

A. Supplementary Material

A.1. S annealing schedule

Acceleration factor α , annealing start iteration i_{min} , annealing end iteration i_{max} , current iteration i_{cur} , $d_{target} = \frac{1}{\alpha}$.

$$d_{cur} = d_{target} + (1 - d_{target}) \left(1 - \frac{i_{cur} - i_{min}}{i_{max} - i_{min}} \right)$$

$$S_{cur} = d_{cur} \cdot D$$

A.2. Proof of Constraint Projection

Proof for Eq. 6 and Eq. 7 is taken and adapted from [38]. Transforming updated parameters $\tilde{\theta} \in \mathbb{R}^D$ into θ , which fulfills the sparsification constraint can be described as a least-squares convex problem:

$$\arg \min_{\theta} \frac{1}{2} \|\tilde{\theta} - \theta\|^2 \text{ s.t. } \sum_{i=1}^D \theta_i = \mathbf{1}^\top \theta \leq S \text{ and } 0 \leq \theta_i \leq 1. \quad (9)$$

This can be solved by the Lagrangian multiplier method:

$$\mathcal{L}(\theta, \lambda) = \frac{1}{2} \|\theta - \tilde{\theta}\|^2 + \lambda(\mathbf{1}^\top \theta - S) \quad (10)$$

$$= \frac{1}{2} \|\theta - (\tilde{\theta} - \lambda \mathbf{1})\|^2 + \lambda(\mathbf{1}^\top \tilde{\theta} - S) - \frac{n}{2} \lambda^2, \quad (11)$$

where $\lambda \geq 0$ and $0 \leq \theta_i \leq 1$. Minimizing w.r.t. θ results in

$$\bar{\theta} = \mathbf{1}_{\tilde{s} - \lambda \mathbf{1} \geq 1} + (\tilde{s} - \lambda \mathbf{1})_{1 > \tilde{s} - \lambda \mathbf{1} > 0}. \quad (12)$$

Thus, for $\lambda \geq 0$

$$\begin{aligned} g(\lambda) &= \mathcal{L}(\bar{\theta}, \lambda) \\ &= \frac{1}{2} \|[\tilde{s} - \lambda \mathbf{1}]_- + [\tilde{\theta} - (\lambda + 1)\mathbf{1}]_+\|^2 \\ &\quad + \lambda(\mathbf{1}^\top \tilde{\theta} - \theta) - \frac{D}{2} \lambda^2 \\ &= \frac{1}{2} \|[\tilde{s} - \lambda \mathbf{1}]_-\|^2 + \frac{1}{2} \|[\tilde{\theta} - (\lambda + 1)\mathbf{1}]_+\|^2 \\ &\quad + \lambda(\mathbf{1}^\top \tilde{\theta} - \theta) - \frac{D}{2} \lambda^2 \end{aligned} \quad (13)$$

and

$$\begin{aligned} g'(\lambda) &= \mathbf{1}^\top [\lambda \mathbf{1} - \tilde{\theta}]_+ + \mathbf{1}^\top [(\lambda + 1)\mathbf{1} - \tilde{\theta}]_- \\ &\quad + (\mathbf{1}^\top \tilde{\theta} - \theta) - D\lambda \\ &= \mathbf{1}^\top \min(\mathbf{1}, \max(\mathbf{0}, \tilde{\theta} - \lambda \mathbf{1})) - S \\ &= [\sum_{i=1}^D \min(1, \max(0, \tilde{\theta}_i - \lambda))] - S. \end{aligned} \quad (14)$$

With $g'(\lambda)$ being a monotone function, λ_1^* a solution for $g'(\lambda) = 0$ can be obtained by e.g. a convex solver or a

bisection method. The maximum of $g(\lambda)$ with $\lambda \geq 0$ is at $\lambda_2^* = \max(0, \lambda_1^*)$. Eventually,

$$\theta^* = \mathbf{1}_{\tilde{s} - \lambda_2^* \mathbf{1} \geq 1} + (\tilde{s} - \lambda_2^* \mathbf{1})_{1 > \tilde{s} - \lambda_2^* \mathbf{1} > 0} \quad (15)$$

$$= \min(\mathbf{1}, \max(\mathbf{0}, \tilde{\theta} - \lambda_2^* \mathbf{1})) \quad (16)$$

$$= \min(\mathbf{1}, \max(\mathbf{0}, \tilde{\theta} - \max(0, \lambda_1^*) \mathbf{1})). \quad (17)$$