

Net-Trim: Convex Pruning of Deep Neural Networks with Performance Guarantee

Alireza Aghasi (Prev. IBM Research, Curr. GSU Business School), Afshin Abdi (Georgia Tech),
Nam Nguyen (IBM Research) & Justin Romberg (Georgia Tech)

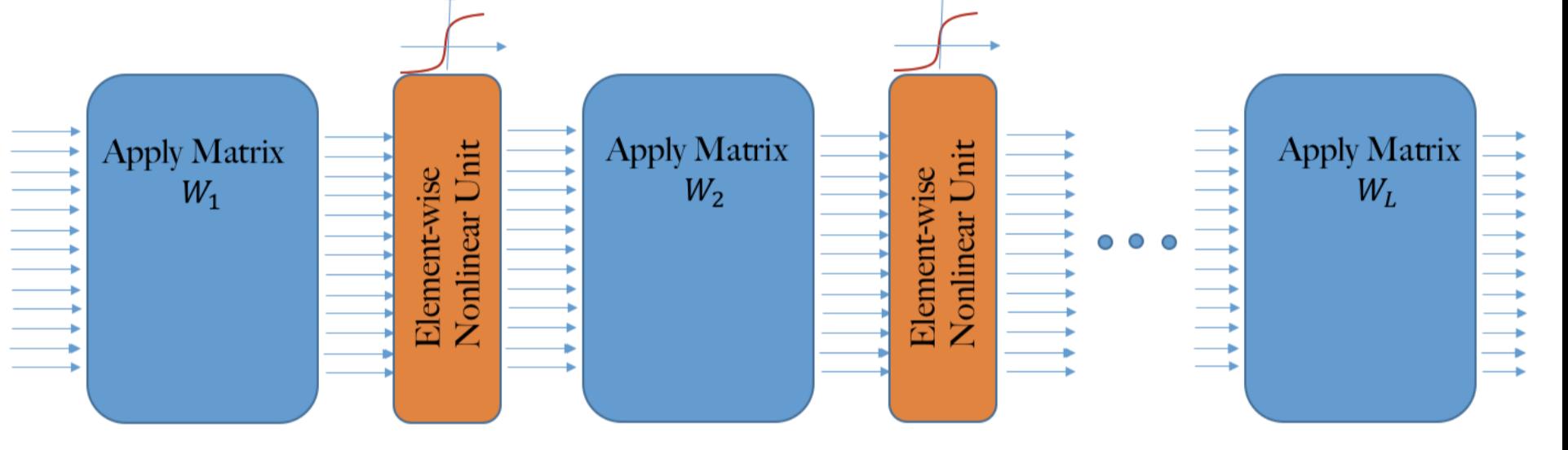
aaghasi@gsu.edu

Abstract

- Presentation and analysis of a new technique for model reduction in deep neural networks:
- Our algorithm **prunes** (sparsifies) a trained network via **convex optimization** for networks with ReLU activation
 - The program seeks a sparse set of weights at each layer that keeps the layer inputs and outputs consistent with the originally trained model
 - We present both **parallel** and **cascade** versions of the algorithm and prove an overall network consistency with the originally trained network
 - **Sample complexity:** we show that if the network response can be described using a maximum number of s non-zero weights per node, these weights can be learned from $\mathcal{O}(s \log N)$ samples
 - We present an **ADMM scheme** to implement the convex pruning scheme

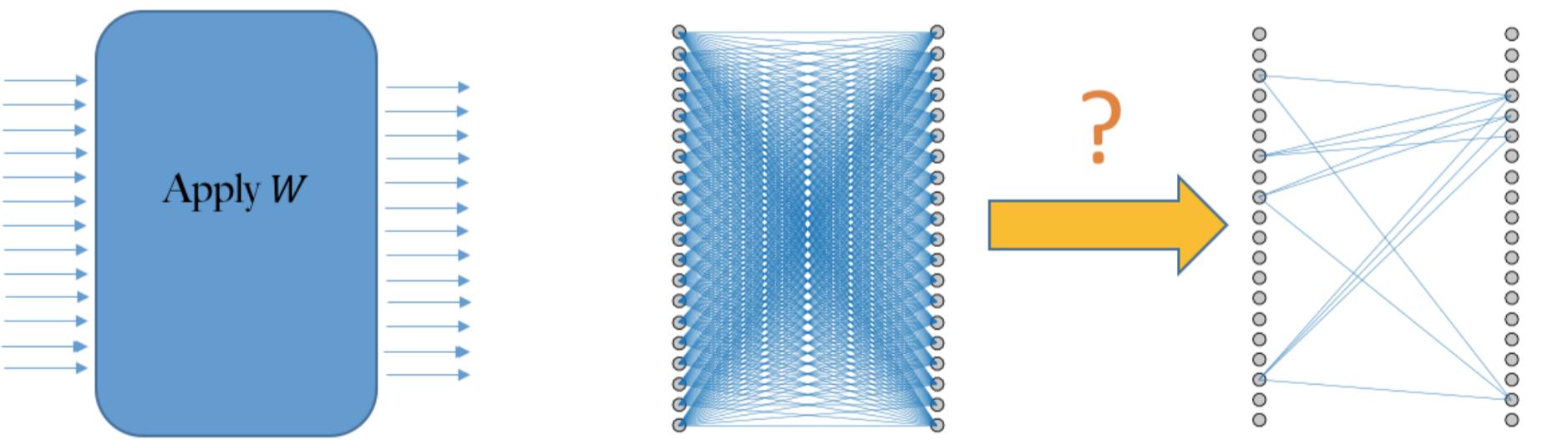
1 Introduction

- Linear models are not flexible enough to represent all systems
- Deep neural networks are nonlinear models, which are able to learn any system of arbitrary complexity!



– In linear models we aim to learn the matrix (or vector) \mathbf{W} , while in neural nets we aim to learn $\mathbf{W}_1, \dots, \mathbf{W}_L$.

- Each layer of a neural net can be represented by a graph
- An edge between two nodes is present when $W_{i,j} \neq 0$



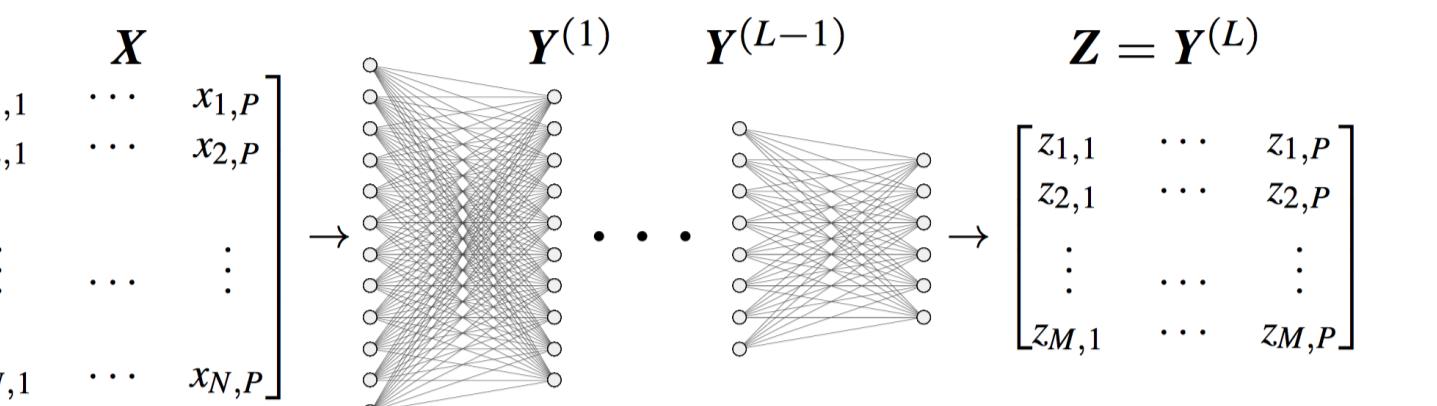
- A similar model reduction as LASSO searches for sparse matrices $\mathbf{W}_1, \dots, \mathbf{W}_L$
- We ideally want to address

$$\min_{\mathbf{W}_1, \dots, \mathbf{W}_L} \sum_{\ell=1}^L \|\mathbf{W}_\ell\|_1 \quad \text{subject to: } \|F(\mathbf{W}_1, \dots, \mathbf{W}_L, \mathbf{X}) - \mathbf{Y}\| \leq \epsilon,$$

where $F(\mathbf{W}_1, \dots, \mathbf{W}_L, \mathbf{X}) = \mathbf{W}_L^T \sigma(\mathbf{W}_{L-1}^T \sigma(\dots \sigma(\mathbf{W}_1^T \mathbf{X}) \dots))$

- The resulting program is highly nonconvex and easily gets trapped in local minima
- Efforts to convexify the neural nets objective have failed and even some proposed convex approximations are not computationally tractable [1]
- Most pruning strategies are based on heuristics, e.g., ℓ_1 penalty, Dropout, DropConnet which rely on random removal of connections/activations [3], [4], compression algorithms which rely on thresholding and retraining [2]
- Due to the complexity of the model, there are only few works with mathematical guarantees about the neural nets performance
- None of the recent works with theoretical guarantees have addressed the pruning problem

2 Network Model & Net-Trim Pruning



- Input: P training samples $x_p, p = 1, \dots, P$, stacked up as: $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_P]$
- Weight Matrices: passing through each layer

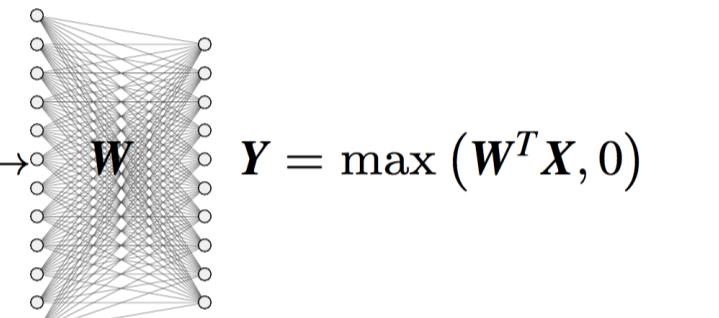
$$\mathbf{Y}^{(\ell)} = \max(\mathbf{W}^\ell \mathbf{Y}^{(\ell-1)}, 0), \quad \ell = 1, \dots, L, \quad \text{where } \mathbf{Y}^{(0)} = \mathbf{X}, \quad \mathbf{Y}^{(L)} = \mathbf{Z}.$$

- Proper Scaling: since $\max(\alpha x, 0) = \alpha \max(x, 0)$ for $\alpha > 0$, we can normalize by:

$$\mathbf{w}_\ell \rightarrow \mathbf{w}_\ell / \|\mathbf{w}_\ell\|_1, \quad \mathbf{Y}^{(\ell+1)} \rightarrow \mathbf{Y}^{(\ell+1)} / \prod_{j=0}^{\ell} \|\mathbf{w}_j\|_1.$$

Definition

A given neural network $\mathcal{T}\mathcal{N}(\{\mathbf{W}_\ell\}_{\ell=1}^L, \mathbf{X})$ is link-normalized when $\|\mathbf{w}_\ell\|_1 = 1$ for every layer $\ell = 1, \dots, L$.



- Inspired by the ℓ_1 recovery problem we consider the following program to prune a layer

$$\hat{\mathbf{W}} = \arg \min_U \|U\|_1 \quad \text{s.t.} \quad \|\max(U^T \mathbf{X}, 0) - \mathbf{Y}\|_F \leq \epsilon.$$

To convexify: knowing that $y_{m,p}$ is either zero or positive, we enforce the $\max(\cdot, 0)$ argument to be negative when $y_{m,p} = 0$, and close to $y_{m,p}$ elsewhere:

$$\left\{ \begin{array}{l} \sum_{m,p: y_{m,p}>0} (u_m^T \mathbf{x}_p - y_{m,p})^2 \leq \epsilon^2 \\ u_m^T \mathbf{x}_p \leq v_{m,p} \end{array} \right. \quad \text{for } m, p : y_{m,p} = 0 \quad \iff \quad U \in \mathcal{C}_\epsilon(\mathbf{X}, \mathbf{Y}, \mathbf{V}).$$

$$\hat{\mathbf{W}} = \arg \min_U \|U\|_1 \quad \text{s.t.} \quad U \in \mathcal{C}_\epsilon(\mathbf{X}, \mathbf{Y}, \mathbf{V}).$$

Algorithm 1 Parallel Net-Trim

- 1: Input: $\mathbf{X}, \epsilon > 0$, and link-normalized $\mathbf{W}_1, \dots, \mathbf{W}_L$
- 2: $\mathbf{Y}^{(0)} \rightarrow \mathbf{Y}$
- 3: for $\ell = 1, \dots, L$ do
- 4: $\mathbf{Y}^{(\ell)} \leftarrow \max(\hat{\mathbf{W}}_\ell \mathbf{Y}^{(\ell-1)}, 0)$
- 5: end for
- 6: for all $\ell = 1, \dots, L$ do
- 7: $\hat{\mathbf{W}}_\ell \leftarrow \arg \min_U \|U\|_1 \quad \text{s.t.} \quad U \in \mathcal{C}_\epsilon(\mathbf{Y}^{(\ell-1)}, \mathbf{Y}^{(\ell)}, \mathbf{0})$ % retraining
- 8: end for
- 9: Output: $\hat{\mathbf{W}}_1, \dots, \hat{\mathbf{W}}_L$

Theorem

Given a link-normalized network $\mathcal{T}\mathcal{N}(\{\mathbf{W}_\ell\}_{\ell=1}^L, \mathbf{X})$ with layer outcomes $\mathbf{Y}^{(\ell)}$, consider retraining each layer individually via

$$\hat{\mathbf{W}}_\ell = \arg \min_U \|U\|_1 \quad \text{s.t.} \quad U \in \mathcal{C}_{\epsilon_\ell}(\mathbf{Y}^{(\ell-1)}, \mathbf{Y}^{(\ell)}, \mathbf{0}).$$

For the retrained network $\mathcal{T}\mathcal{N}(\{\hat{\mathbf{W}}_\ell\}_{\ell=1}^L, \mathbf{X})$ with layer outcomes $\hat{\mathbf{Y}}^{(\ell)} = \max(\hat{\mathbf{W}}_\ell^T \hat{\mathbf{Y}}^{(\ell-1)}, 0)$,

$$\|\hat{\mathbf{Y}}^{(\ell)} - \mathbf{Y}^{(\ell)}\|_F \leq \sum_{j=1}^L \epsilon_j.$$

Corollary

Using Algorithm 1 ($\epsilon_1 = \dots = \epsilon_L = \epsilon$), the ultimate network outcome obeys

$$\|\hat{\mathbf{Y}}^{(L)} - \mathbf{Y}^{(L)}\|_F \leq L\epsilon.$$

Cascade Scheme:

$$\begin{aligned} \hat{\mathbf{W}}_{XY} &= \arg \min_U \|U\|_1 \quad \text{s.t.} \quad U \in \mathcal{C}_\epsilon(\mathbf{X}, \mathbf{Y}, \mathbf{0}) \\ \hat{\mathbf{W}}_{YZ} &= \arg \min_U \|U\|_1 \quad \text{s.t.} \quad U \in \mathcal{C}_{\epsilon_2}(\hat{\mathbf{Y}}, \mathbf{Z}, \mathbf{V}) \end{aligned}$$

To retrain the second layer, we ideally expect the following program

$$\min_U \|U\|_1 \quad \text{s.t.} \quad U \in \mathcal{C}_\epsilon(\hat{\mathbf{Y}}, \mathbf{Z}, \mathbf{V}).$$

However, there is no feasibility guarantee, i.e., there exists $\mathbf{W} = [w_1, \dots, w_2]$ such that

$$\left\{ \begin{array}{l} \sum_{m,p: z_{m,p}>0} (w_m^T \hat{\mathbf{y}}_p - z_{m,p})^2 \leq \epsilon^2 \\ w_m^T \hat{\mathbf{y}}_p \leq 0 \end{array} \right. \quad \text{for } m, p : z_{m,p} = 0$$

If instead of $\hat{\mathbf{Y}}$ the program was parameterized by \mathbf{Y} , a natural feasible point would have been \mathbf{W}_{YZ} . So we can plug in \mathbf{W}_{YZ} and slack the constraints to warrant feasibility.

We address the slacked program: $\min_U \|U\|_1 \quad \text{s.t.} \quad U \in \mathcal{C}_{\epsilon_2}(\hat{\mathbf{Y}}, \mathbf{Z}, \mathbf{V})$. Here, $\mathbf{V} = \mathbf{W}_{YZ}^T \hat{\mathbf{Y}}$ and

$$\epsilon_2 = \gamma \sum_{m,p: z_{m,p}>0} (w_m^T \hat{\mathbf{y}}_p - z_{m,p})^2.$$

The constants $\gamma \geq 1$ (referred to as the *inflation rate*) is a free parameter, which controls the sparsity level.

Algorithm 2 Cascade Net-Trim

- 1: Input: $\mathbf{X}, \epsilon > 0, \gamma > 1$ and link-normalized $\mathbf{W}_1, \dots, \mathbf{W}_L$
- 2: $\mathbf{Y} \leftarrow \max(\mathbf{W}^T \mathbf{X}, 0)$
- 3: $\hat{\mathbf{W}}_1 \leftarrow \arg \min_U \|U\|_1 \quad \text{s.t.} \quad U \in \mathcal{C}_\epsilon(\mathbf{X}, \mathbf{Y}, \mathbf{0})$
- 4: $\hat{\mathbf{Y}} \leftarrow \max(\hat{\mathbf{W}}_1^T \mathbf{Y}, 0)$
- 5: for $\ell = 2, \dots, L$ do
- 6: $\mathbf{Y} \leftarrow \max(\hat{\mathbf{W}}_{\ell-1}^T \mathbf{Y}, 0)$
- 7: $\epsilon \leftarrow (\gamma \sum_{m,p: y_{m,p}>0} (w_{\ell,m}^T \hat{\mathbf{y}}_p - y_{m,p})^2)^{1/2}$ % $w_{\ell,m}$ is the m -th column of \mathbf{W}_ℓ
- 8: $\hat{\mathbf{W}}_\ell \leftarrow \arg \min_U \|U\|_1 \quad \text{s.t.} \quad U \in \mathcal{C}_\epsilon(\mathbf{Y}, \mathbf{Y}, \mathbf{W}_{\ell-1}^T \mathbf{Y})$
- 9: end for
- 10: Output: $\hat{\mathbf{W}}_1, \dots, \hat{\mathbf{W}}_L$

Theorem

Given a link-normalized network $\mathcal{T}\mathcal{N}(\{\mathbf{W}_\ell\}_{\ell=1}^L, \mathbf{X})$ with layer outcomes $\mathbf{Y}^{(\ell)}$, consider a cascade retraining where

$$\epsilon_\ell^2 = \gamma \sum_{m,p: y_{m,p}^{(\ell)}>0} (w_{\ell,m}^T \hat{\mathbf{y}}_p^{(\ell)} - y_{m,p}^{(\ell)})^2,$$

$\hat{\mathbf{Y}}^{(\ell)} = \max(\hat{\mathbf{W}}_\ell^T \hat{\mathbf{Y}}^{(\ell-1)}, 0)$, $\hat{\mathbf{Y}}^{(1)} = \max(\hat{\mathbf{W}}_1^T \mathbf{Y}, 0)$ and $\gamma > 1$. For $\mathcal{T}\mathcal{N}(\{\hat{\mathbf{W}}_\ell\}_{\ell=1}^L, \mathbf{X})$ being the retrained network

$$\|\hat{\mathbf{Y}}^{(L)} - \mathbf{Y}^{(L)}\|_F \leq \epsilon_1 \sqrt{\prod_{j=2}^L \gamma_j}.$$

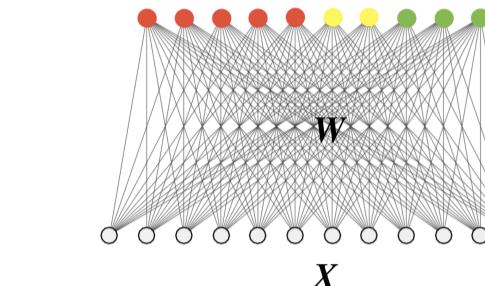
Corollary

Using Algorithm 2, the ultimate network outcome obeys

$$\|\hat{\mathbf{Y}}^{(L)} - \mathbf{Y}^{(L)}\|_F \leq \gamma^{\frac{L-1}{2}} \epsilon.$$

4 Additional Remarks & Experiments

- For each layer the retraining can be performed on separate clusters of the output nodes in parallel



- The $\mathcal{O}(s \log N/s)$ sample complexity is generalizable to all the layers under certain well conditioning of the input covariance matrix

- We have been able to present an **ADMM** scheme to address the Net-Trim core program (the details of the program are available in the paper and the **code is available online**)

In the following example, we classify 200 points in the 2D plane with a network of two hidden layers (200 neurons each) and Net-Trim generates a 93% smaller network with similar performance

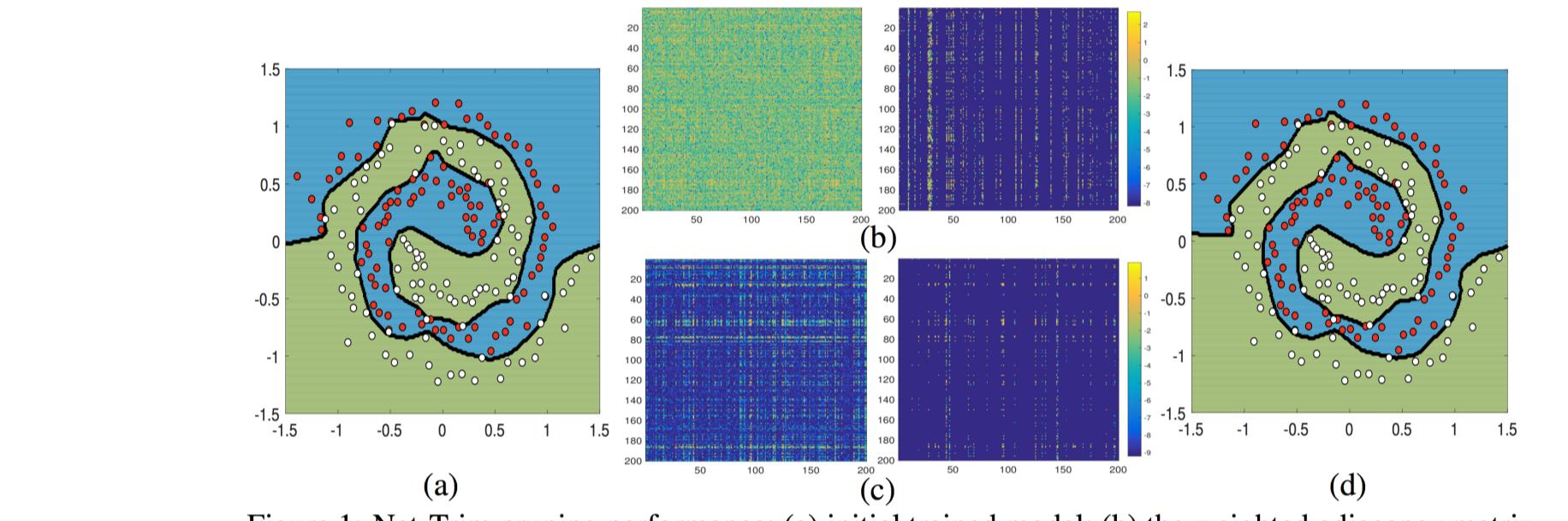


Figure 1: Net-Trim pruning performance: (a) initial trained model; (b) the weighted adjacency matrix relating the two hidden layers before (left) and after (right) the application of Net-Trim; (c) left: the adjacency matrix after training the network with Dropout and ℓ_1 regularization; right: via Net-Trim; (d) the retrained classifier

The following experiments compare Net-Trim against HPTD in [2]

Table 1: The test accuracy of different models before and after Net-Trim (NT) and HPTD [14]. Without a fine training (FT) step, Net-Trim produces pruned networks in the majority of cases more accurate than HPTD and with no risk of poor local minima. Adding an additional FT step makes Net-Trim consistently prominent

	NN2-10K	NN3-10K	NN3-60K	NN3-100K	NN3-200K
Init. Mod. Acc. (%)	95.59	97.58	98.18	98.37	99.11
Total Pruning (%)	75.87	75.82	77.48	76.45	77.52
NT + FT Acc. (%)	1.98	1.98	1.77	1.98	1.98
NT No FT Acc. (%)	95.47	97.55	98.1	93.91	93.74
HPTD No FT Acc. (%)	93.3	97.34	8.92	19.17	55.92
HPTD No FT Acc. (%)	95.85	97.67	98.12	98.35	99.21
NT + FT Acc. (%)	93.56	97.32	EMT	98.16	EMT
HPTD + FT Acc. (%)	95.61	EMT	97.96	EMT	99.01
HPTD + FT Acc. (%)	95.61	EMT	97.96	EMT	99.00

