



Approximate Empirical Bayes for Deep Neural Networks

Han Zhao*, Yao-Hung Hubert Tsai*, Ruslan Salakhutdinov and Geoff Gordon
Machine Learning Department, Carnegie Mellon University



Summary

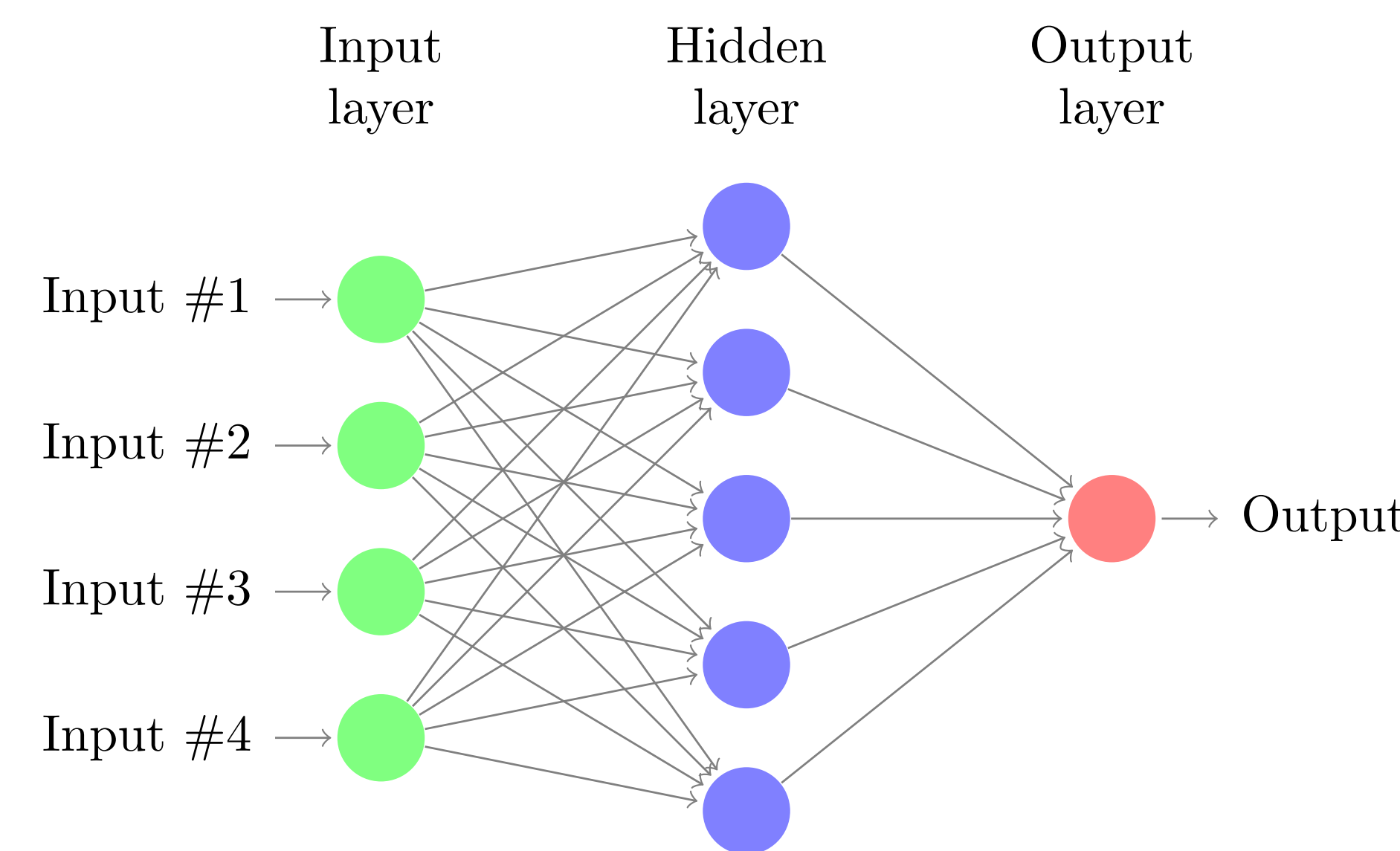
Learning Parameters in Neural Networks:

- Maximum likelihood estimation: rich neural networks overfit on small datasets.
- Regularizations: Weight decay, early stopping, dropout, DeCov, spectral regularization, etc.

Contributions:

- We propose an approximate empirical Bayes (AEB) framework for learning neural networks.
- We give a block coordinate ascent algorithm to optimize the weight matrix.

Motivating Example: A simple two-layer feed-forward neural network:



Forward propagation: $\hat{y} = \mathbf{a}^T \mathbf{h}$, $\mathbf{h} = \sigma(W\mathbf{x})$.

- Input \mathbf{x} , hidden layer \mathbf{h} , a single output $y \in \mathbb{R}$.
- Weights: W and \mathbf{a} .
- ReLU activation $\sigma(\cdot)$, ℓ_2 loss as objective function.

Backward propagation:

- $W \leftarrow W - \alpha(\hat{y} - y)(\mathbf{a} \odot \mathbf{h}')\mathbf{x}^T$.
- The gradient matrix is rank one.
- Rows/Columns of the updates are correlated with each other.

It is beneficial to learn from the experience of others.

– Bradley Efron

The Model and Algorithm

Approximate Empirical Bayes:

Introduce a matrix-variate normal distribution over the weight matrix:

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

- $\Sigma_r \in \mathbb{S}_{++}^p$: Covariance matrix over row vectors.
- $\Sigma_c \in \mathbb{S}_{++}^d$: Covariance matrix over column vectors.
- $\mathcal{MN}(A, \Sigma_r, \Sigma_c)$ is a matrix-variate normal with mean $\text{vec}(A)$ and covariance $\Sigma_r \otimes \Sigma_c$.

The true empirical Bayes approach is intractable for neural networks, hence we approximate it with the following optimization formulation:

$$\max_W \max_{\Sigma_r, \Sigma_c} \log p(\mathcal{D} | W) + \log p(W | \Sigma_r, \Sigma_c)$$

Plugging in the prior distribution leads to a constrained optimization problem:

$$\begin{aligned} \min_{W, \mathbf{a}} \min_{\Omega_r, \Omega_c} & \frac{1}{2} |\hat{y}(W, \mathbf{a}) - y|^2 + \lambda \|\Omega_r^{1/2} W \Omega_c^{1/2}\|_F^2 \\ & - \lambda (d \log \det(\Omega_r) + p \log \det(\Omega_c)) \\ \text{subject to} & uI_p \preceq \Omega_r \preceq vI_p, uI_d \preceq \Omega_c \preceq vI_d \end{aligned}$$

- $\Omega_r := \Sigma_r^{-1}$ and $\Omega_c := \Sigma_c^{-1}$ are the corresponding precision matrices.

Block Coordinate Ascent Algorithm:

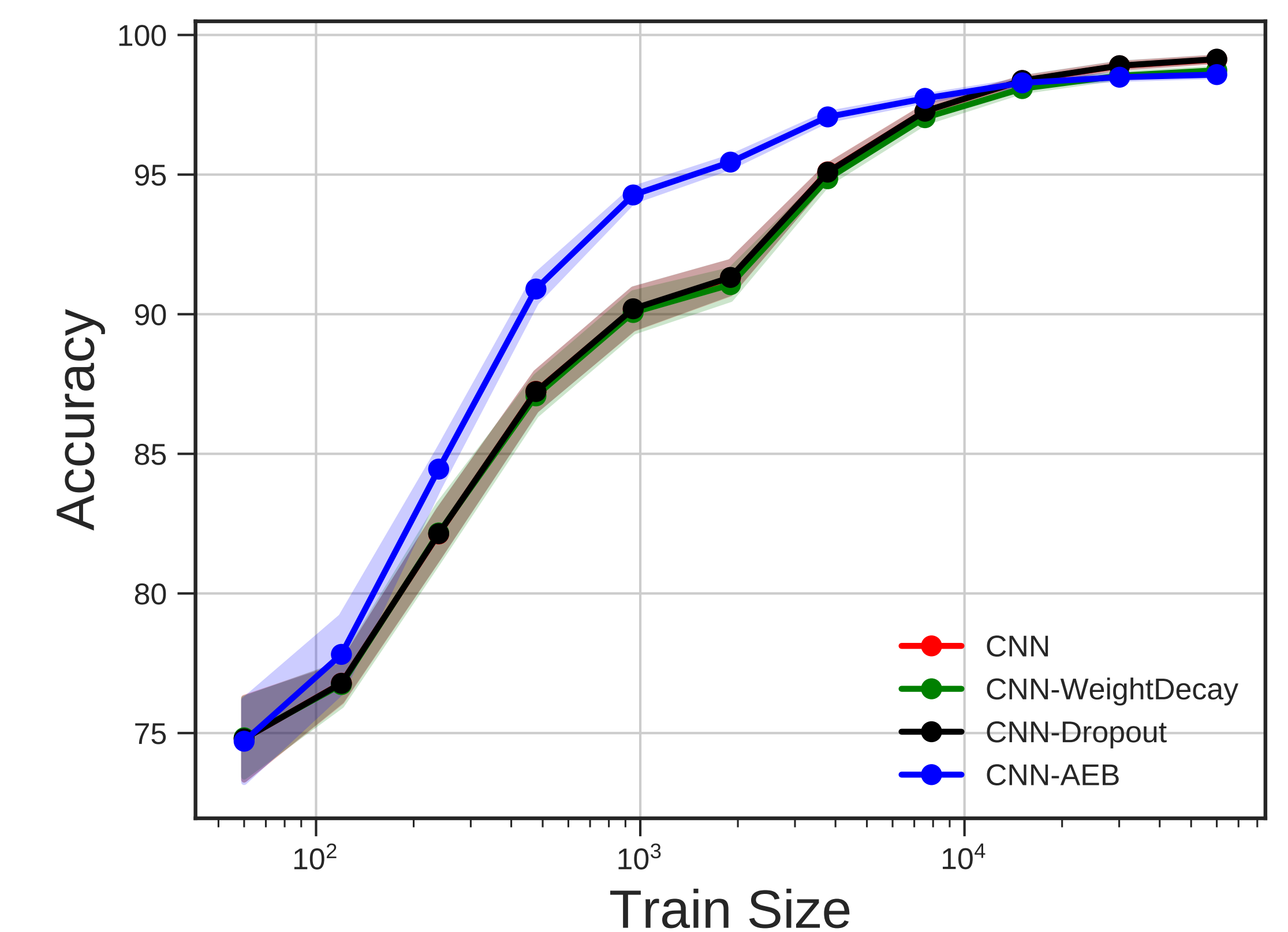
Input: Initial value $\mathbf{w}^0 := \{\mathbf{a}^{(0)}, W^{(0)}\}$, $\Omega_r^{(0)}$ and $\Omega_c^{(0)}$, first-order optimization algorithm \mathcal{A} , constants $0 < u \leq v$.

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

Experiments

Multi-class classification: MNIST dataset.

Network structure: CONV_{5×5×1×10}-CONV_{5×5×10×20}-FC_{320×50}-FC_{50×10}. Only impose prior distribution on the weight matrix of the last layer.



Multi-task regression: SARCOS dataset.

Goal: Map from a 21-dimensional input space (7 joint positions, 7 joint velocities, 7 joint accelerations) to the corresponding 7 joint torques. The training set and test set contain 44,484 and 4,449 examples. Network structure: FC_{21×256}-FC_{256×100}-FC_{100×7}.

