

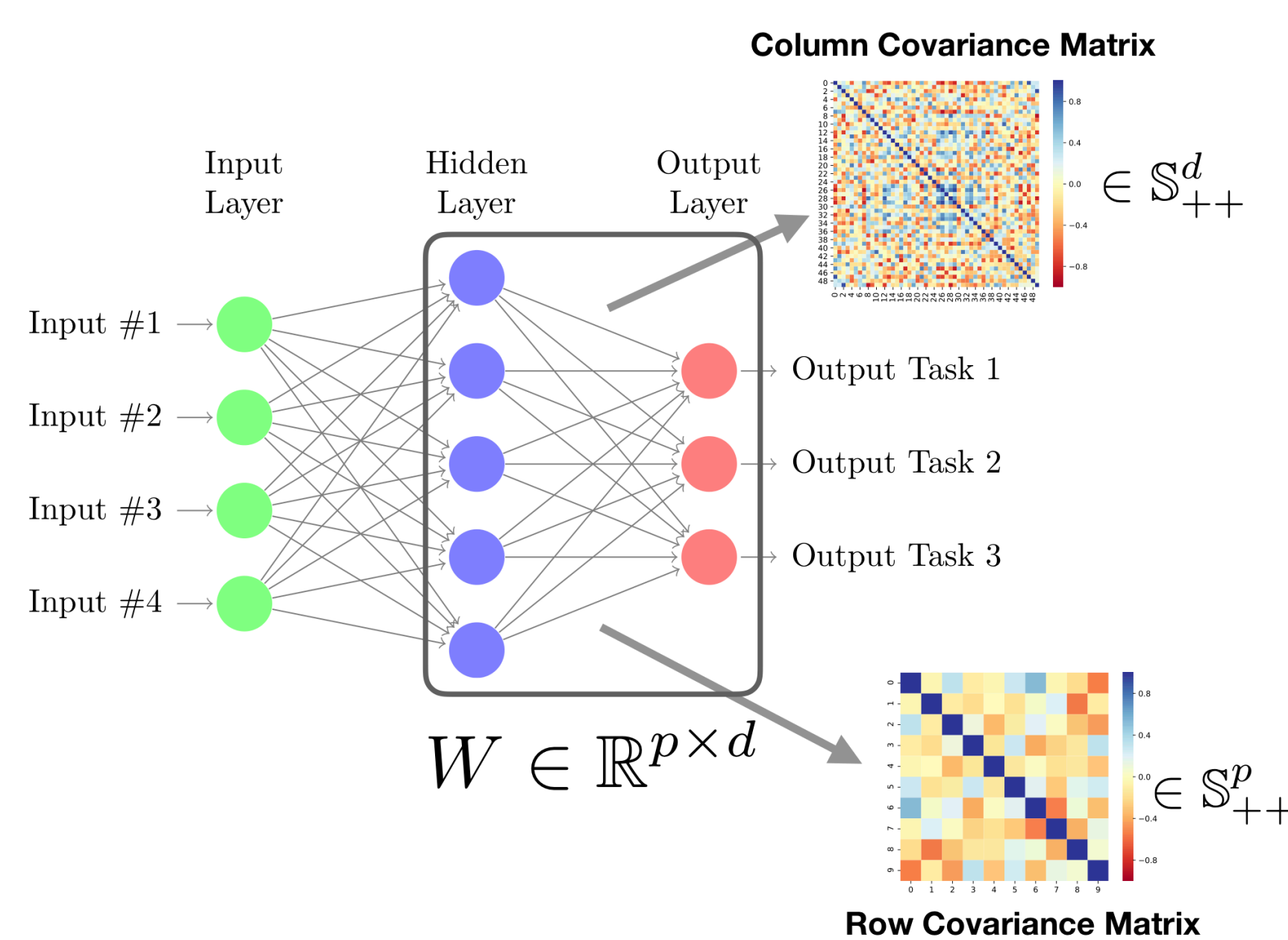
## Overview

**Q: How to effectively regularize neural networks given few available training data?**

- **AdaReg**: an approximate empirical Bayes method for regularizing NN training on small datasets
- Learn the preconditioning matrix adaptively from the data
- Significant improvement in terms of spectral norm and stable rank, leading to smaller generalization error

## Motivation

### Two-layer Neural Network



Hidden layer  $\mathbf{h} \in \mathbb{R}^p$  with output  $\hat{y} \in \mathbb{R}$  for regression

$$\hat{y} = \mathbf{a}^\top \mathbf{h}, \mathbf{h} = g(W\mathbf{x}), W \in \mathbb{R}^{p \times d}$$

Loss function

$$\ell(W, \mathbf{a}) = \frac{1}{2}(\hat{y} - y)^2$$

Update for weight matrix  $W$

$$W \leftarrow W - \gamma(\hat{y} - y)(\mathbf{a} \circ \mathbf{h}')\mathbf{x}^\top$$

$\mathbf{h}'$ : component-wise derivative of  $\mathbf{h}$  w.r.t. its input

$(\hat{y} - y)(\mathbf{a} \circ \mathbf{h}')\mathbf{x}^\top$ : gradient is always rank-1

→ rows/columns of  $W$  are correlated!

**A: Update each row/column by taking information from other rows/columns into consideration!**

## Prior of Parameters

### Multivariate Normal Distribution

$$W \sim \mathcal{MN}(\mathbf{0}_{p \times d}, \Sigma_r, \Sigma_c)$$

$\Sigma_r \in \mathbb{S}_{++}^p, \Sigma_c \in \mathbb{S}_{++}^d$ : row and column covariance matrices

$$p(W | \Sigma_r, \Sigma_c) = \frac{\exp\left(-\text{Tr}(\Sigma_r^{-1} W \Sigma_c^{-1} W^\top) / 2\right)}{(2\pi)^{pd/2} \det(\Sigma_r)^{d/2} \det(\Sigma_c)^{p/2}}$$

### Determine the Parameters in the Prior

1) Empirical Bayes: estimate the parameters of the prior from data

$$\hat{\Sigma}_r, \hat{\Sigma}_c = \arg \max_{\Sigma_r, \Sigma_c} p(\mathcal{D} | \Sigma_r, \Sigma_c) = \arg \max_{\Sigma_r, \Sigma_c} \int p(\mathcal{D} | W) p(W | \Sigma_r, \Sigma_c) dW$$

Intractable.

2) Iterative maximization of the joint distribution: sequence of MAP

$$W^{(t+1)} = \arg \max_W \log p(\mathcal{D} | W) + \log p(W | \Sigma_r^{(t)}, \Sigma_c^{(t)})$$

$$\Sigma_r^{(t+1)} = \arg \max_{\Sigma_r} \log p(\mathcal{D} | W^{(t+1)}) + \log p(W^{(t+1)} | \Sigma_r, \Sigma_c^{(t)})$$

$$\Sigma_c^{(t+1)} = \arg \max_{\Sigma_c} \log p(\mathcal{D} | W^{(t+1)}) + \log p(W^{(t+1)} | \Sigma_r^{(t+1)}, \Sigma_c)$$

Approximate empirical Bayes and tractable.

## Optimization

$$\min_{W, \mathbf{a}} \min_{\Omega_r, \Omega_c} \frac{1}{2n} \sum_{i \in [n]} (\hat{y}_i(\mathbf{x}_i; W, \mathbf{a}) - y_i)^2 + \lambda \|\Omega_r^{1/2} W \Omega_c^{1/2}\|_F^2 - \lambda(d \log \det(\Omega_r) + p \log \det(\Omega_c))$$

subject to  $uI_p \preceq \Omega_r \preceq vI_p, uI_d \preceq \Omega_c \preceq vI_d$

$\Omega_r := \Sigma_r^{-1}, \Omega_c := \Sigma_c^{-1}$ : precision matrices

### Solving $\Omega_r$

$$\min_{\Omega_r} \text{Tr}(\Omega_r W \Omega_c W^\top) - d \log \det(\Omega_r) + \mathbb{I}_{\mathcal{C}}(\Omega_r) \text{ with } \mathcal{C} := \{A \in \mathbb{S}_{++}^p \mid uI_p \preceq A \preceq vI_p\}$$

$$\rightarrow 0 \in \partial \left( \frac{1}{d} \text{Tr}(\Omega_r W \Omega_c W^\top) - \log \det(\Omega_r) + \mathbb{I}_{\mathcal{C}}(\Omega_r) \right) = W \Omega_c W^\top / d - \Omega_r^{-1} + \mathcal{N}_{\mathcal{C}}(\Omega_r) \rightarrow W \Omega_c W^\top / d - \Omega_r^{-1} \in \mathcal{N}_{\mathcal{C}}(\Omega_r^{-1})$$

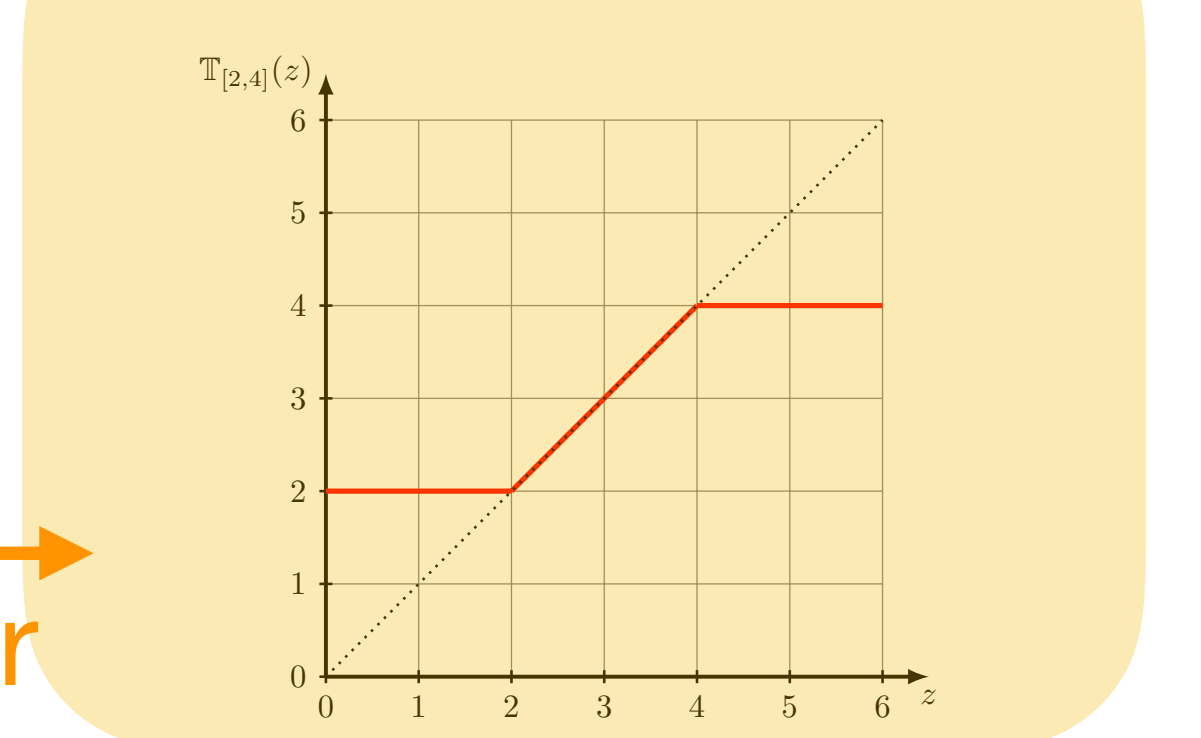
→ Optimal  $\Omega_r^{-1}$  is the Euclidean projection of  $W \Omega_c W^\top / d$  onto  $\mathcal{C}$

### Algorithm 1 Block Coordinate Descent for Adaptive Regularization

**Input:** Initial value  $\phi^{(0)} := \{\mathbf{a}^{(0)}, W^{(0)}\}, \Omega_r^{(0)} \in \mathbb{S}_{++}^p$  and  $\Omega_c^{(0)} \in \mathbb{S}_{++}^d$ , first-order optimization algorithm  $\mathfrak{A}$ .

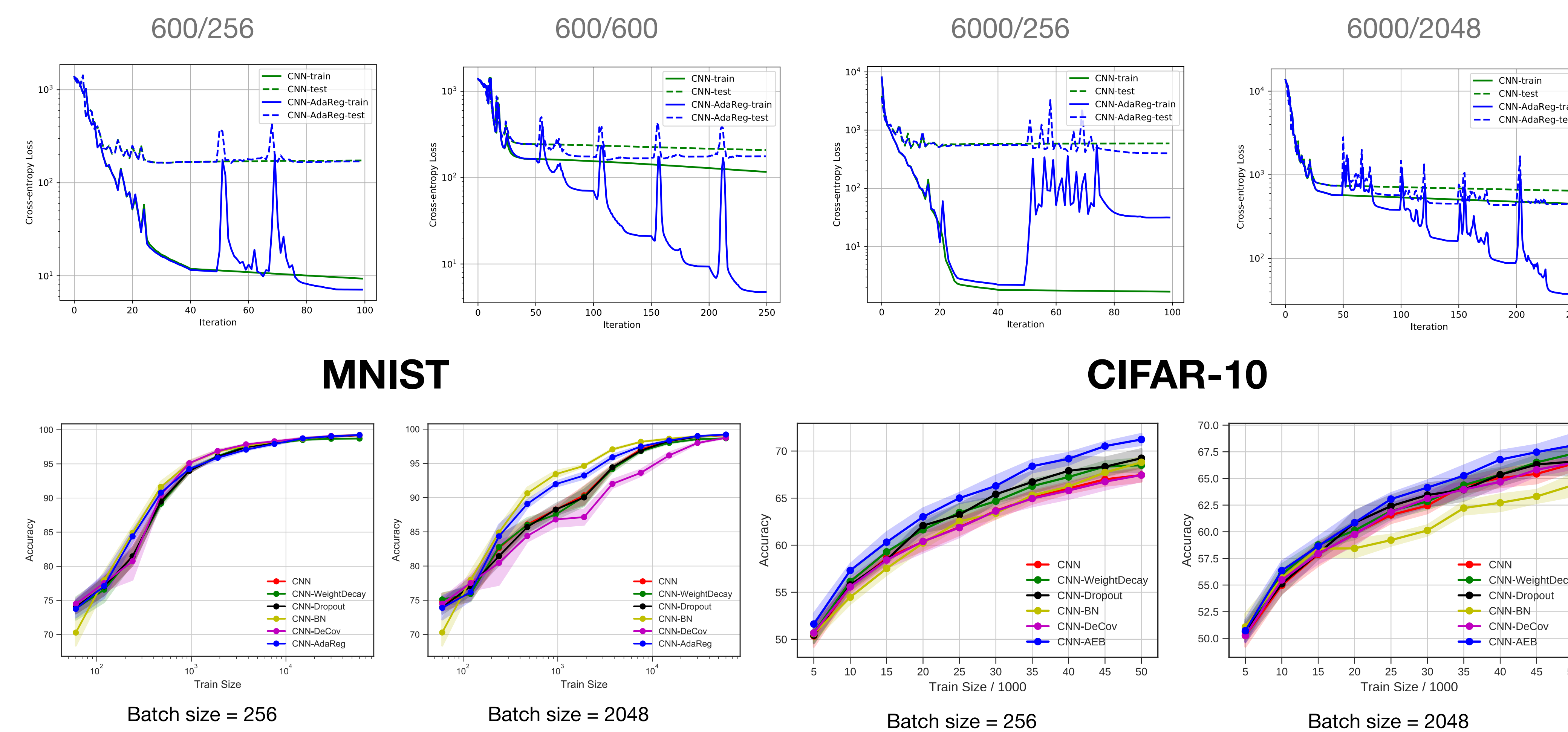
- 1: **for**  $t = 1, \dots, \infty$  until convergence **do**
- 2: Fix  $\Omega_r^{(t-1)}, \Omega_c^{(t-1)}$ , optimize  $\phi^{(t)}$  by backpropagation and algorithm  $\mathfrak{A}$
- 3:  $\Omega_r^{(t)} \leftarrow \text{INVTHRESHOLD}(W^{(t)} \Omega_c^{(t-1)} W^{(t)T}, d, u, v)$
- 4:  $\Omega_c^{(t)} \leftarrow \text{INVTHRESHOLD}(W^{(t)T} \Omega_r^{(t)} W^{(t)}, p, u, v)$
- 5: **end for**
- 6: **procedure** INVTHRESHOLD( $\Delta, m, u, v$ )
- 7: Compute SVD:  $Q \text{diag}(\mathbf{r}) Q^\top = \text{SVD}(\Delta)$
- 8: Hard thresholding  $\mathbf{r}' \leftarrow \mathbb{T}_{[u, v]}(m/\mathbf{r})$
- 9: **return**  $Q \text{diag}(\mathbf{r}') Q^\top$
- 10: **end procedure**

$$\mathbb{T}_{[u, v]}(x) := \max\{u, \min\{v, x\}\}$$



## Experiments

### AdaReg optimization trajectory on MNIST (train size/batch size)

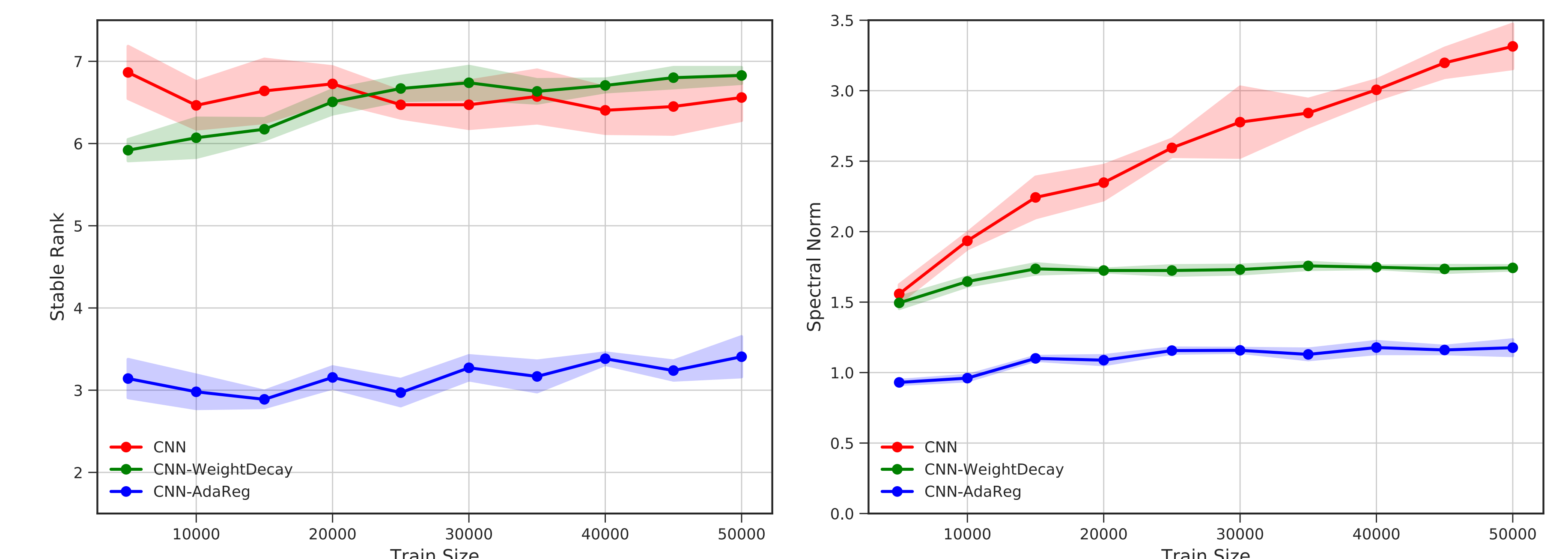


### Stable Rank and Spectral Norm for Generalization Error (Neyshabur et al. ICLR'17)

$$\text{Generalization Error} = O\left(\sqrt{\prod_{j=1}^L \|W_j\|_2^2 \sum_{j=1}^L \text{srnk}(W_j)/n}\right)$$

$$\text{srnk}(W) := \|W\|_F^2 / \|W\|_2^2$$

$$1 \leq \text{srnk}(W) \leq \text{rank}(W)$$



CIFAR-10