

Contributions:

1. We extend prior work (Lin et al. (2020)) on dynamic pruning during training and lift its limitation of a hand-designed sparsity schedule.
2. We only consider a weight for pruning if it's no longer receiving reliable gradients. To measure this point, we propose to use the Mini-Batch Gradient Signal-To-Noise Ratio.
3. We benchmark our approach on different datasets and archs.

Algorithm:

Algorithm 1 *The detailed training procedure of pruning guided by gradient SNR.*

Require: uncompressed model weights $\theta \in \mathbb{R}^d$, pruned weights: $\hat{\theta}$, mask: $\mathbf{m}_{prune} \in \{0,1\}^d$; SNR exp. avg.: γ , burn-in steps: bin_{steps} , mask: $\mathbf{m} \in \{0,1\}^d$; training iterations: T .

- 1: **for** $t = 0, \dots, T$ **do**
- 2: **if** $t > bin_{steps}$ **then** ▷ trigger mask update, by default after $bin_{steps} = 1$ epoch
- 3: compute SNR mask $\mathbf{m}_{snr} \leftarrow \{snr(\theta_i^t) > 1 \mid i \text{ in } |\theta_t|\}$ ▷ let sp_{snr} be the sparsity of the resulting mask
- 4: compute MAG mask $\mathbf{m}_{mag} \leftarrow \{|\theta_i^t| > \gamma \mid i \text{ in } |\theta_t|\}$ ▷ γ cut-off weight mag. acc to sp_{snr}
- 5: compute PRUNE mask $\mathbf{m}_{prune} \leftarrow \mathbf{m}_{mag} \wedge \mathbf{m}_{snr}$ ▷ only prune, if \mathbf{m}_{mag} and \mathbf{m}_{snr} agree to prune
- 6: **end if**
- 7: $\hat{\theta}_t \leftarrow \mathbf{m}_{prune} \odot \theta_t$ ▷ apply resulting mask
- 8: compute (mini-batch) gradient $\nabla L_B(\hat{\theta})$ ▷ forward/backward pass with pruned weights $\hat{\theta}_t$
- 9: update $\mathbf{m}_{snr,t}$ ▷ update SNR exp mov. avg. per weight
- 10: $\theta_{t+1} \leftarrow$ gradient update $\nabla L_B(\hat{\theta})$ to θ_t ▷ via arbitrary optimizer (e.g. SGD with momentum)
- 11: **end for**

Ensure: θ_T and $\hat{\theta}_T$

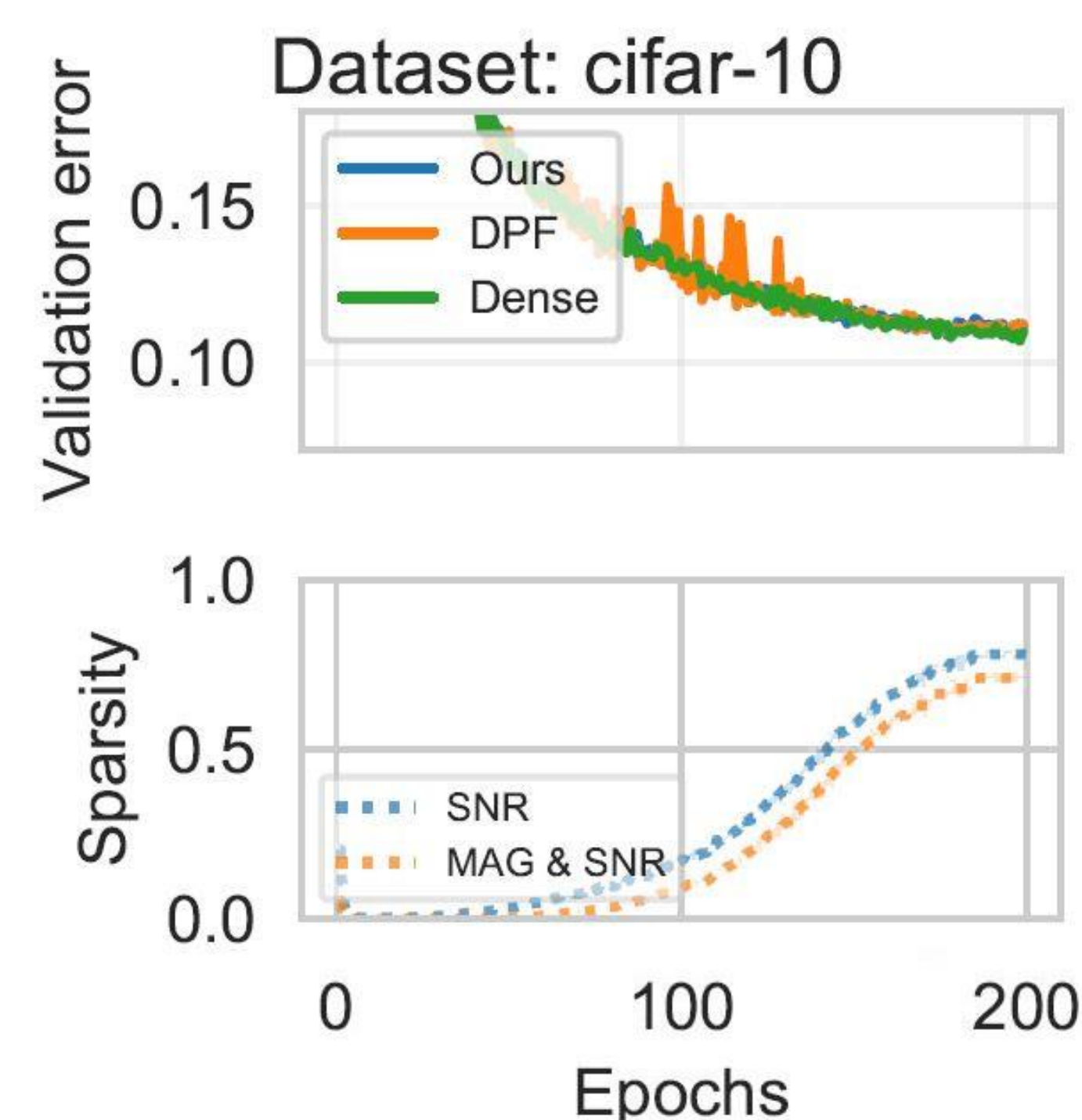
Background:

1. Mini-Batch Gradient Signal-to-Noise Ratio (Mahsereci et al. (2017))
- Measures reliability of a mini-batch gradient with respect to noise in mini-batch gradient.

$$snr(\theta_k) := \frac{\nabla L_{B,k}^2(\theta)}{\Sigma_{B,kk}(\theta)} = \begin{cases} > 1 & \text{Reliable gradient} \\ \leq 1 & \text{Noisy gradient} \end{cases}$$

2. Dynamic Pruning with Feedback (Lin et al. (2020)):
 - Increasingly prunes weights based on magnitude during training.
 - Gradients computed also for pruned weights:
 - Pruned weights can become unpruned.

Experiments:



- Sparsity depends on dataset and NN.
- Performs on par with DPF and normal dense training, but with less hyperparameters than DPF.
- Pruning hardly influences validation error (unnecessary weights removed)

References:

- Mahsereci, Maren, et al. "Early stopping without a validation set." arXiv preprint arXiv:1703.09580 (2017).
- Tao Lin, et al. "Dynamic Model Pruning with Feedback." International Conference on Learning Representations.