

CSC 2541: Neural Net Training Dynamics

Lecture 11 - Bilevel Optimization

Roger Grosse

University of Toronto, Winter 2021

Today

- Most of this course considered the optimization setting: minimizing a single cost function
- Last lecture considered differential games, where two or more “players” or “agents” simultaneously minimize/maximize different functions
 - Goal was to find a Nash equilibrium (no player can improve its utility by deviating from its current action, given the other players’ actions)
- Now we consider **bilevel optimization**: minimize a cost function defined in terms of the optimal solution to another cost function
 - In game theory, this is a **Stackelberg game**, or **leader-follower game**. The difference is that one player moves first.
 - The analogous solution concept is a **Stackelberg equilibrium**.

Bilevel Optimization

- In **bilevel (or nested) optimization**, the **outer (or upper) objective** is defined in terms of the optimal solution to an **inner (or lower) objective**.

$$\boldsymbol{\lambda}^* \in \arg \min_{\boldsymbol{\lambda}} \mathcal{J}_{\text{out}}(\boldsymbol{\lambda}, \mathbf{w}^*) \quad \text{s.t.} \quad \mathbf{w}^* \in \arg \min_{\mathbf{w}} \mathcal{J}_{\text{in}}(\boldsymbol{\lambda}, \mathbf{w}),$$

where $\boldsymbol{\lambda}$ are the **outer variables** and \mathbf{w} are the **inner variables**.

- I'll use hyperparameter optimization as a running example. Here, $\boldsymbol{\lambda}$ are the hyperparameters, and \mathbf{w} are the network weights. \mathcal{J}_{out} is the validation loss, and \mathcal{J}_{in} is the regularized training loss.
- Here, I write \in because the optimum may not be unique.

Bilevel Optimization

- For most of this lecture, I'll make the simplifying assumption that the optima are unique. In this case,

$$\boldsymbol{\lambda}^* = \arg \min_{\boldsymbol{\lambda}} \mathcal{J}_{\text{out}}(\boldsymbol{\lambda}, \mathbf{w}^*) \quad \text{s.t.} \quad \mathbf{w}^* = \arg \min_{\mathbf{w}} \mathcal{J}_{\text{in}}(\boldsymbol{\lambda}, \mathbf{w}).$$

- Assuming uniqueness, we can define the **best-response function**, or **rational reaction function**,

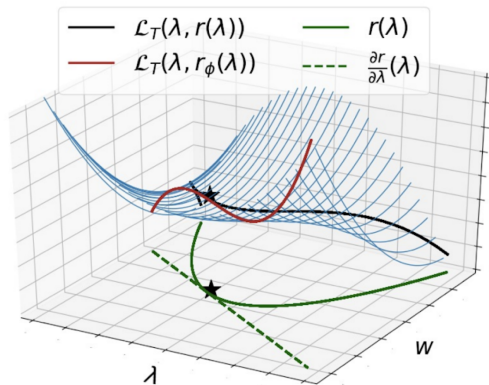
$$\mathbf{w}^* = r(\boldsymbol{\lambda}) = \arg \min_{\mathbf{w}} \mathcal{J}_{\text{in}}(\boldsymbol{\lambda}, \mathbf{w})$$

- The **Implicit Function Theorem (IFT)** proves existence under conditions which we won't worry about
- We can rewrite the optimization problem as:

$$\boldsymbol{\lambda}^* = \arg \min_{\boldsymbol{\lambda}} \mathcal{J}_{\text{out}}(\boldsymbol{\lambda}, \mathbf{r}(\boldsymbol{\lambda})) \quad \text{where} \quad \mathbf{r}(\boldsymbol{\lambda}) = \arg \min_{\mathbf{w}} \mathcal{J}_{\text{in}}(\boldsymbol{\lambda}, \mathbf{w}).$$

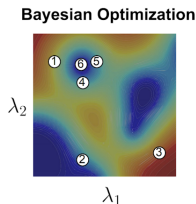
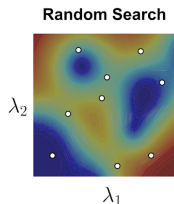
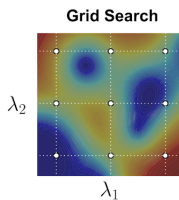
Bilevel Optimization

Example best-response function (green)



Black-Box Approaches

- One approach is to treat the outer objective as a **black box**: we can query function values, but not gradients, Hessians, etc.
- Each query: train the network and measure the validation loss
- The simplest algorithms are **non-adaptive**, like **grid search** and **random search**
- There are also **adaptive** algorithms which make use of information from past evaluations, like **Bayesian optimization**
- Drawbacks: can't use gradient information, each query is expensive
- I won't cover black-box methods since they don't raise any new NNTD issues



Hypergradient

Hypergradient

- Gradient-based optimizers are usually much more efficient than black-box ones
- To do gradient descent on λ , we need the **total gradient** of \mathcal{J}_{val} with respect to λ . This is often called the **hypergradient**, to distinguish it from the inner gradient.

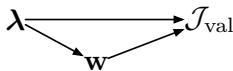
$$\overbrace{\frac{d}{d\lambda} [\mathcal{J}_{\text{val}}(\lambda, \mathbf{r}(\lambda))]}^{\text{total gradient}} = \underbrace{\frac{\partial \mathcal{J}_{\text{val}}}{\partial \lambda}(\lambda, \mathbf{r}(\lambda))}_{\text{direct gradient}} + \underbrace{\left(\overbrace{\frac{\partial \mathbf{r}}{\partial \lambda}(\lambda)}^{\text{response Jacobian}} \right)^\top \frac{\partial \mathcal{J}_{\text{val}}}{\partial \mathbf{w}}(\lambda, \mathbf{r}(\lambda))}_{\text{response gradient}}$$

- This is just the Chain Rule. In CSC2516 backprop notation,

$$\overline{\mathcal{J}_{\text{val}}} = 1$$

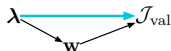
$$\overline{\mathbf{w}} = \frac{\partial \mathcal{J}_{\text{val}}}{\partial \mathbf{w}}^\top \overline{\mathcal{J}_{\text{val}}}$$

$$\overline{\lambda} = \frac{\partial \mathcal{J}_{\text{val}}}{\partial \lambda}^\top \overline{\mathcal{J}_{\text{val}}} + \frac{\partial \mathbf{w}}{\partial \lambda}^\top \overline{\mathbf{w}}$$



Hypergradient: Direct Term

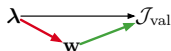
$$\frac{d}{d\boldsymbol{\lambda}} [\mathcal{J}_{\text{val}}(\boldsymbol{\lambda}, \mathbf{r}(\boldsymbol{\lambda}))] = \frac{\partial \mathcal{J}_{\text{val}}}{\partial \boldsymbol{\lambda}}(\boldsymbol{\lambda}, \mathbf{r}(\boldsymbol{\lambda})) + \left(\frac{\partial \mathbf{r}}{\partial \boldsymbol{\lambda}}(\boldsymbol{\lambda}) \right)^\top \frac{\partial \mathcal{J}_{\text{val}}}{\partial \mathbf{w}}(\boldsymbol{\lambda}, \mathbf{r}(\boldsymbol{\lambda}))$$



- For optimizing regularization hyperparameters, the direct term is typically not very interesting
 - Regularizers like dropout or data augmentation aren't applied at validation time
 - Therefore, the direct term is $\mathbf{0}$
- Example where we'd use a direct term: $\boldsymbol{\lambda}$ parameterizes a neural net architecture (# layers, # units, etc.), and we want to penalize the amount of memory or the number of arithmetic operations

Hypergradient: Response Term

$$\frac{d}{d\boldsymbol{\lambda}} [\mathcal{J}_{\text{val}}(\boldsymbol{\lambda}, \mathbf{r}(\boldsymbol{\lambda}))] = \frac{\partial \mathcal{J}_{\text{val}}}{\partial \boldsymbol{\lambda}}(\boldsymbol{\lambda}, \mathbf{r}(\boldsymbol{\lambda})) + \left(\frac{\partial \mathbf{r}}{\partial \boldsymbol{\lambda}}(\boldsymbol{\lambda}) \right)^\top \frac{\partial \mathcal{J}_{\text{val}}}{\partial \mathbf{w}}(\boldsymbol{\lambda}, \mathbf{r}(\boldsymbol{\lambda}))$$



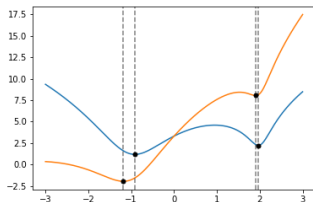
- The response term is more interesting: it says how changing $\boldsymbol{\lambda}$ influences the \mathcal{J}_{val} by way of changing \mathbf{w}^*
- The **response Jacobian** $\partial \mathbf{r} / \partial \boldsymbol{\lambda}$ measures how the optimal solution changes due to infinitesimal perturbations to $\boldsymbol{\lambda}$

Hypergradient: Response Term

- Formula for the response gradient (also given in Lecture 2):

$$\frac{\partial \mathbf{r}}{\partial \boldsymbol{\lambda}}(\boldsymbol{\lambda}) = - \underbrace{\left(\frac{\partial^2 \mathcal{J}_{\text{tr}}}{\partial \mathbf{w}^2}(\boldsymbol{\lambda}, \mathbf{r}(\boldsymbol{\lambda})) \right)^{-1}}_{=\mathbf{H}^{-1}} \frac{\partial^2 \mathcal{J}_{\text{tr}}}{\partial \boldsymbol{\lambda} \partial \mathbf{w}}(\boldsymbol{\lambda}, \mathbf{r}(\boldsymbol{\lambda}))$$

- Sanity check from Lecture 2:



$\mathcal{J}(w; \lambda) = g(w) + \lambda w$ for $\lambda = 0$ and $\lambda = 3$

- How on earth do we wind up with \mathbf{H}^{-1} ?

Hypergradient: Response Term

- We can derive the response Jacobian using a neat trick called **implicit differentiation**
- At a minimum $\mathbf{r}(\boldsymbol{\lambda})$ of the inner objective, $\partial \mathcal{J}_{\text{tr}} / \partial \mathbf{w} = \mathbf{0}$
- Since this holds for *any* $\boldsymbol{\lambda}$, the total derivative w.r.t. $\boldsymbol{\lambda}$ must be $\mathbf{0}$:

$$\frac{d}{d\boldsymbol{\lambda}} \left[\frac{\partial \mathcal{J}_{\text{tr}}}{\partial \mathbf{w}}(\boldsymbol{\lambda}, \mathbf{r}(\boldsymbol{\lambda})) \right] = \mathbf{0}$$

- We expand this total derivative using the Chain Rule:

$$\frac{d}{d\boldsymbol{\lambda}} \left[\frac{\partial \mathcal{J}_{\text{tr}}}{\partial \mathbf{w}}(\boldsymbol{\lambda}, \mathbf{r}(\boldsymbol{\lambda})) \right] = \underbrace{\frac{\partial^2 \mathcal{J}_{\text{tr}}}{\partial \boldsymbol{\lambda} \partial \mathbf{w}}(\boldsymbol{\lambda}, \mathbf{r}(\boldsymbol{\lambda}))}_{\text{direct term}} + \underbrace{\left(\frac{\partial \mathbf{r}}{\partial \boldsymbol{\lambda}}(\boldsymbol{\lambda}) \right)^\top \frac{\partial^2 \mathcal{J}_{\text{tr}}}{\partial \mathbf{w}^2}(\boldsymbol{\lambda}, \mathbf{r}(\boldsymbol{\lambda}))}_{\text{response term}}$$

- Set this equal to $\mathbf{0}$ and solve for $\partial \mathbf{r} / \partial \boldsymbol{\lambda}$:

$$\frac{\partial \mathbf{r}}{\partial \boldsymbol{\lambda}}(\boldsymbol{\lambda}) = - \left(\frac{\partial^2 \mathcal{J}_{\text{tr}}}{\partial \mathbf{w}^2}(\boldsymbol{\lambda}, \mathbf{r}(\boldsymbol{\lambda})) \right)^{-1} \frac{\partial^2 \mathcal{J}_{\text{tr}}}{\partial \boldsymbol{\lambda} \partial \mathbf{w}}(\boldsymbol{\lambda}, \mathbf{r}(\boldsymbol{\lambda}))$$

Computing the Hypergradient

Computing the Hypergradient

- Want to compute:

$$\frac{d}{d\boldsymbol{\lambda}} [\mathcal{J}_{\text{val}}(\boldsymbol{\lambda}, \mathbf{r}(\boldsymbol{\lambda}))] = \frac{\partial \mathcal{J}_{\text{val}}}{\partial \boldsymbol{\lambda}}(\boldsymbol{\lambda}, \mathbf{r}(\boldsymbol{\lambda})) + \left(\frac{\partial \mathbf{r}}{\partial \boldsymbol{\lambda}}(\boldsymbol{\lambda}) \right)^\top \frac{\partial \mathcal{J}_{\text{val}}}{\partial \mathbf{w}}(\boldsymbol{\lambda}, \mathbf{r}(\boldsymbol{\lambda}))$$

- The **direct term** and $\partial \mathcal{J}_{\text{val}} / \partial \mathbf{w}$ are easy to compute
- The hard part is multiplying by the **response Jacobian**, which requires the **inverse Hessian**:

$$\frac{\partial \mathbf{r}}{\partial \boldsymbol{\lambda}} = \underbrace{\left(\frac{\partial^2 \mathcal{J}_{\text{tr}}}{\partial \mathbf{w}^2} \right)^{-1}}_{=\mathbf{H}^{-1}} \frac{\partial^2 \mathcal{J}_{\text{tr}}}{\partial \mathbf{w} \partial \boldsymbol{\lambda}}$$

- The hypergradient is almost always estimated in one of two ways:
 - 1 Approximately solve the linear system using an iterative algorithm (e.g. CG), like many examples from this class
 - 2 Unroll the inner optimization, and backprop through it as if it were a neural net

Computation: Solving the Linear System

- **Approach 1:** Iteratively solve the linear system
- First, optimize the inner objective to convergence
 - In practice, we usually settle for approximate solutions, but the theoretical justification is unclear
- We can solve the linear system using algorithms like CG, computing the Hessian-vector products in the usual way (see Lecture 2)

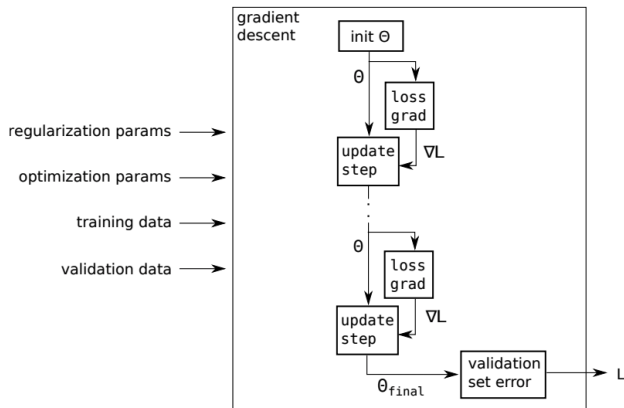
Computation: Solving the Linear System

Examples:

- Bengio (2000) used implicit differentiation to optimize ML hyperparameters (exact solution for small models)
 - Lorraine et al. (2020): optimizing millions of hyperparameters
- Influence functions (Koh and Liang, 2017): see Lecture 2, and student presentation next week
- Optimization layers (Amos and Kolter, 2017): neural net layers defined implicitly in terms of the solution to an optimization problem
 - generalized the IFT trick to constrained optimization
- Deep equilibrium models (Bai et al., 2019): student presentation next week
- Implicit MAML (Rajeswaran et al., 2019): a variant of MAML that uses implicit differentiation
 - The original MAML uses unrolling (covered next)

Computation: Unrolling

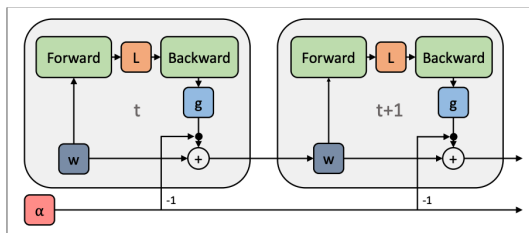
- **Approach 2:** Unroll the inner optimization
- Computation graph for gradient descent:



- Backprop through this graph in the usual way

Computation: Unrolling

- In contrast to implicit differentiation, unrolling can be used to tune optimization hyperparameters (this is known as [meta-optimization](#))
- Here's the computation graph for adapting the learning rate (known as [meta-descent](#))



- Just because you can do this doesn't mean it's a good idea. We'll see in just a bit what actually happens.

Figure: Wu et al., 2018, "Understanding short-horizon bias"

Computation: Unrolling

Examples:

- Domke (2012): learning energy-based models
- Maclaurin et al. (2015): hyperparameter optimization (student presentation next week)
 - This was also the paper that introduced Autograd, the predecessor to JAX
- MAML (Finn et al., 2017): student presentation next week
- adapting learning rates (Baydin et al., 2018)
- “Learning to learn by gradient descent by gradient descent” (Andrychowicz et al., 2016): tried to use unrolling to learn an optimization algorithm (represented as an RNN)
- differentiable neural architecture search (DARTS) (Liu et al., 2019) (unrolls only one iteration???)

Implicit Differentiation vs. Unrolling

Implicit Differentiation vs. Unrolling

Which method to use?

- Lorraine et al. (2020) related the two methods to each other.
- Suppose we unroll the inner optimization and train it to convergence. Then we use **truncated backprop through time**, which only backprops through the last K time steps
- They showed that this method is equivalent to approximately solving the IFT system by doing gradient descent on a quadratic objective

$$\frac{1}{2} \Delta \mathbf{w}^\top \mathbf{H} \Delta \mathbf{w} + \nabla_\lambda \nabla_{\mathbf{w}} \mathcal{J}_{\text{tr}}(\mathbf{w})^\top \Delta \mathbf{w}$$

- This is also equivalent to **Neumann iterations**, a method for solving linear systems

Implicit Differentiation vs. Unrolling

- **Educated guess:** implicit differentiation using Neumann iterations will behave similarly to unrolling gradient descent, for problems where they're both applicable
 - I'm not aware of any rigorous investigation of this
 - Implicit differentiation using CG should converge more efficiently for deterministic inner objectives
 - For stochastic inner objectives (e.g. neural net training), stochastic Neumann iterations (SGD on the quadratic) should be more efficient than (batch) CG in practice, if it's noise dominated rather than curvature dominated (Lecture 7)
 - This was the approach taken by Koh et al. (2017) for influence functions

Implicit Differentiation vs. Unrolling

- Implicit differentiation uses much less memory than unrolling
 - Unrolling requires storing the individual iterates (parameter vectors) along the optimization trajectory
 - A big piece of Maclaurin et al. (2015)'s work on hyperparameter optimization was a scheme for cheaply storing the parameter vectors

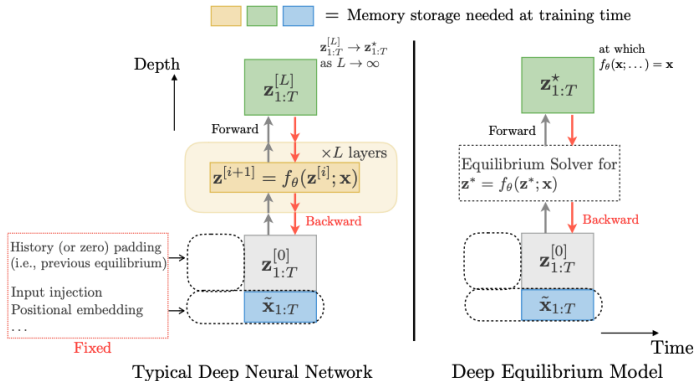


Figure: Bai et al., "Deep equilibrium models"

Implicit Differentiation vs. Unrolling

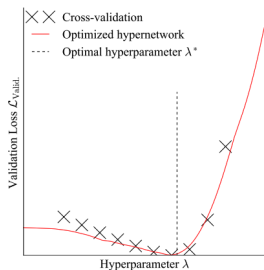
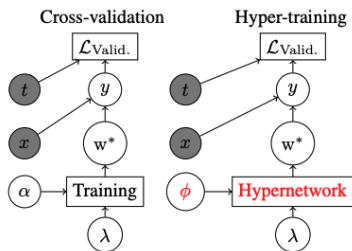
To summarize:

- Unrolled SGD is closely related to implicit differentiation with Neumann iterations
- Some advantages of implicit differentiation:
 - Lower memory costs
 - Can use faster-converging algorithms for solving the linear system (e.g. CG, Broyden's method)
- Some advantages of unrolling:
 - Trivial to implement (in JAX)
 - Can adapt optimization hyperparameters
 - Still makes sense if the inner minimization is approximate
- In practice, they're largely interchangeable, in cases where they're both applicable
- Both methods have the drawback that you have to do an inner optimization for every outer update

Hypernetworks

Hypernetworks

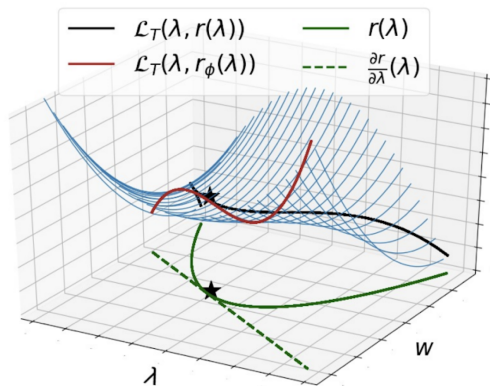
- Another approach we're starting to explore is to learn a **hypernetwork** which tries to approximate the best-response function
 - Takes in λ , outputs w
- At any time, we have a guess of the optimal λ . We want the hypernet to be accurate in the vicinity of λ



Figures: Lorraine and Duvenaud, “Stochastic hyperparameter optimization through hypernetworks”

Hypernetworks

Best-response function approximated in a parametric form (linear):



Hypernetworks

- Suppose we've somehow learned a hypernetwork \mathbf{r}_ϕ . We can compute the hypergradient by computing the derivatives of:

$$\mathcal{J}_{\text{val}}(\boldsymbol{\lambda}) = \mathcal{J}_{\text{val}}(\boldsymbol{\lambda}, \mathbf{r}_\phi(\boldsymbol{\lambda}))$$

This is just ordinary backprop. No inverse Hessian required!

- Computing the gradient requires the value and Jacobian of $\mathbf{r}_\phi(\boldsymbol{\lambda})$. Our goal in training the hypernetwork is to make sure these are accurate at $\boldsymbol{\lambda}$.
- In contrast to implicit differentiation and unrolling, this lets us amortize the cost of computing the hypergradient.

Hypernetworks

- Training iteration for the hypernetwork:
 - Sample perturbed hyperparameters $\lambda' = \lambda + \epsilon$, $\epsilon \sim \mathcal{N}(\mathbf{0}, \Sigma)$
 - Sample a training batch (index i)
 - Do the gradient update:

$$\phi \leftarrow \phi - \alpha \nabla_{\phi} \mathcal{J}_{\text{tr}}(\lambda', \mathbf{r}_{\phi}(\lambda'))$$

- Note: the perturbation scale Σ is important

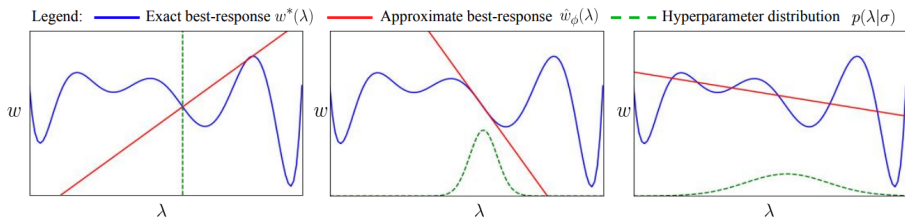


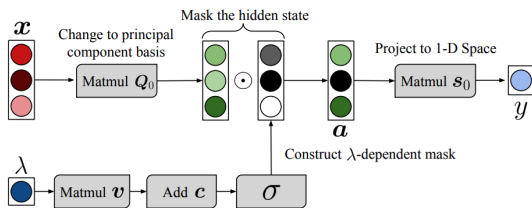
Figure: MacKay et al., “Self-tuning networks”

Self-Tuning Networks

- **Problem:** the hypernetwork is an extremely high-dimensional mapping
 - Input dimension is H (number of hyperparameters), output dimension is D (number of network parameters)
- We can linearize the network, but it's still size HD
- **Self-tuning networks (STNs)** (MacKay et al., 2019) are the first scalable approach to bilevel optimization using hypernetworks. The trick is a compact and efficient representation of the hypernet

Self-Tuning Networks

- **Some inspiration:** the following architecture exactly represents the global best-response function for linear regression



- We have an ordinary network (the **base network**) whose activations are **modulated** based on λ
- This modulation can be equivalently interpreted as rescaling the rows of the weight matrix:

$$\mathbf{Q}(\lambda) = \sigma(\lambda \mathbf{v} + \mathbf{c}) \odot_{\text{row}} \mathbf{Q}_0$$

- This is essentially how STN layers are defined (details in the paper)

Self-Tuning Networks

- The hyperparameters λ and the hypernetwork parameters ϕ are trained jointly
- This converts the bilevel optimization problem into a simultaneous game (as in Lecture 10)
- The perturbation scale σ is also adapted simultaneously

Algorithm 1 STN Training Algorithm

Initialize: Best-response approximation parameters ϕ , hyperparameters λ , learning rates $\{\alpha_i\}_{i=1}^3$

while not converged **do**

for $t = 1, \dots, T_{train}$ **do**

$\epsilon \sim p(\epsilon|\sigma)$

$\phi \leftarrow \phi - \alpha_1 \frac{\partial}{\partial \phi} f(\lambda + \epsilon, \hat{\mathbf{w}}_\phi(\lambda + \epsilon))$

for $t = 1, \dots, T_{valid}$ **do**

$\epsilon \sim p(\epsilon|\sigma)$

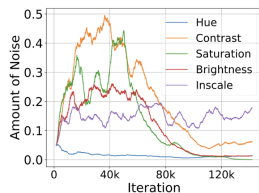
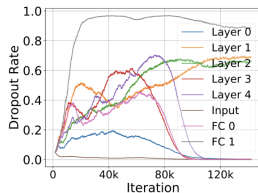
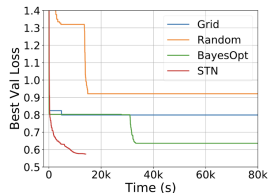
$\lambda \leftarrow \lambda - \alpha_2 \frac{\partial}{\partial \lambda} (F(\lambda + \epsilon, \hat{\mathbf{w}}_\phi(\lambda + \epsilon)) - \tau \mathbb{H}[p(\epsilon|\sigma)])$

$\sigma \leftarrow \sigma - \alpha_3 \frac{\partial}{\partial \sigma} (F(\lambda + \epsilon, \hat{\mathbf{w}}_\phi(\lambda + \epsilon)) - \tau \mathbb{H}[p(\epsilon|\sigma)])$

Self-Tuning Networks

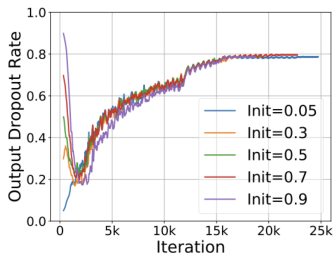
Training 15 hyperparameters of a CIFAR-10 classifier

- layer-specific dropout
- input noise
- discrete data augmentation (cutout)
- continuous data augmentation (perturb hue, saturation, brightness, contrast)



Self-Tuning Networks

- Adapting the dropout rate for an LSTM PTB language model
- Yields a schedule that seems to outperform any particular hyperparameter
- Note that the initialization is unlearned quickly



Method	Val	Test
$p = 0.68$, Fixed	85.83	83.19
$p = 0.68$ w/ Gaussian Noise	85.87	82.29
$p = 0.68$ w/ Sinusoid Noise	85.29	82.15
$p = 0.78$ (Final STN Value)	89.65	86.90
STN	82.58	79.02
LSTM w/ STN Schedule	82.87	79.93

Δ -STN

- Remember how it's a good idea to center the inputs? (Lecture 1)
- Our original STN used an uncentered parameterization of the hypernetwork:

$$\mathbf{w} = \mathbf{r}_\phi(\boldsymbol{\lambda}) = \Phi\boldsymbol{\lambda} + \phi_0$$

- The Δ -STN (Bae and Grosse, 2020) makes several algorithmic improvements, including a centered parameterization:

$$\mathbf{r}_\phi(\boldsymbol{\lambda}) = \Phi(\boldsymbol{\lambda} - \boldsymbol{\lambda}_0) + \mathbf{w}_0$$

- In Lecture 1, we understood this trick in terms of conditioning and outlier eigenvalues.
 - That explanation still applies, but for hypernetworks uncentering causes an even bigger problem.

- STN parameterization:

$$\mathbf{r}_\phi(\boldsymbol{\lambda}) = \Phi \boldsymbol{\lambda} + \phi_0$$

- Gradient descent update for Φ :

$$\Phi \leftarrow \Phi - \alpha [\nabla_{\mathbf{w}} \mathcal{J}_{\text{tr}}(\boldsymbol{\lambda}, \mathbf{w})] \boldsymbol{\lambda}^\top$$

- So early in training, approximately $\Phi \propto -\mathbf{g} \boldsymbol{\lambda}^\top$, where \mathbf{g} is the weight gradient.
 - The response Jacobian mistakenly thinks that adjusting $\boldsymbol{\lambda}$ in the direction $\boldsymbol{\lambda}$ will cause $\mathbf{w}^*(\boldsymbol{\lambda})$ to move opposite the gradient direction!

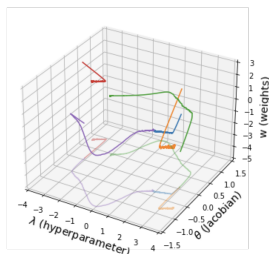
Illustrative Example:

$$\mathcal{J}_{\text{val}}(\lambda, w) = \frac{1}{10}\lambda^2 + w$$

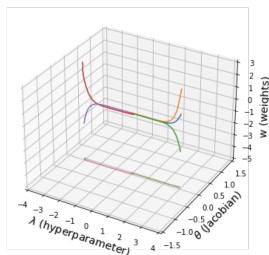
$$\mathcal{J}_{\text{tr}}(\lambda, w) = w^2$$

- This is an easy problem, since the inner objective doesn't depend on λ .
- The response Jacobian is 0, so λ and w can be optimized separately.
- Optimum: $\lambda = w = 0$.

Uncentered parameterization:



Centered parameterization:



Centering eliminates some pathological choices of hyperparameters early in training

