zeta} r(t|z = 0) dt}_{\text{0-component}} + q \underbrace{\int_{t=-\infty}^{\zeta} r(t|z = 1) dt}_{\text{1-component}} = u, \quad (3)$$

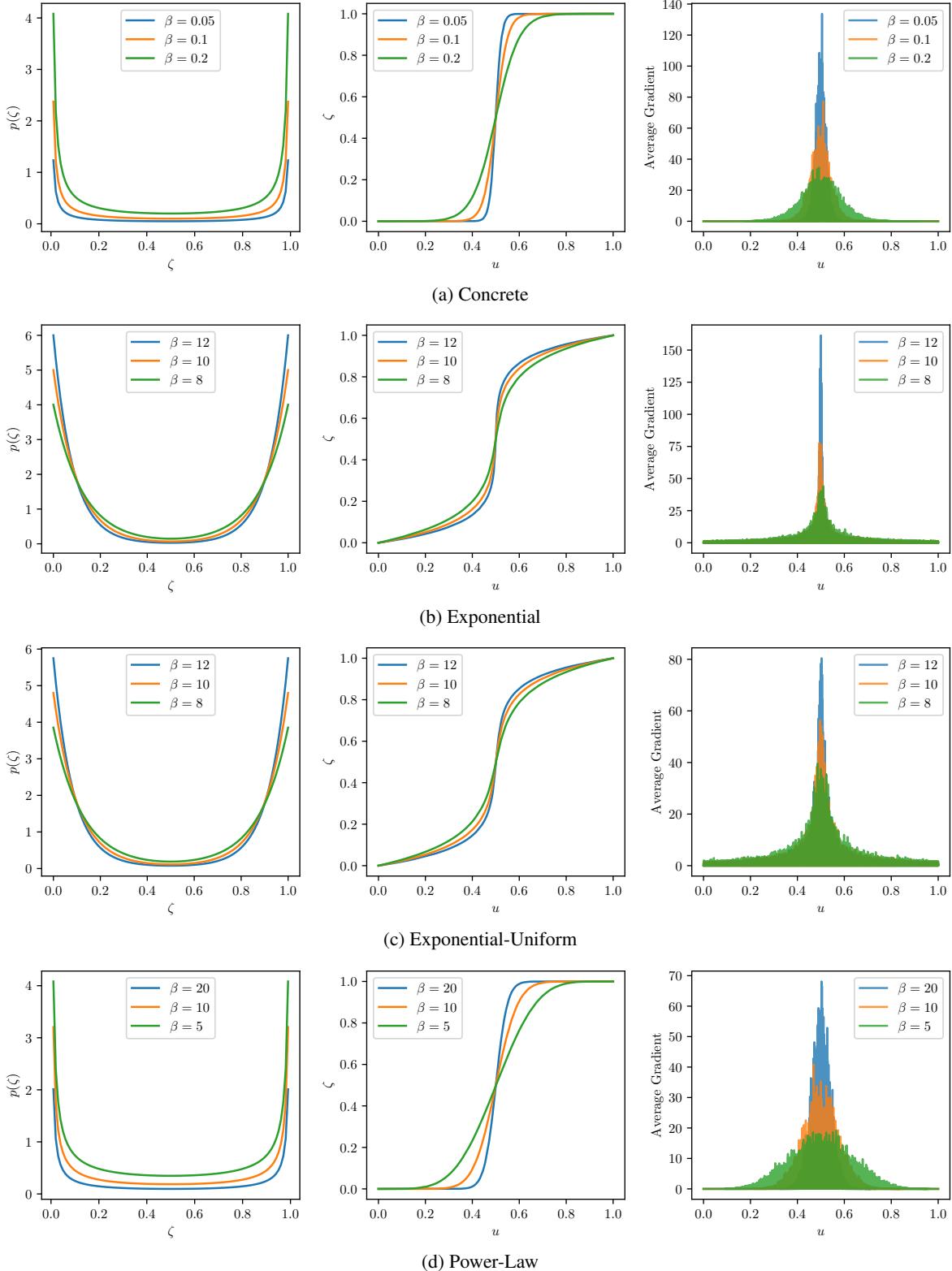


Figure 3: (a) indicating the Concrete distribution and (b)(c)(d) corresponding to different mixture distributions. We visualize the PDF (first column) and the inverse CDF (second column) of the mixture distribution $q(\zeta|\theta) = \sum_{z_i} q(z_i|\theta)r(\zeta|z_i)$ where $q(z=1|\theta) = 0.5$. In the third column, we take the gradient of ζ w.r.t. q from the inverse CDF. Observing the third column, we see that the variance of the gradient increases with increasing β . Further, (c) and (d) generally have lower gradient variance estimates than that of (b).

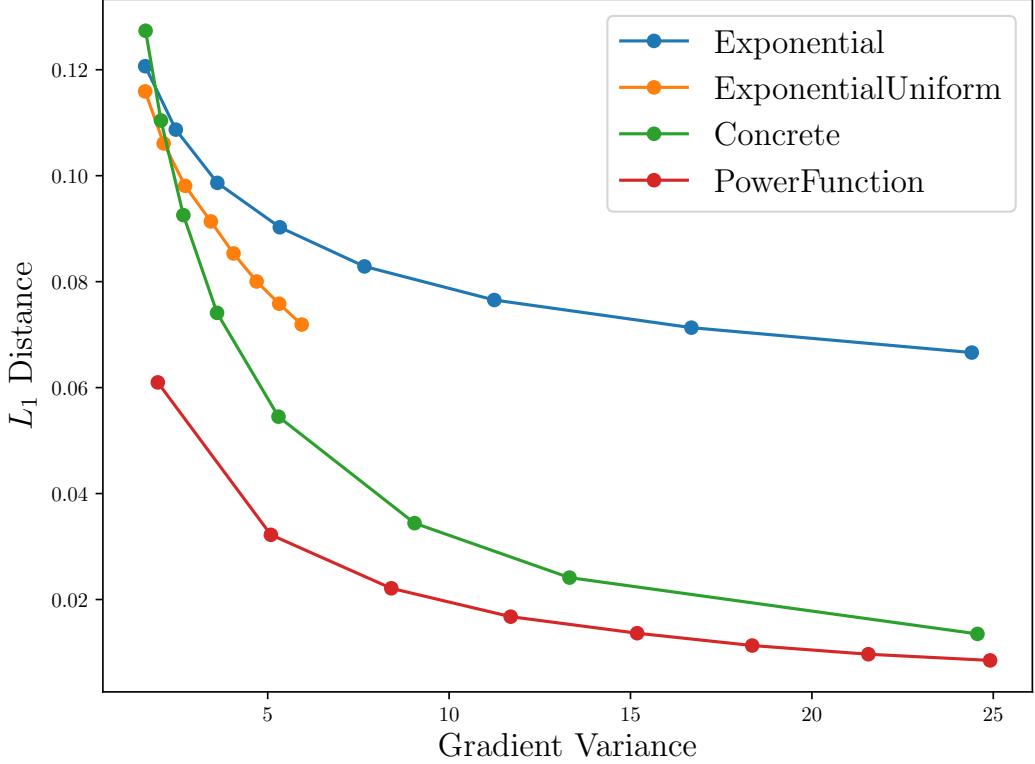


Figure 4: Visualization on the gradient variance (horizontal axis) and the average L_1 distance (vertical axis) of the MDRs and the Concrete. Dots per MDR curve indicate different β values increasing from left to right, whereas dots in the Concrete curve refer to different β values decreasing from left to right.

where $u \sim \text{Uniform}(0, 1)$. Here, we take the gradient w.r.t. q from both sides of Eq. (3) giving

$$\frac{\partial \text{CDF}^{\text{mix}}(\zeta)}{\partial q} = \frac{\partial u}{\partial q} \quad (4)$$

$$\Rightarrow \frac{\partial(1-q)}{\partial q} \int_{t=-\infty}^{\zeta} r(t|z=0) dt + (1-q) \frac{\partial}{\partial q} \left(\int_{t=-\infty}^{\zeta} r(t|z=0) dt \right) \quad (5)$$

$$+ \frac{\partial q}{\partial q} \int_{t=-\infty}^{\zeta} r(t|z=1) dt + q \frac{\partial}{\partial q} \left(\int_{t=-\infty}^{\zeta} r(t|z=1) dt \right) \quad (6)$$

$$= 0, \quad (7)$$

where Eq. (5) and Eq. (6) are given by respectively differentiating the 0- and the 1- components in Eq. (3). Differentiating the mixture weights and computing the integration over the PDFs, we can rewrite the first term in Eq. (5) as

$$\frac{\partial(1-q)}{\partial q} \int_{t=-\infty}^{\zeta} r(t|z=0) dt = -R(\zeta|z=0) \quad (8)$$

and the first term in Eq. (6) as

$$\frac{\partial q}{\partial q} \int_{t=-\infty}^{\zeta} r(t|z=1) dt = R(\zeta|z=1), \quad (9)$$

where R are the CDFs of the mixture components. As for the second term in Eq. (5), we apply the general form of Leibniz integral rule to it:

$$\begin{aligned} \frac{\partial}{\partial q} \left(\int_{t=-\infty}^{\zeta} r(t|z=0) dt \right) &= r(\zeta|z=0) \frac{\partial \zeta}{\partial q} - r(-\infty|z=0) \frac{\partial(-\infty)}{\partial q} \\ &\quad + \int_{t=-\infty}^{\zeta} \frac{\partial}{\partial q} r(t|z=0) dt. \end{aligned} \quad (10)$$

Observing that the derivatives of the second and the third terms on the right-hand side of Eq. (10) equal to zero, we can summarize the above as

$$\frac{\partial}{\partial q} \left(\int_{t=-\infty}^{\zeta} r(t|z=0) dt \right) = r(\zeta|z=0) \frac{\partial \zeta}{\partial q}, \quad (11)$$

and hence the second term of Eq. (6) becomes

$$\frac{\partial}{\partial q} \left(\int_{t=-\infty}^{\zeta} r(t|z=1) dt \right) = r(\zeta|z=1) \frac{\partial \zeta}{\partial q}. \quad (12)$$

Now, we are ready to substituting Eqs. (8), (9), (11) and (12) into Eqs. (5) and (6), giving

$$\frac{\partial \text{CDF}^{\text{mix}}(\zeta)}{\partial q} = -R(\zeta|z=0) + (1-q)r(\zeta|z=0) \frac{\partial \zeta}{\partial q} \quad (13)$$

$$+ R(\zeta|z=1) + q r(\zeta|z=1) \frac{\partial \zeta}{\partial q} = 0. \quad (14)$$

After some rearrangements, we have

$$\frac{\partial \zeta}{\partial q} [(1-q)r(\zeta|z=0) + qr(\zeta|z=1)] + R(\zeta|z=1) - R(\zeta|z=0) = 0, \quad (15)$$

and the gradient w.r.t. q can be obtained as

$$\frac{\partial \zeta}{\partial q} = \frac{R(\zeta|z=0) - R(\zeta|z=1)}{(1-q)r(\zeta|z=0) + qr(\zeta|z=1)}. \quad (16)$$

As shown in Eq. (16), given the PDF and the CDF of the mixture components, we can take gradients of ζ w.r.t. q . Reparameterization for the MDR is highly required when the inverse CDF of the mixture distribution is intractable. In that case, we can not sample ζ through inverse transform sampling, instead we optimize the MDR through backpropagating $\frac{\partial \zeta}{\partial q}$.

With this reparameterization method, we can also take the gradient w.r.t. other model parameters (e.g. β) in a similar way, i.e.:

$$\frac{\partial \text{CDF}^{\text{mix}}(\zeta)}{\partial \beta} = \frac{\partial u}{\partial \beta} \quad (17)$$

$$\Rightarrow \frac{\partial(1-q)}{\partial \beta} \int_{t=-\infty}^{\zeta} r(t|z=0) dt + (1-q) \frac{\partial}{\partial \beta} \left(\int_{t=-\infty}^{\zeta} r(t|z=0) dt \right) \quad (18)$$

$$+ \frac{\partial q}{\partial \beta} \int_{t=-\infty}^{\zeta} r(t|z=1) dt + q \frac{\partial}{\partial \beta} \left(\int_{t=-\infty}^{\zeta} r(t|z=1) dt \right) \quad (19)$$

$$= 0. \quad (20)$$

Observing that $\frac{\partial(1-q)}{\partial\beta} = 0$ and $\frac{\partial q}{\partial\beta} = 0$, we can rewrite the first term in Eq. (18) as

$$\frac{\partial(1-q)}{\partial\beta} \int_{t=-\infty}^{\zeta} r(t|z=0) dt = 0, \quad (21)$$

and the first term in Eq. (6) as

$$\frac{\partial q}{\partial\beta} \int_{t=-\infty}^{\zeta} r(t|z=1) dt = 0. \quad (22)$$

Considering the second term in Eq. (18), we also apply the general form of Leibniz integral rule giving

$$\begin{aligned} \frac{\partial}{\partial\beta} \left(\int_{t=-\infty}^{\zeta} r(t|z=0) dt \right) &= r(\zeta|z=0) \frac{\partial\zeta}{\partial\beta} - r(-\infty|z=0) \frac{\partial(-\infty)}{\partial\beta} \\ &\quad + \int_{t=-\infty}^{\zeta} \frac{\partial}{\partial\beta} r(t|z=0) dt. \end{aligned} \quad (23)$$

It is trivial that $\frac{\partial(-\infty)}{\partial\beta} = 0$. However, the third term on the right-hand side of Eq. (10) equals to zero since $r(t|z=0)$ does not depend on q , whereas the third term on the right-hand side of Eq. (23) contains the dependency between $r(t|z=0)$ and β . To sum up, Eq. (23) can be rewritten as

$$\frac{\partial}{\partial\beta} \left(\int_{t=-\infty}^{\zeta} r(t|z=0) dt \right) = r(\zeta|z=0) \frac{\partial\zeta}{\partial\beta} + \int_{t=-\infty}^{\zeta} \frac{\partial}{\partial\beta} r(t|z=0) dt. \quad (24)$$

The same rule can be applied to the second term of Eq. (19) giving

$$\frac{\partial}{\partial\beta} \left(\int_{t=-\infty}^{\zeta} r(t|z=1) dt \right) = r(\zeta|z=1) \frac{\partial\zeta}{\partial\beta} + \int_{t=-\infty}^{\zeta} \frac{\partial}{\partial\beta} r(t|z=1) dt. \quad (25)$$

Finally, after some rearrangements, we have

$$\frac{\partial\zeta}{\partial\beta} [(1-q)r(\zeta|z=0) + qr(\zeta|z=1)] + \frac{\partial}{\partial\beta} R(\zeta|z=0) + \frac{\partial}{\partial\beta} R(\zeta|z=1) = 0, \quad (26)$$

and the gradient w.r.t. β can now be written as

$$\frac{\partial\zeta}{\partial\beta} = -\frac{(1-q)\frac{\partial}{\partial\beta} R(\zeta|z=0) + q\frac{\partial}{\partial\beta} R(\zeta|z=1)}{(1-q)r(\zeta|z=0) + qr(\zeta|z=1)}. \quad (27)$$

C. Visualization of the emulated binary gates ζ

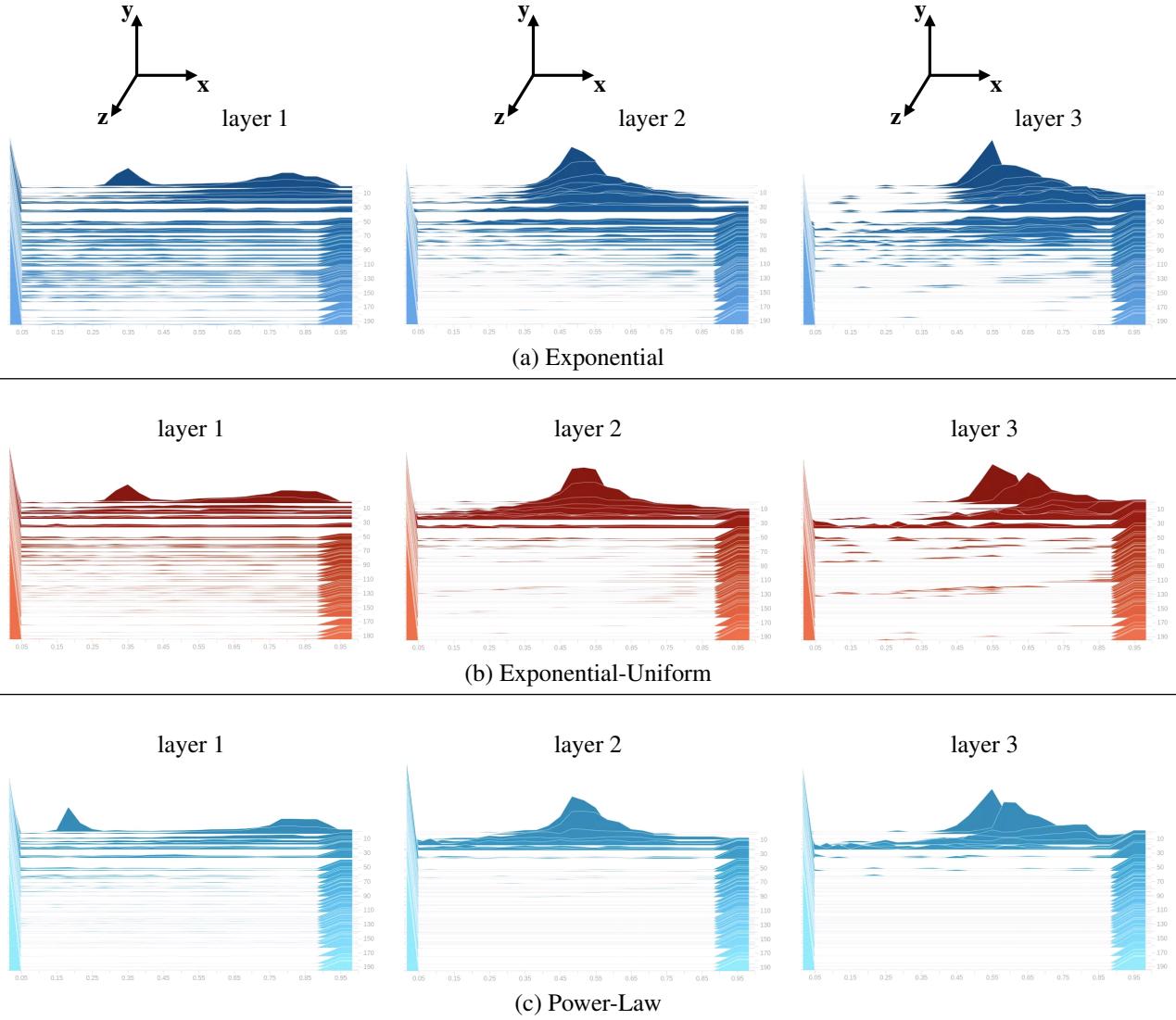


Figure 5: Visualization of the emulated binary gates ζ using LeNet-300-100 adapted on MNIST. Each column indicates the corresponding layer of LeNet-300-100. Each row refers to the mixture distribution used. All of the mixture distributions converge to Bernoulli-like distributions yet still allow for reparameterization. The x- and y- axes in each subfigure represent the histogram of the drawn ζ , while the z-axis indicates the training epochs. As we can see, ζ drawn from the MDR converges to $\delta(\zeta = 0)$ and $\delta(\zeta = 1)$ with extremely high mixing proportion at the 0-component.

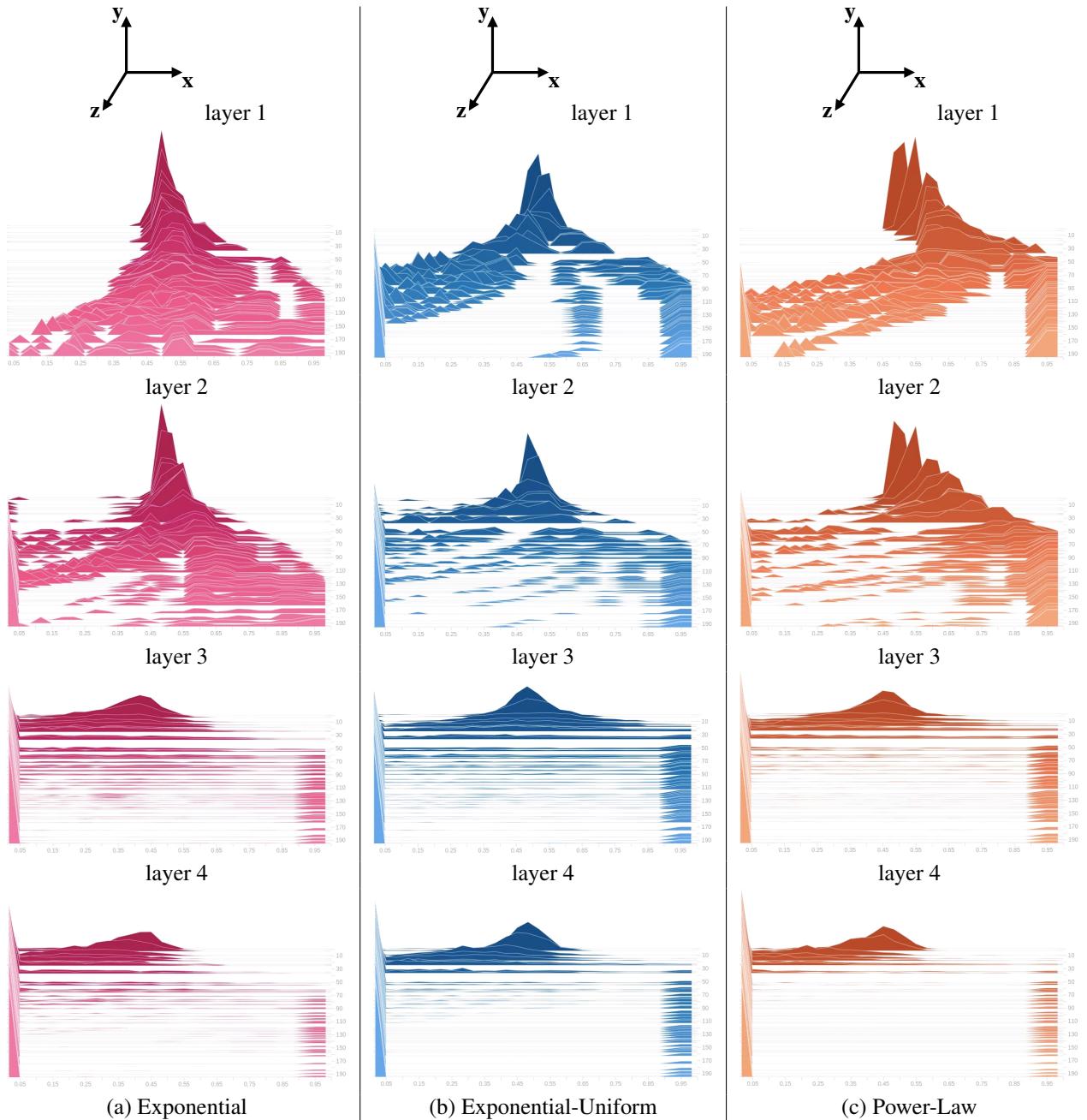


Figure 6: Visualization of the emulated binary gates ζ using LeNet-5-Caffe adapted on MNIST. Each row indicates the corresponding layer of LeNet-5-Caffe. Each column refers to the mixture distribution used. The x- and y- axes in each subfigure represent the histogram of the drawn ζ , while the z-axis indicates the training epochs. Note that the reported pruned-architecture of MDR-Exp is inferred by the trained architecture, not the estimate ζ^* . As we can observe, the convolutional layers converge much slower than the fc layers; however, ζ in general still converges to $\delta(\zeta = 0)$ and $\delta(\zeta = 1)$ at the end of training, with high mixing proportion of zero.