



香港中文大學
The Chinese University of Hong Kong

CENG5030

Part 2-5: CNN Inaccurate Speedup-3 — Pruning

Bei Yu

(Latest update: March 26, 2019)

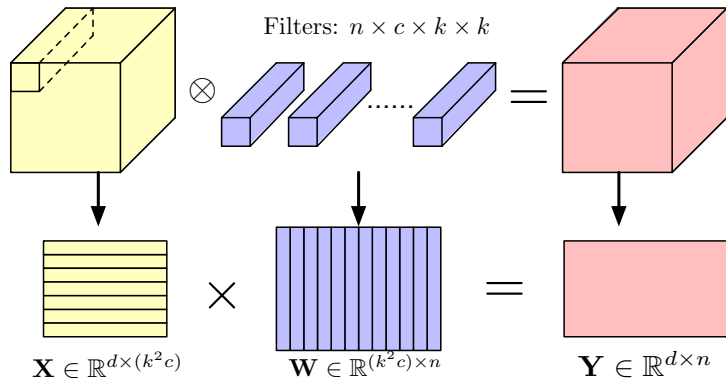
Spring 2019

These slides contain/adapt materials developed by

- ▶ Yuzhe Ma et al. (2018). “A Unified Approximation Framework for Non-Linear Deep Neural Networks”. In: *arXiv preprint arXiv:1807.10119*
- ▶ Yihui He, Xiangyu Zhang, and Jian Sun (2017). “Channel Pruning for Accelerating Very Deep Neural Networks”. In: *Proc. ICCV*



Im2col (Image2Column) Convolution



- ▶ Transform convolution to **matrix multiplication**
- ▶ **Unified** calculation for both convolution and fully-connected layers



Matrix Approximation or Matrix Regression?

$\mathbf{X} \in \mathbb{R}^{d \times (k^2 c)} \times \mathbf{W} \in \mathbb{R}^{(k^2 c) \times n} = \mathbf{Y} \in \mathbb{R}^{d \times n}$

- ▶ Matrix approximation: $\mathbf{W} \approx \mathbf{W}'$
- ▶ Matrix regression: $\mathbf{Y} = \mathbf{W} \cdot \mathbf{X} \approx \mathbf{W}' \cdot \mathbf{X}$



Compression Approach 1: Sparsity

$\mathbf{X} \in \mathbb{R}^{d \times (k^2 c)}$ \times $\mathbf{S} \in \mathbb{R}^{(k^2 c) \times n}$ $=$ $\mathbf{Y} \in \mathbb{R}^{d \times n}$

Sparse DNN

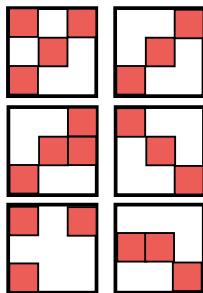
- ▶ *Sparsification*: weight pruning;
- ▶ *Compression*: compressed sparse format for storage;
- ▶ *Potential acceleration*: sparse matrix multiplication algorithm.

Exploring the Granularity of Sparsity that is Hardware-friendly

4 types of pruning granularity



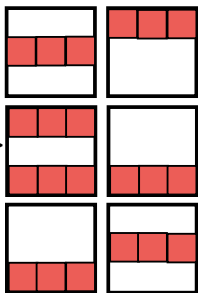
irregular sparsity



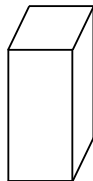
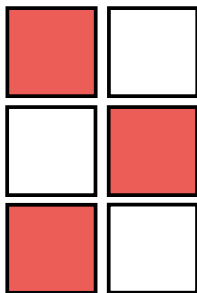
[Han et al, NIPS'15]



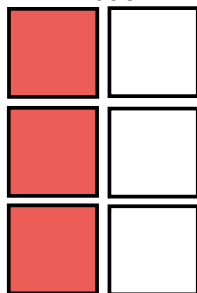
regular sparsity



more regular sparsity

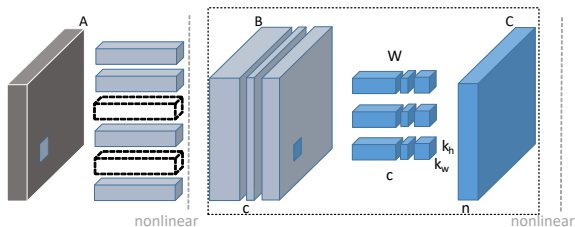


fully-dense model



[Molchanov et al, ICLR'17]

Our approach



We aim to reduce the width of feature map B, while minimizing the reconstruction error on feature map C. Our optimization algorithm performs within the dotted box, which does not involve nonlinearity. This figure illustrates the situation that two channels are pruned for feature map B. Thus corresponding channels of filters W can be removed. Furthermore, even though not directly optimized by our algorithm, the corresponding filters in the previous layer can also be removed (marked by dotted filters). c, n : number of channels for feature maps B and C, $k_h \times k_w$: kernel size.



Optimization

Formally, to prune a feature map with c channels, we consider applying $n \times c \times k_h \times k_w$ convolutional filters \mathbf{W} on $N \times c \times k_h \times k_w$ input volumes \mathbf{X} sampled from this feature map, which produces $N \times n$ output matrix \mathbf{Y} . Here, N is the number of samples, n is the number of output channels, and k_h, k_w are the kernel size. For simple representation, bias term is not included in our formulation. To prune the input channels from c to desired c' ($0 \leq c' \leq c$), while minimizing reconstruction error, we formulate our problem as follow:

$$\begin{aligned} \arg \min_{\beta, \mathbf{W}} \frac{1}{2N} \left\| \mathbf{Y} - \sum_{i=1}^c \beta_i \mathbf{X}_i \mathbf{W}_i^T \right\|_F^2 \\ \text{subject to } \|\beta\|_0 \leq c' \end{aligned} \quad (1)$$

$\|\cdot\|_F$ is Frobenius norm. \mathbf{X}_i is $N \times k_h k_w$ matrix sliced from i th channel of input volumes \mathbf{X} , $i = 1, \dots, c$. \mathbf{W}_i is $n \times k_h k_w$ filter weights sliced from i th channel of \mathbf{W} . β is coefficient vector of length c for channel selection, and β_i is i th entry of β . Notice that, if $\beta_i = 0$, \mathbf{X}_i will be no longer useful, which could be safely pruned from feature map. \mathbf{W}_i could also be removed.



Solving this ℓ_0 minimization problem in Eqn. 1 is NP-hard. we relax the ℓ_0 to ℓ_1 regularization:

$$\begin{aligned} \arg \min_{\beta, W} \frac{1}{2N} \left\| Y - \sum_{i=1}^c \beta_i X_i W_i^\top \right\|_F^2 + \lambda \|\beta\|_1 \\ \text{subject to } \|\beta\|_0 \leq c', \forall i \|W_i\|_F = 1 \end{aligned} \quad (2)$$

λ is a penalty coefficient. By increasing λ , there will be more zero terms in β and one can get higher speed-up ratio. We also add a constrain $\forall i \|W_i\|_F = 1$ to this formulation, which avoids trivial solution. Now we solve this problem in two folds. First, we fix W , solve β for channel selection. Second, we fix β , solve W to reconstruct error.



Optimization

(i) **The subproblem of β :** In this case, \mathbf{W} is fixed. We solve β for channel selection.

$$\hat{\beta}^{LASSO}(\lambda) = \arg \min_{\beta} \frac{1}{2N} \left\| \mathbf{Y} - \sum_{i=1}^c \beta_i \mathbf{Z}_i \right\|_F^2 + \lambda \|\beta\|_1 \quad (3)$$

subject to $\|\beta\|_0 \leq c'$

Here $\mathbf{Z}_i = \mathbf{X}_i \mathbf{W}_i^\top$ (size $N \times n$). We will ignore i th channels if $\beta_i = 0$.

(ii) **The subproblem of \mathbf{W} :** In this case, β is fixed. We utilize the selected channels to minimize reconstruction error. We can find optimized solution by least squares:

$$\arg \min_{\mathbf{W}'} \left\| \mathbf{Y} - \mathbf{X}' (\mathbf{W}')^\top \right\|_F^2 \quad (4)$$

Here $\mathbf{X}' = [\beta_1 \mathbf{X}_1 \ \beta_2 \mathbf{X}_2 \ \dots \ \beta_i \mathbf{X}_i \ \dots \ \beta_c \mathbf{X}_c]$ (size $N \times ck_h k_w$). \mathbf{W}' is $n \times ck_h k_w$ reshaped \mathbf{W} , $\mathbf{W}' = [\mathbf{W}_1 \ \mathbf{W}_2 \ \dots \ \mathbf{W}_i \ \dots \ \mathbf{W}_c]$. After obtained result \mathbf{W}' , it is reshaped back to \mathbf{W} . Then we assign $\beta_i \leftarrow \beta_i \|\mathbf{W}_i\|_F$, $\mathbf{W}_i \leftarrow \mathbf{W}_i / \|\mathbf{W}_i\|_F$. Constrain $\forall i \|\mathbf{W}_i\|_F = 1$ satisfies.



Compression Approach 2: Low-Rank

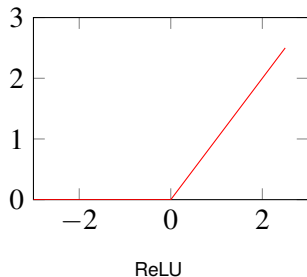
The diagram shows a sequence of matrix operations. On the left is a yellow matrix \mathbf{X} with 5 rows and 4 columns. This is followed by a multiplication symbol \times and a blue matrix \mathbf{U} with 5 rows and 2 columns. Another multiplication symbol \times follows, then a purple matrix \mathbf{V} with 2 rows and 4 columns. An equals sign $=$ leads to a red matrix \mathbf{Y} with 5 rows and 4 columns.

$$\mathbf{X} \in \mathbb{R}^{d \times (k^2 c)} \quad \mathbf{U} \in \mathbb{R}^{(k^2 c) \times r} \quad \mathbf{V} \in \mathbb{R}^{r \times n} \quad \mathbf{Y} \in \mathbb{R}^{d \times n}$$

Low-rank DNN

- ▶ *Low-rank approximation*: matrix decomposition or tensor decomposition.
- ▶ *Compression and acceleration*: less storage required and less FLOP in computation.

Non-linearity Approximation



- ▶ Activation unit: ReLU
- ▶ Error more sensitive to positive response;
- ▶ Enlarge the solution space.

$$\min_{\mathbf{W}} \sum_{i=1}^N \|\mathbf{W}\mathbf{X}_i - \mathbf{Y}_i\|_F \rightarrow \min_{\mathbf{W}} \sum_{i=1}^N \|r(\mathbf{W}\mathbf{X}_i) - \mathbf{Y}_i\|_F$$

- ▶ \mathbf{X} : input feature map
- ▶ \mathbf{Y} : output feature map

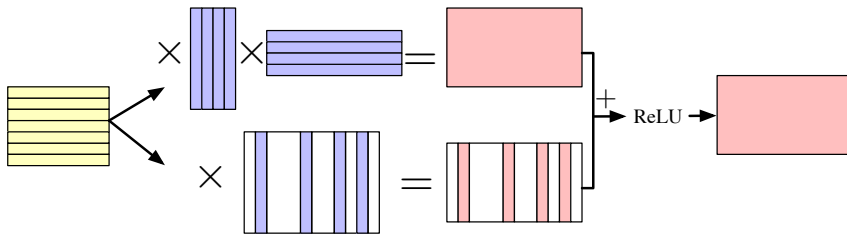


Reading List

- ▶ Xiangyu Zhang et al. (2015). “Efficient and accurate approximations of nonlinear convolutional networks”. In: *Proc. CVPR*, pp. 1984–1992
- ▶ Hao Zhou, Jose M Alvarez, and Fatih Porikli (2016). “Less is more: Towards compact cnns”. In: *Proc. ECCV*, pp. 662–677
- ▶ Yihui He, Xiangyu Zhang, and Jian Sun (2017). “Channel Pruning for Accelerating Very Deep Neural Networks”. In: *Proc. ICCV*
- ▶ Xiyu Yu et al. (2017). “On compressing deep models by low rank and sparse decomposition”. In: *Proc. CVPR*, pp. 7370–7379



Proposed Unified Structure



- ▶ Simultaneous low-rank approximation and network sparsification;
- ▶ Non-linearity is taken into account.
- ▶ Acceleration is achieved with structured sparsity.

Formulation

Given a pre-trained network, the goal is to minimize the reconstruction error of the response in each layer after activation, using sparse component and low-rank component.

$$\begin{aligned} \min_{\mathbf{A}, \mathbf{B}} \quad & \sum_{i=1}^N \|\mathbf{Y}_i - r((\mathbf{A} + \mathbf{B})\mathbf{X}_i)\|_F, \\ \text{s.t.} \quad & \|\mathbf{A}\|_0 \leq S, \\ & \text{rank}(\mathbf{B}) \leq L. \end{aligned}$$

- ▶ \mathbf{X} : input feature map
- ▶ \mathbf{Y} : output feature map

Not easy to solve: l_0 minimization and rank minimization are NP-hard.



Relaxation

$$\min_{\mathbf{A}, \mathbf{B}} \sum_{i=1}^N \|\mathbf{Y}_i - r((\mathbf{A} + \mathbf{B})\mathbf{X}_i)\|_F^2 + \lambda_1 \|\mathbf{A}\|_{2,1} + \lambda_2 \|\mathbf{B}\|_*$$

- ▶ The l_0 constraint is relaxed by $l_{2,1}$ norm such that the zero elements in \mathbf{A} appear column-wise;
- ▶ The rank constraint on \mathbf{B} is relaxed by nuclear norm of \mathbf{B} , which is the sum of the singular values;
- ▶ Apply alternating direction method of multipliers (ADMM) to solve it;



Alternating Direction Method of Multipliers (ADMM)

Reformulating the problem with an auxiliary variable \mathbf{M} ,

$$\begin{aligned} \min_{\mathbf{A}, \mathbf{B}, \mathbf{M}} \quad & \sum_{i=1}^N \|\mathbf{Y}_i - r(\mathbf{M}\mathbf{X}_i)\|_F^2 + \lambda_1 \|\mathbf{A}\|_{2,1} + \lambda_2 \|\mathbf{B}\|_*, \\ \text{s.t.} \quad & \mathbf{A} + \mathbf{B} = \mathbf{M}. \end{aligned}$$

Then the augmented Lagrangian function is

$$\begin{aligned} & L_t(\mathbf{A}, \mathbf{B}, \mathbf{M}, \mathbf{\Lambda}) \\ = & \sum_{i=1}^N \|\mathbf{Y}_i - r(\mathbf{M}\mathbf{X}_i)\|_F^2 + \lambda_1 \|\mathbf{A}\|_{2,1} + \lambda_2 \|\mathbf{B}\|_* + \langle \mathbf{\Lambda}, \mathbf{A} + \mathbf{B} - \mathbf{M} \rangle + \frac{t}{2} \|\mathbf{A} + \mathbf{B} - \mathbf{M}\|_F^2, \end{aligned}$$



Alternating Direction Method of Multipliers (ADMM)

Iteratively solve with following rules. All of them can be solved efficiently.

$$\left\{ \begin{array}{l} \mathbf{A}_{k+1} = \underset{\mathbf{A}}{\operatorname{argmin}} \lambda_1 \|\mathbf{A}\|_{2,1} + \frac{t}{2} \left\| \mathbf{A} + \mathbf{B}_k - \mathbf{M}_k + \frac{\boldsymbol{\Lambda}_k}{t} \right\|_F^2, \\ \mathbf{B}_{k+1} = \underset{\mathbf{B}}{\operatorname{argmin}} \lambda_2 \|\mathbf{B}\|_* + \frac{t}{2} \left\| \mathbf{B} + \mathbf{A}_{k+1} - \mathbf{M}_k + \frac{\boldsymbol{\Lambda}_k}{t} \right\|_F^2, \\ \mathbf{M}_{k+1} = \underset{\mathbf{M}}{\operatorname{argmin}} \sum_{i=1}^N \|\mathbf{Y}_i - r(\mathbf{M}\mathbf{X}_i)\|_F^2 + \langle \boldsymbol{\Lambda}_k, \mathbf{A}_{k+1} + \mathbf{B}_{k+1} - \mathbf{M} \rangle + \frac{t}{2} \|\mathbf{A}_{k+1} + \mathbf{B}_{k+1} - \mathbf{M}\|_F^2, \\ \boldsymbol{\Lambda}_{k+1} = \boldsymbol{\Lambda}_k + t(\mathbf{A}_{k+1} + \mathbf{B}_{k+1} - \mathbf{M}_{k+1}). \end{array} \right.$$



Solving $l_{2,1}$ -norm

$$\min_{\mathbf{A}} \lambda_1 \|\mathbf{A}\|_{2,1} + \frac{t}{2} \left\| \mathbf{A} + \mathbf{B}_k - \mathbf{M}_k + \frac{\mathbf{\Lambda}_k}{t} \right\|_F^2$$

Closed Form Update Rule¹

$$\begin{aligned} \mathbf{A}_{k+1} &= \text{prox}_{\frac{\lambda_1}{t} \|\cdot\|_{2,1}} \left(\mathbf{M}_k - \mathbf{B}_k - \frac{\mathbf{\Lambda}_k}{t} \right), \\ \mathbf{C} &= \mathbf{M}_k - \mathbf{B}_k - \frac{\mathbf{\Lambda}_k}{t}, \\ [\mathbf{A}_{k+1}]_{:,i} &= \begin{cases} \frac{\|[\mathbf{C}]_{:,i}\|_2 - \frac{\lambda_1}{t}}{\|[\mathbf{C}]_{:,i}\|_2} [\mathbf{C}]_{:,i}, & \text{if } \|[\mathbf{C}]_{:,i}\|_2 > \frac{\lambda_1}{t}; \\ 0, & \text{otherwise.} \end{cases} \end{aligned}$$

¹G. Liu et al., "Robust recovery of subspace structures by low-rank representation", *TPAMI*, 2013.



Solving nuclear norm

$$\min_{\mathbf{B}} \lambda_2 \|\mathbf{B}\|_* + \frac{t}{2} \left\| \mathbf{B} + \mathbf{A}_{k+1} - \mathbf{M}_k + \frac{\mathbf{\Lambda}_k}{t} \right\|_F^2$$

Closed Form Update Rule²

$$\mathbf{B}_{k+1} = \text{prox}_{\frac{\lambda_2}{t} \|\cdot\|_*} \left(\mathbf{M}_k - \mathbf{A}_{k+1} - \frac{\mathbf{\Lambda}_k}{t} \right),$$

$$\mathbf{D} = \mathbf{M}_k - \mathbf{A}_{k+1} - \frac{\mathbf{\Lambda}_k}{t},$$

$$\mathbf{B}_{k+1} = \mathbf{U} \mathbf{D}_{\frac{\lambda_2}{t}}(\mathbf{\Sigma}) \mathbf{V}, \quad \text{where } \mathbf{D}_{\frac{\lambda_2}{t}}(\mathbf{\Sigma}) = \text{diag}(\{(\sigma_i - \frac{\lambda_2}{t})_+\}).$$

²J-F. Cai et al., "A singular value thresholding algorithm for matrix completion", *SIOPT*, 2010.



Comparison on *CIFAR-10* dataset

Model	Method	Accuracy ↓	CR	Speed-up
VGG-16	Original	0.00%	1.00	1.00
	ICLR'17 ³	0.06%	2.70	1.80
	Ours	0.40%	4.44	2.20
NIN	Original	0.00%	1.00	1.00
	ICLR'16 ⁴	1.43%	1.54	1.50
	IJCAI'18 ⁵	1.43%	1.45	-
	Ours	0.41%	2.77	1.70

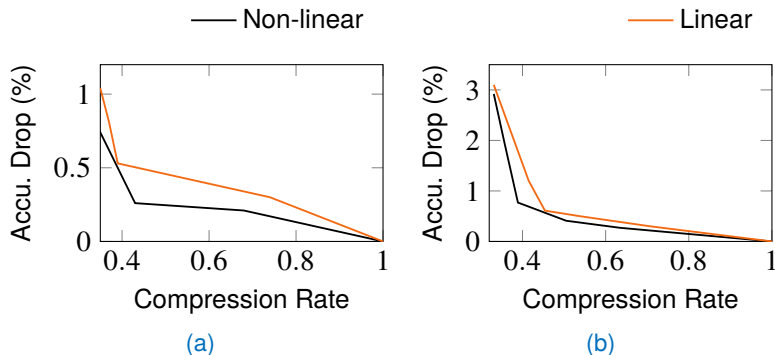
³Hao Li et al. (2017). “Pruning filters for efficient convnets”. In: *Proc. ICLR*.

⁴Cheng Tai et al. (2016). “Convolutional neural networks with low-rank regularization”. In: *Proc. ICLR*.

⁵Shiva Prasad Kasiviswanathan, Nina Narodytska, and Hongxia Jin (2018). “Network Approximation using Tensor Sketching”. In: *Proc. IJCAI*, pp. 2319–2325.



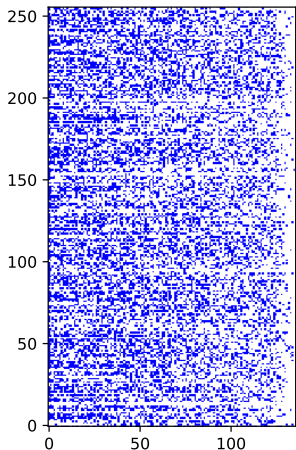
Preliminary Results



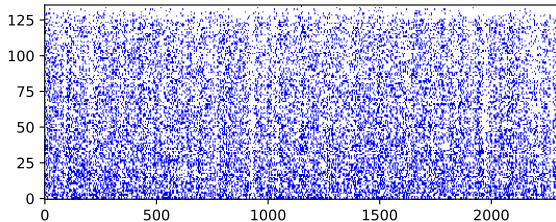
Comparison of reconstructing linear response and non-linear response: (a) layer conv2-1; (b) layer conv3-1.



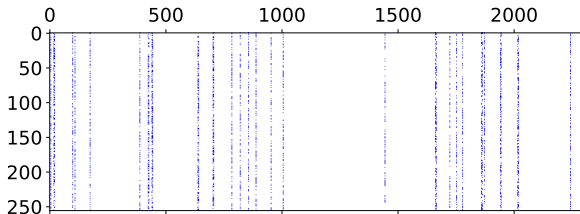
Approximation Example



(a)



(b)



(c)

Approximated filters of `conv3-1`. Blue dots have non-zero values. Low-rank filter B with rank 136 is decomposed into UV , both of which have rank 136. (a) Matrix U ; (b) Matrix V . (c) Column-wise sparse filter A .



Comparison on *ImageNet* dataset

Model	Method	Top-5 Accu.↓	CR	Speed-up
AlexNet	Original	0.00%	1.00	1.00
	ICLR'16 ⁶	0.37%	5.00	1.82
	ICLR'16 ⁷	1.70%	5.46	1.81
	CVPR'18 ⁸	1.43%	1.50	-
	Ours	1.27%	5.56	1.10
GoogleNet	Original	0.00%	1.00	1.00
	ICLR'16 ¹²	0.42%	2.84	1.20
	ICLR'16 ¹³	0.24%	1.28	1.23
	CVPR'18 ¹⁴	0.21%	1.50	-
	Ours	0.00%	2.87	1.35

⁶Cheng Tai et al. (2016). "Convolutional neural networks with low-rank regularization". In: *Proc. ICLR*.

⁷Yong-Deok Kim et al. (2016). "Compression of deep convolutional neural networks for fast and low power mobile applications". In: *Proc. ICLR*.

⁸Ruichi Yu et al. (2018). "NISP: Pruning networks using neuron importance score propagation". In: *Proc. CVPR*.

