



Hyperparameter Optimization

Hyperparameter optimization can be formulated as a bilevel optimization problem.

 $\min_{\boldsymbol{\lambda}} \mathcal{L}_V(\boldsymbol{\lambda}, \mathsf{w}^*) \text{ subject to } \mathsf{w}^* = \arg\min \mathcal{L}_T(\boldsymbol{\lambda}, \mathsf{w})$

Goal: Find the optimal hyperparameters λ^* that minimizes the validation objective after training.

Parametric Best-Response Function



Best-response function maps hyperparameters to the optimal parameters w^{*}: $\mathsf{r}(oldsymbol{\lambda}) = \mathsf{arg} \min \mathcal{L}_{\mathcal{T}}(oldsymbol{\lambda},\mathsf{w})$

Substituting the best-response function in the validation objective converts the bilevel problem to a single-level problem.

Idea: Learn a **parametric best-response function** r_{ϕ} with hypernetwork and **jointly** optimize hypernetwork and hyperparameters.

$$\min_{\phi} \mathbb{E}_{\varepsilon \sim \rho(\varepsilon \mid \sigma)} \left[\mathcal{L}_{T}(\boldsymbol{\lambda} + \varepsilon, \mathsf{r}_{\phi}(\boldsymbol{\lambda} + \varepsilon)) \right], \quad \min_{\boldsymbol{\lambda}} \mathcal{L}_{V}(\boldsymbol{\lambda}, \mathsf{r}_{\phi}(\boldsymbol{\lambda} + \varepsilon)) \right]$$

Self-Tuning Networks (STNs)

Self-Tuning Networks (STNs) parameterize the best-response functions (with further structure imposed) as:

$$\mathsf{r}_{\phi}(oldsymbol{\lambda}) = \Phi oldsymbol{\lambda} + \phi_0$$
 .

Method: Jointly optimize hypernetwork and hyperparameters with structured best-response function.

Contributions

Introduce an improved hypernetwork architecture that efficiently optimizes hyperparameters online.

- Improve convergence and stability in training Self-Tuning Networks (STNs).
- Accurate approximation of the best-response Jacobian instead of the full best-response function.
- Achieve better generalization in less time, compared to existing approaches to bilevel optimization.
- Open-sourced (Δ -)Self-Tuning Networks: https://github.com/pomonam/Self-Tuning-Networks

Δ -STN: Efficient Bilevel Optimization for Neural Networks using Structured Response Jacobians

Juhan Bae, Roger Grosse

University of Toronto, Vector Institute

Fix 1: Centered Hypernetwork

Use a **centered parameterization** of the hypernetwork (with further structure imposed for neural networks):

 $\mathsf{r}_{m{ heta}}(m{\lambda},m{\lambda}_0) = m{\Theta}(m{\lambda}-m{\lambda}_0) + \mathsf{w}_0$





– –▶ : Forward Pass with Linearization

Use a **modified update rule** to separately train parameters and best-response Jacobian.

 $\arg\min \mathcal{L}_{\mathcal{T}}(\boldsymbol{\lambda}, \mathsf{w}_0), \quad \arg\min \mathbb{E}_{\varepsilon \sim p(\varepsilon | \sigma)} [\mathcal{L}_{\mathcal{T}}]$

Advantages:

- Improve the conditioning of the Gauss-Newton Hessian.
- Fix undesirable bilevel optimization dynamics.
- Eliminate any bias to the optimal w_0^* induced by the perturbation.
- Enhance the stability and convergence in training STNs.

Fix 2: Direct Approximation of Best-Response Jacobian

It is unnecessary to account for the **nonlinear effect** of large changes to w_0 , as the hypernetwork is only used to estimate the best-response Jacobian.



Directly approximate the best-response Jacobian instead of learning the full best-response function.

Network prediction with linearized hypernetwork: $\mathsf{y}' = f(\mathsf{x}, \mathsf{r}_{\theta}(\boldsymbol{\lambda}_0 + \boldsymbol{\varepsilon}, \boldsymbol{\lambda}_0), \boldsymbol{\lambda}_0 + \boldsymbol{\varepsilon}, \xi) \approx f(\mathsf{x}, \mathsf{r}_{\theta}(\boldsymbol{\lambda}_0, \boldsymbol{\lambda}_0), \boldsymbol{\lambda}_0 + \boldsymbol{\varepsilon}, \xi) + \mathsf{J}_{\mathsf{yw}}\Delta\mathsf{w}$

where $J_{yw} = \frac{\partial y}{\partial w}$ is the weight-output Jacobian. This relationship can also be written with $\Delta y \approx J_{yw}\Delta w = J_{yw}\Theta \varepsilon$. The linearization can be efficiently computed using forward mode automatic differentiation

Advantage: Accurate approximation of the best-response Jacobian.

 $(\boldsymbol{\lambda}))$

- λ_0 : "current hyperparameters", w₀: "current weights"

$$-(oldsymbol{\lambda}+arepsilon,\mathsf{r}_{oldsymbol{ heta}}(oldsymbol{\lambda}+arepsilon,oldsymbol{\lambda}))]$$

Idea: Linearize the network around the current parameters $r(\boldsymbol{\lambda}_0) = w_0$.

 $= f(\mathsf{x},\mathsf{w}_0,\boldsymbol{\lambda}_0+\boldsymbol{\varepsilon},\xi) + \mathsf{J}_{\mathsf{vw}}\boldsymbol{\Theta}\boldsymbol{\varepsilon},$

Training Algorithm for \triangle -STNs

Method: Jointly optimize the **structured** hypernetwork and hyperparameters. **Initialize:** hypernetwork $\boldsymbol{\theta} = \{w_0, \boldsymbol{\Theta}\}$, hyperparameters $\boldsymbol{\lambda}$. while not converged do for $t = 1, ..., T_{train}$ do $arepsilon \sim p(arepsilon | oldsymbol{\sigma})$ $\begin{vmatrix} \mathsf{w}_{0} \leftarrow \mathsf{w}_{0} - \alpha_{1} \nabla_{\mathsf{w}_{0}} (\mathcal{L}_{T}(\boldsymbol{\lambda}, \mathsf{r}_{\theta}(\boldsymbol{\lambda}, \boldsymbol{\lambda}))) \\ \boldsymbol{\Theta} \leftarrow \boldsymbol{\Theta} - \alpha_{1} \nabla_{\boldsymbol{\Theta}} (\mathcal{L}_{T}(\boldsymbol{\lambda} + \boldsymbol{\varepsilon}, \mathsf{r}_{\theta}(\boldsymbol{\lambda} + \boldsymbol{\varepsilon}, \boldsymbol{\lambda}))) \end{vmatrix}$ \triangleright Linearization with fwd mode autodiff. for $t = 1, ..., T_{valid}$ do $arepsilon \sim p(arepsilon | oldsymbol{\sigma})$ $\boldsymbol{\lambda} \leftarrow \boldsymbol{\lambda} - lpha_2 \nabla_{\boldsymbol{\lambda}} (\mathcal{L}_V(\boldsymbol{\lambda} + \boldsymbol{\varepsilon}, \mathsf{r}_{\boldsymbol{\theta}}(\boldsymbol{\lambda} + \boldsymbol{\varepsilon}, \boldsymbol{\lambda}_0)))|_{\boldsymbol{\lambda} = \boldsymbol{\lambda}_0}$ $\boldsymbol{\sigma} \leftarrow \boldsymbol{\sigma} - \alpha_3 \nabla_{\boldsymbol{\sigma}} (\mathcal{L}_V(\boldsymbol{\lambda} + \boldsymbol{\varepsilon}, \mathsf{r}_{\boldsymbol{\theta}}(\boldsymbol{\lambda} + \boldsymbol{\varepsilon}, \boldsymbol{\lambda})) - \tau \mathbb{H}[\boldsymbol{\rho}(\boldsymbol{\varepsilon} | \boldsymbol{\sigma})])$

Linear Regression on UCI dataset

Tune L^2 regularization and input dropout rate.



Image Classification & Language Modelling Tasks

Tune \sim 30 regularization hyperparameters (e.g. dropout rate, number of Cutout holes, random translation/degrees) on various datasets and architectures.

Dataset	Network	RS	BO	STN	Centered STN	Δ -STN
MNIST	MLP	0.043 (0.042)	0.042 (0.043)	0.043 (0.041)	0.041 (0.039)	0.040 (0.038)
FMNIST	SimpleCNN	0.206 (0.214)	0.217 (0.215)	0.196 (0.218)	0.191 (0.212)	0.189 (0.209)
CIFAR10	AlexNet	0.631 (0.671)	0.594 (0.598)	0.474 (0.488)	0.431 (0.450)	0.425 (0.446)
	VGG16	0.566 (0.595)	0.421 (0.446)	0.330 (0.354)	0.286 (0.321)	0.272 (0.296)
	ResNet18	0.264 (0.298)	0.230 (0.267)	0.266 (0.312)	0.222 (0.258)	0.204 (0.238)
PTB	LSTM	84.81 (81.46)	72.13 (69.29)	70.67 (67.78)	69.40 (66.67)	68.63 (66.26)

Figure: Final validation/test performance on image classification and language modelling tasks.



Figure: A comparison of the best validation loss achieved by random search, Bayesian optimization, STNs, and Δ -STNs over time for AlexNet. Hyperparameter schedules found by Δ -STNs.





