

Supplemental Note: Net-Trim Implementation for Convolutional Layers

Alireza Aghasi*

Afshin Abdi†

Justin Romberg†

Abstract

This short note presents the details of implementing the the Net-Trim algorithm for convolutional layers. Net-Trim requires a least-squares solve at each iteration and the main idea is to employ the conjugate gradient (CG) method in operator form. This way the least squares step can be efficiently carried out via a series of forward/adjoint operations.

1 Net-Trim in Operator Form

For $\mathbf{W} \in \mathbb{R}^{N \times M}$, $\mathbf{X}^{in} \in \mathbb{R}^{N \times P}$, and $\Omega \subseteq \{1, \dots, M\} \times \{1, \dots, P\}$ the Net-Trim central program in the matrix form is

$$\underset{\mathbf{W}}{\text{minimize}} \quad \|\mathbf{W}\|_1 \quad \text{subject to} \quad \begin{cases} \|(\mathbf{W}^\top \mathbf{X}^{in} - \mathbf{X}^{out})_\Omega\|_F \leq \epsilon \\ (\mathbf{W}^\top \mathbf{X}^{in})_{\Omega^c} \leq \mathbf{V}_{\Omega^c} \end{cases} . \quad (1)$$

When $\mathbf{X}^{in} \in \mathbb{T}_{in}$, $\mathbf{W} \in \mathbb{T}_w$ and $\mathbf{X}^{out} \in \mathbb{T}_{out}$ are tensors, and Ω is a set of indices on the tensor elements, our central program takes the following form:

$$\underset{\mathbf{W}}{\text{minimize}} \quad \|\mathbf{W}\|_1 \quad \text{subject to} \quad \begin{cases} \|(\mathcal{A}_{\mathbf{X}^{in}}(\mathbf{W}) - \mathbf{X}^{out})_\Omega\|_F \leq \epsilon \\ (\mathcal{A}_{\mathbf{X}^{in}}(\mathbf{W}))_{\Omega^c} \leq \mathbf{V}_{\Omega^c} \end{cases} , \quad (2)$$

where $\|\cdot\|_1$ and $\|\cdot\|_F$ naturally apply to the vectorized tensors. To clarify the notation, given a tensor \mathbf{Z} and an index set Ω , \mathbf{Z}_Ω is a tensor of similar size as \mathbf{Z} , which takes identical values as \mathbf{Z} on Ω and takes zero values on Ω^c .

The operator $\mathcal{A}_{\mathbf{X}^{in}} : \mathbb{T}_w \rightarrow \mathbb{T}_{out}$ is a linear operator that is parameterized by \mathbf{X}^{in} . For instance in convolutional layers it is a tensor convolution operator with one of the operands being \mathbf{X}^{in} . Throughout the text we assume that $\mathcal{A}_{\mathbf{X}^{in}}^* : \mathbb{T}_{out} \rightarrow \mathbb{T}_w$ is the adjoint operator. The adjoint operator needs to satisfy the following property:

$$\forall \mathbf{W} \in \mathbb{T}_w, \forall \mathbf{Z} \in \mathbb{T}_{out} : \quad \langle \mathcal{A}_{\mathbf{X}^{in}}(\mathbf{W}), \mathbf{Z} \rangle_{\mathbb{T}_{out}} = \langle \mathbf{W}, \mathcal{A}_{\mathbf{X}^{in}}^*(\mathbf{Z}) \rangle_{\mathbb{T}_w} .$$

This is the fundamental property that we verify when forming the adjoint operator.

* (Corresponding Author) Robinson College of Business, Georgia State University, Atlanta, GA. Email: aaghasi@gsu.edu

† School of Electrical and Computer Engineering, Georgia Tech, Atlanta, GA. Emails: {abdi,jrom}@ece.gatech.edu.

Program (2) can be equivalently cast as

$$\begin{aligned} & \underset{\substack{\mathbf{W}^{(1)} \in \mathbb{T}_{out} \\ \mathbf{W}^{(2)}, \mathbf{W}^{(3)} \in \mathbb{T}_w}}{\text{minimize}} & f_1(\mathbf{W}^{(1)}) + f_2(\mathbf{W}^{(2)}) & \text{subject to} & \begin{cases} \mathbf{W}^{(1)} = \mathcal{A}_{\mathbf{X}^{in}}(\mathbf{W}^{(3)}) \\ \mathbf{W}^{(2)} = \mathbf{W}^{(3)} \end{cases}, \end{aligned} \quad (3)$$

where

$$f_1(\mathbf{W}) = \mathcal{I}_{\|\mathbf{W}_{\Omega} - \mathbf{X}_{\Omega}^{out}\|_F \leq \epsilon}(\mathbf{W}) + \mathcal{I}_{\mathbf{W}_{\Omega^c} \leq \mathbf{V}_{\Omega^c}}(\mathbf{W}), \text{ and } f_2(\mathbf{W}) = \|\mathbf{W}\|_1,$$

and $\mathcal{I}_C(\cdot)$ represents the indicator function of the set C :

$$\mathcal{I}_C(\mathbf{W}) = \begin{cases} 0 & \mathbf{W} \in C \\ +\infty & \mathbf{W} \notin C \end{cases}.$$

For the convex program (3), the ADMM update for each variable at the k -th iteration follows the standard forms

$$\mathbf{W}_{k+1}^{(1)} = \arg \min_{\mathbf{W}} f_1(\mathbf{W}) + \frac{\rho}{2} \left\| \mathbf{W} + \mathbf{U}_k^{(1)} - \mathcal{A}_{\mathbf{X}^{in}}(\mathbf{W}_k^{(3)}) \right\|_F^2, \quad (4)$$

$$\mathbf{W}_{k+1}^{(2)} = \arg \min_{\mathbf{W}} f_2(\mathbf{W}) + \frac{\rho}{2} \left\| \mathbf{W} + \mathbf{U}_k^{(2)} - \mathbf{W}_k^{(3)} \right\|_F^2, \quad (5)$$

$$\mathbf{W}_{k+1}^{(3)} = \arg \min_{\mathbf{W}} \frac{\rho}{2} \left\| \mathbf{W}_{k+1}^{(1)} + \mathbf{U}_k^{(1)} - \mathcal{A}_{\mathbf{X}^{in}}(\mathbf{W}) \right\|_F^2 + \frac{\rho}{2} \left\| \mathbf{W}_{k+1}^{(2)} + \mathbf{U}_k^{(2)} - \mathbf{W} \right\|_F^2, \quad (6)$$

and the dual updates are performed via

$$\mathbf{U}_{k+1}^{(1)} = \mathbf{U}_k^{(1)} + \mathbf{W}_{k+1}^{(1)} - \mathcal{A}_{\mathbf{X}^{in}}(\mathbf{W}_k^{(3)}), \quad \mathbf{U}_{k+1}^{(2)} = \mathbf{U}_k^{(2)} + \mathbf{W}_{k+1}^{(2)} - \mathbf{W}_{k+1}^{(3)}.$$

Obtaining closed-form expressions for the programs (4) and (5) is identical to the Net-Trim derivation in ([AAR18]) and we bring them here for completeness. For $\mathbf{W}^{(1)}$ update we have

$$\left(\mathbf{W}_{k+1}^{(1)} \right)_{\Omega} = \arg \min_{\mathbf{W}_{\Omega}: \|\mathbf{W}_{\Omega} - \mathbf{X}_{\Omega}^{out}\|_F \leq \epsilon} \frac{\rho}{2} \left\| \mathbf{W}_{\Omega} - \left(\mathcal{A}_{\mathbf{X}^{in}}(\mathbf{W}_k^{(3)}) - \mathbf{U}_k^{(1)} \right) \right\|_F^2,$$

and the closed form expression is obtained using the fact that

$$\arg \min_{\mathbf{W}_{\Omega}: \|\mathbf{W}_{\Omega} - \mathbf{Z}_{\Omega}\|_F \leq \epsilon} \frac{\rho}{2} \left\| \mathbf{W}_{\Omega} - \mathbf{Y}_{\Omega} \right\|_F^2 = \begin{cases} \mathbf{Y}_{\Omega} & \text{if } \|\mathbf{Y}_{\Omega} - \mathbf{Z}_{\Omega}\|_F \leq \epsilon \\ \mathbf{Z}_{\Omega} + \epsilon \frac{\mathbf{Y}_{\Omega} - \mathbf{Z}_{\Omega}}{\|\mathbf{Y}_{\Omega} - \mathbf{Z}_{\Omega}\|_F} & \text{else} \end{cases}.$$

Also

$$\left(\mathbf{W}_{k+1}^{(1)} \right)_{\Omega^c} = \arg \min_{\mathbf{W}_{\Omega^c}: \mathbf{W}_{\Omega^c} \leq \mathbf{V}_{\Omega^c}} \frac{\rho}{2} \left\| \mathbf{W}_{\Omega^c} - \left(\mathcal{A}_{\mathbf{X}^{in}}(\mathbf{W}_k^{(3)}) - \mathbf{U}_k^{(1)} \right) \right\|_F^2,$$

to obtain which we use the fact that

$$\arg \min_{\mathbf{W}_{\Omega^c}: \mathbf{W}_{\Omega^c} \leq \mathbf{V}_{\Omega^c}} \frac{\rho}{2} \|\mathbf{W}_{\Omega^c} - \mathbf{Y}_{\Omega^c}\|_F^2 = \mathbf{Y}_{\Omega^c} - (\mathbf{Y}_{\Omega^c} - \mathbf{V}_{\Omega^c})^+.$$

Finally, the solution to (5) is the standard soft thresholding operator (e.g., see §4.4.3 of [BPC⁺11]), which reduces the update to

$$\mathbf{W}_{k+1}^{(2)} = S_{1/\rho} \left(\mathbf{W}_k^{(3)} - \mathbf{U}_k^{(2)} \right), \text{ where } S_c(w) = \begin{cases} w - c & w > c \\ 0 & |w| \leq c \\ w + c & w < -c \end{cases}.$$

Note that when $S_{1/\rho}$ applies to a tensor, it applies to each element individually. After combining the steps above, we propose Algorithm 1 as a computational scheme to address the Net-Trim in operator form.

<p>Algorithm 1: Implementation of the Net-Trim In Operator Form</p> <pre> Input: $\mathbf{X}^{in} \in \mathbb{T}_{in}$, $\mathbf{X}^{out} \in \mathbb{T}_{out}$, Ω, \mathbf{V}_{Ω}, ϵ, ρ initialize $\mathbf{U}^{(1)} \in \mathbb{T}_{out}$, $\mathbf{U}^{(2)} \in \mathbb{T}_w$ and $\mathbf{W}^{(3)} \in \mathbb{T}_w$ % all initializations can be with $\mathbf{0}$ while not converged do $\mathbf{Y} \leftarrow \mathcal{A}_{\mathbf{X}^{in}}(\mathbf{W}^{(3)}) - \mathbf{U}^{(1)}$ if $\ \mathbf{Y}_{\Omega} - \mathbf{X}_{\Omega}^{out}\ _F \leq \epsilon$ then $\mathbf{W}_{\Omega}^{(1)} \leftarrow \mathbf{Y}_{\Omega}$ else $\mathbf{W}_{\Omega}^{(1)} \leftarrow \mathbf{X}_{\Omega}^{out} + \epsilon \ \mathbf{Y}_{\Omega} - \mathbf{X}_{\Omega}^{out}\ _F^{-1} (\mathbf{Y}_{\Omega} - \mathbf{X}_{\Omega}^{out})$ end $\mathbf{W}_{\Omega^c}^{(1)} \leftarrow \mathbf{Y}_{\Omega^c} - (\mathbf{Y}_{\Omega^c} - \mathbf{V}_{\Omega^c})^+$ $\mathbf{W}^{(2)} \leftarrow S_{1/\rho}(\mathbf{W}^{(3)} - \mathbf{U}^{(2)})$ % $S_{1/\rho}$ applies to each element of the tensor $\mathbf{W}^{(3)} \leftarrow \arg \min_{\mathbf{W}} \frac{1}{2} \ \mathcal{A}_{\mathbf{X}^{in}}(\mathbf{W}) - (\mathbf{W}^{(1)} + \mathbf{U}^{(1)})\ _F^2 + \frac{1}{2} \ \mathbf{W} - (\mathbf{W}^{(2)} + \mathbf{U}^{(2)})\ _F^2$ $\mathbf{U}^{(1)} \leftarrow \mathbf{U}^{(1)} + \mathbf{W}^{(1)} - \mathcal{A}_{\mathbf{X}^{in}}(\mathbf{W}^{(3)})$ $\mathbf{U}^{(2)} \leftarrow \mathbf{U}^{(2)} + \mathbf{W}^{(2)} - \mathbf{W}^{(3)}$ end Output: $\mathbf{W}^{(3)}$ </pre>
--

The only undiscussed part in Algorithm 1 is the update for $\mathbf{W}^{(3)}$ which we address using an operator form of the conjugate gradient algorithm explained in the next section.

2 Least Squares Update Using an Operator Conjugate Gradient

In this section we address the minimization

$$\underset{\mathbf{W} \in \mathbb{T}_w}{\text{minimize}} \frac{1}{2} \|\mathcal{A}_{\mathbf{X}^{in}}(\mathbf{W}) - \mathbf{B}\|_F^2 + \frac{1}{2} \|\mathbf{W} - \mathbf{C}\|_F^2, \quad (7)$$

where $\mathcal{A}_{\mathbf{X}^{in}} : \mathbb{T}_w \rightarrow \mathbb{T}_{out}$, and $\mathbf{B} \in \mathbb{T}_{out}$ and $\mathbf{C} \in \mathbb{T}_w$ are tensors. The minimizer to (7) can be found by taking a derivative and setting it to zero, i.e.,

$$\mathcal{A}_{\mathbf{X}^{in}}^* (\mathcal{A}_{\mathbf{X}^{in}} (\mathbf{W}) - \mathbf{B}) + \mathbf{W} - \mathbf{C} = \mathbf{0}, \quad (8)$$

or

$$\mathcal{A}_{\mathbf{X}^{in}}^* (\mathcal{A}_{\mathbf{X}^{in}} (\mathbf{W})) + \mathbf{W} = \mathcal{A}_{\mathbf{X}^{in}}^* (\mathbf{B}) + \mathbf{C}. \quad (9)$$

Solving (9) for \mathbf{W} is efficiently possible via the method of conjugate gradient. The following algorithm outlines the process of solving (9), which is a variant of the original CG algorithm (e.g., see [CGL]) modified to address (9).

Algorithm 2: Least Squares Update in the Net-Trim Using Conjugate Gradient

```

initialize  $\mathbf{W}_0 = \mathbf{0}$ ;  $\mathbf{R}_0 = \mathcal{A}_{\mathbf{X}^{in}}^* (\mathbf{B}) + \mathbf{C}$ ;  $\mathbf{P}_0 = \mathbf{R}_0$ 
for  $k = 1, \dots, K_{\max}$  do
     $\mathbf{T}_{k-1} = \mathcal{A}_{\mathbf{X}^{in}} (\mathbf{P}_{k-1})$ 
     $\alpha_k = \frac{\|\mathbf{R}_{k-1}\|_F^2}{\|\mathbf{T}_{k-1}\|_F^2 + \|\mathbf{P}_{k-1}\|_F^2}$ 
     $\mathbf{W}_k = \mathbf{W}_{k-1} + \alpha_k \mathbf{P}_{k-1}$ 
     $\mathbf{R}_k = \mathbf{R}_{k-1} - \alpha_k (\mathcal{A}_{\mathbf{X}^{in}}^* (\mathbf{T}_{k-1}) + \mathbf{P}_{k-1})$ 
     $\beta_k = \frac{\|\mathbf{R}_k\|_F^2}{\|\mathbf{R}_{k-1}\|_F^2}$ 
     $\mathbf{P}_k = \mathbf{R}_k + \beta_k \mathbf{P}_{k-1}$ 
end
Output:  $\mathbf{W}_{K_{\max}}$ 

```

References

- [AAR18] Alireza Aghasi, Afshin Abdi, and Justin Romberg. Fast convex pruning of deep neural networks. *arXiv preprint arXiv:1806.06457*, 2018.
- [BPC⁺11] S. Boyd, N. Parikh, E. Chu, B. Peleato, and J. Eckstein. Distributed optimization and statistical learning via the alternating direction method of multipliers. *Foundations and Trends® in Machine Learning*, 3(1):1–122, 2011.
- [CGL] Conjugate Gradient Algorithm, howpublished = <https://math.aalto.fi/opetus/inv/cgalgorithm.pdf>, note = Accessed: 2018-07-18.