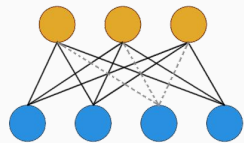


Sparsity

Training, storing and deploying NNs can be very expensive. Fortunately, their performance is **robust** to parameter pruning.

A method for obtaining efficient neural networks is by training them to encourage **sparsity** during training.



$$\|\theta\|_0 = \sum_{j=1}^{|\theta|} \mathbb{I}\{\theta_j \neq 0\}$$

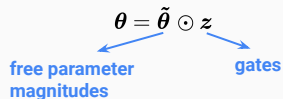
Regularization via L_0 penalties

(Louizos et al., 2018)

Penalize the number of active parameters.

$$\mathcal{R}(\theta) = \frac{1}{N} \left(\sum_{i=1}^N \ell(h(x_i; \theta), y_i) \right) + \lambda \|\theta\|_0$$

Re-parametrize with **differentiable** stochastic gates based on concrete distributions.



$$\begin{aligned} \mathcal{R}(\tilde{\theta}, \phi) &\triangleq \mathbb{E}_{\mathbf{z}} |_{\phi} \left[\mathcal{R}(\tilde{\theta} \odot \mathbf{z}) \right] \\ &= \mathbb{E}_{\mathbf{z}} |_{\phi} \left[\frac{1}{N} \sum_{i=1}^N \ell(h(x_i; \tilde{\theta} \odot \mathbf{z}), y_i) \right] + \lambda \mathbb{E}_{\mathbf{z}} |_{\phi} [\|\mathbf{z}\|_0] \end{aligned}$$

Regularization via L_0 constraints

$$\begin{aligned} \min_{\tilde{\theta}, \phi} \quad & f_{\text{obj}}(\tilde{\theta}, \phi) \triangleq \mathbb{E}_{\mathbf{z}} |_{\phi} \left[\frac{1}{N} \sum_{i=1}^N \ell(h(x_i; \tilde{\theta} \odot \mathbf{z}), y_i) \right] \\ \text{subject to} \quad & \mathfrak{g}_{\text{const}}(\phi) \triangleq \mathbb{E}_{\mathbf{z}} |_{\phi} [\|\mathbf{z}\|_0] \leq \epsilon \cdot |\theta| \\ & \text{expected \# of active params} \quad \text{In } [0, 1] \quad \text{\# of network params} \end{aligned}$$

Consider the associated Lagrangian and min-max game:

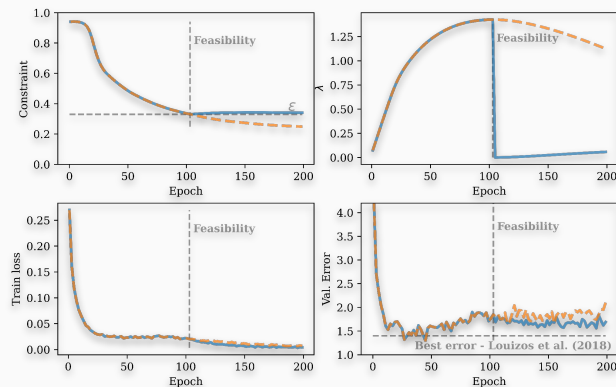
$$\begin{aligned} \mathcal{L}(\tilde{\theta}, \phi, \lambda) &\triangleq f_{\text{obj}}(\tilde{\theta}, \phi) + \lambda \left(\frac{\mathfrak{g}_{\text{const}}(\phi)}{|\theta|} - \epsilon \right) \\ \tilde{\theta}^*, \phi^*, \lambda^* &= \underset{\tilde{\theta}, \phi}{\text{argmin}} \underset{\lambda \geq 0}{\text{argmax}} \mathcal{L}(\tilde{\theta}, \phi, \lambda) \end{aligned}$$

Constraints can be liberating

- ▲ ϵ has **straightforward semantics**: the maximum proportion of active gates. Application specific requirements on sparsity can be incorporated into it.
- ▲ Tuning the penalization λ to get a satisfactory model may require running several experiments. Even harder when introducing various sources of regularization.
- ▲ It is **transparent** whether a model is respecting the sparsity constraints.



Training dynamics



Predictive performance

Architecture	Approach	Pruned	Best error (%)
MLP	Penalized [†] : $\lambda = 0.1/N$	219-214-100	1.4
	Constrained : $\epsilon = 33\%$	198-233-100	1.4
LeNet	Penalized [†] : $\lambda = 0.1/N$	20-25-45-462	0.9
	Constrained : $\epsilon = 10\%$	20-21-34-407	0.53

[†]C. Louizos, M. Welling, and D. P. Kingma. Learning Sparse Neural Networks through L_0 Regularization. *ICLR*, 2018.

The constrained approach is **not universally** better than the penalized approach!

It **can** provide more flexibility and interpretability without compromising predictive performance.