

A. Proof of Eq 7.

$$\begin{aligned}
Z_k^{(l+1)} &= \sum_{i=1}^{N_l} X_i^{(l)} \otimes W_{k,i}^{(l+1)} \\
&= \sum_{i=1, i \neq j}^{N_l} X_i^{(l)} \otimes W_{k,i}^{(l+1)} + X_j^{(l)} \otimes W_{k,j}^{(l+1)} \\
&\approx \sum_{i=1, i \neq j}^{N_l} B(Z_i^{(l)}) \otimes W_{k,i}^{(l+1)} + \\
&\quad B\left(\sum_{i=1, i \neq j}^{N_l} \hat{s}_i \times Z_i^{(l)}\right) \otimes W_{k,j}^{(l+1)} \\
&= \sum_{i=1, i \neq j}^{N_l} B(Z_i^{(l)}) \otimes (W_{k,i}^{(l+1)} + \hat{s}_i \times W_{k,j}^{(l+1)}), \tag{1}
\end{aligned}$$

B. Proof of Eq 10.

$$\begin{aligned}
e_p &= Z_k^{(l+1)} - \hat{Z}_k^{(l+1)} \\
&= (X_j^{(l)} - \sum_{i=1, i \neq j}^{N_l} \hat{s}_i \times X_i^{(l)}) \otimes W_{k,j}^{(l+1)} \\
&= (B(Z_j^{(l)}) - \sum_{i=1, i \neq j}^{N_l} \hat{s}_i \times B(Z_i^{(l)})) \otimes W_{k,j}^{(l+1)} \\
&= \left\{ \left\{ \frac{\gamma_j(Z_j^{(l)} - \mu_j)}{\sigma_j} + \beta_j \right\} - \right. \\
&\quad \left. \sum_{i=1, i \neq j}^{N_l} \hat{s}_i \left\{ \frac{\gamma_i(Z_i^{(l)} - \mu_i)}{\sigma_i} + \beta_i \right\} \right\} \otimes W_{k,j}^{(l+1)} \\
&= \left\{ \frac{\gamma_j}{\sigma_j} (X^{(l-1)} \otimes W_j^{(l)}) - \right. \\
&\quad \left. \sum_{i=1, i \neq j}^{N_l} \hat{s}_i \frac{\gamma_i}{\sigma_i} (X^{(l-1)} \otimes W_i^{(l)}) + \beta_j - \frac{\gamma_j \mu_j}{\sigma_j} + \right. \\
&\quad \left. \sum_{i=1, i \neq j}^{N_l} \hat{s}_i \left(\frac{\gamma_i \mu_i}{\sigma_i} - \beta_i \right) \right\} \otimes W_{k,j}^{(l+1)} \\
&= \left\{ \frac{\gamma_j}{\sigma_j} \{ X^{(l-1)} \otimes (W_j^{(l)} - \sum_{i=1, i \neq j}^{N_l} \hat{s}_i \frac{\gamma_i \sigma_j}{\sigma_i} W_i^{(l)}) \} + \right. \\
&\quad \left. (\beta_j - \frac{\gamma_j \mu_j}{\sigma_j}) - \left(\sum_{i=1, i \neq j}^{N_l} \hat{s}_i (\beta_i - \frac{\gamma_i \mu_i}{\sigma_i}) \right) \right\} \otimes W_{k,j}^{(l+1)} \tag{2}
\end{aligned}$$

C. Proof of Eq 11.

$$\begin{aligned}
e_p &= \Theta(B(Z_j^{(l)})) - \sum_{i=1, i \neq j}^{N_l} \hat{s}_i \times \Theta(B(Z_i^{(l)})) \\
&= \text{Max}(B(Z_j^{(l)}), 0) - \text{Max}\left(\sum_{i=1, i \neq j}^{N_l} \hat{s}_i \times B(Z_i^{(l)}), 0\right) \\
&= \frac{1}{2}(B(Z_j^{(l)}) + |B(Z_j^{(l)})| - \sum_{i=1, i \neq j}^{N_l} \hat{s}_i \times B(Z_i^{(l)}) - \\
&\quad | \sum_{i=1, i \neq j}^{N_l} \hat{s}_i \times B(Z_i^{(l)}) |) \\
&\leq \frac{1}{2}(B(Z_j^{(l)}) - \sum_{i=1, i \neq j}^{N_l} \hat{s}_i \times B(Z_i^{(l)}) + \\
&\quad |B(Z_j^{(l)}) - \sum_{i=1, i \neq j}^{N_l} \hat{s}_i \times B(Z_i^{(l)})|) \\
&= \frac{1}{2}(A + |A|), \tag{3}
\end{aligned}$$

D. Proof of Eq 13.

$$\begin{aligned}
e_q &= Z_k^{(l+1)} - \tilde{Z}_k^{(l+1)} \\
&= (X_m^{(l)} - \tilde{s}_m \tilde{X}_m^{(l)}) \otimes W_{k,j}^{(l+1)} \\
&= (B(Z_m^{(l)}) - \tilde{s}_m B(\tilde{Z}_m^{(l)})) \otimes W_{k,m}^{(l+1)} \\
&= \left\{ \left\{ \frac{\gamma_m(Z_m^{(l)} - \mu_m)}{\sigma_m} + \beta_m \right\} - \right. \\
&\quad \left. \tilde{s}_m \left\{ \frac{\gamma_m(\tilde{Z}_m^{(l)} - \mu_m)}{\sigma_m} + \beta_m \right\} \right\} \otimes W_{k,m}^{(l+1)} \\
&= \left(\frac{\gamma_m Z_m^{(l)}}{\sigma_m} - \tilde{s}_m \frac{\gamma_m \tilde{Z}_m^{(l)}}{\sigma_m} + \tilde{s}_m \frac{\gamma_m \mu_m}{\sigma_m} - \frac{\gamma_m \mu_m}{\sigma_m} \right. \\
&\quad \left. + \beta_m - \tilde{s}_m \beta_m \right) \otimes W_{k,m}^{(l+1)} \\
&= \left\{ \left(\frac{\gamma_m W_m^{(l)}}{\sigma_m} - \tilde{s}_m \frac{\gamma_m \tilde{W}_m^{(l)}}{\sigma_m} \right) \otimes X^{(l-1)} + \tilde{s}_m \frac{\gamma_m \mu_m}{\sigma_m} \right. \\
&\quad \left. - \frac{\gamma_m \mu_m}{\sigma_m} + \beta_m - \tilde{s}_m \beta_m \right\} \otimes W_{k,m}^{(l+1)} \tag{4}
\end{aligned}$$

E. Effectiveness of UDFC.

To better understand the behavior of UDFC in the CNN architecture, we visualize the loss landscape and weights offset of ResNet-56 on CIFAR-10.

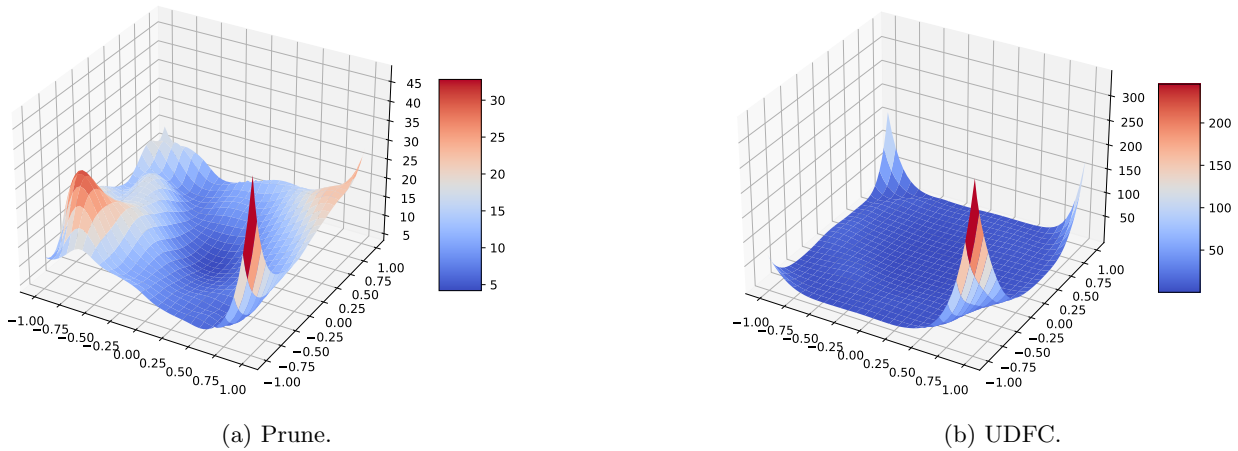


Figure 1: The loss landscape of ResNet-50 before and after compensation.

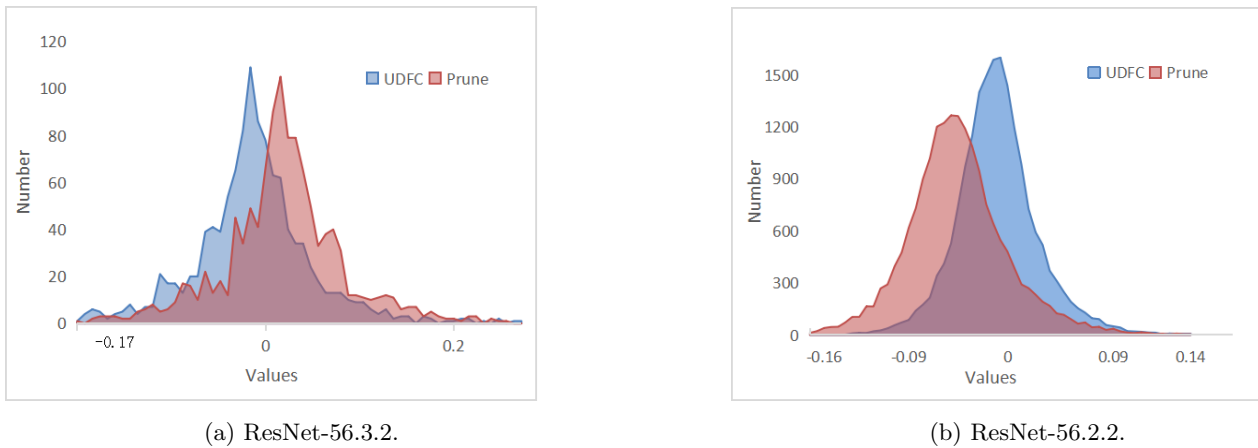


Figure 2: The loss weights offset of Prune and UDFC on ResNet-56.

Figure 1a shows the landscape of Prune and Figure 1b shows the landscape of UDFC. It is clear that the Prune method has a coarser plane of loss landscape than our method. This is due to the fact that pruning without compensation destroys the integrity of the original model. The loss landscape of UDFC becomes smoothed out, proving the effectiveness of our method. As shown in Figure 2, the compressed model produces noisy and divergent weights distribution. In contrast, our method reverts the weight average to zero, which is closer to original model.