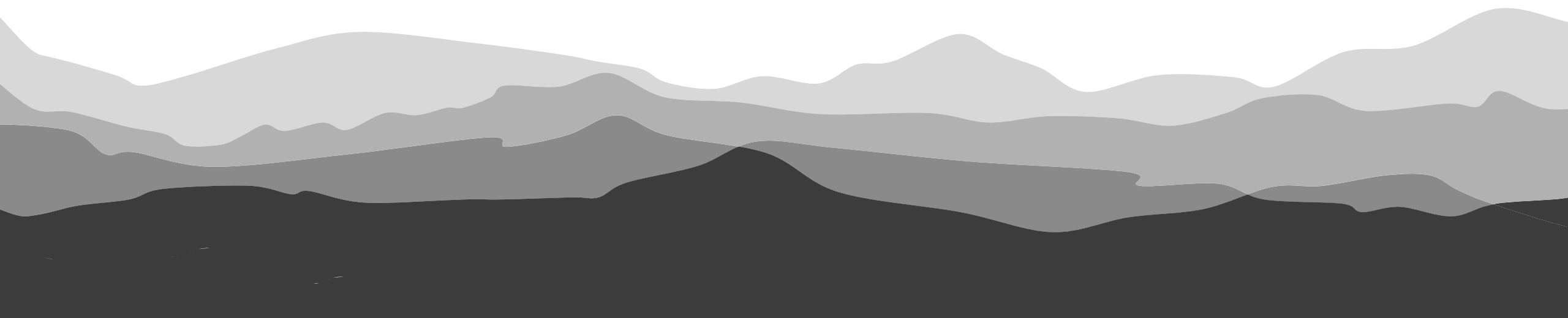


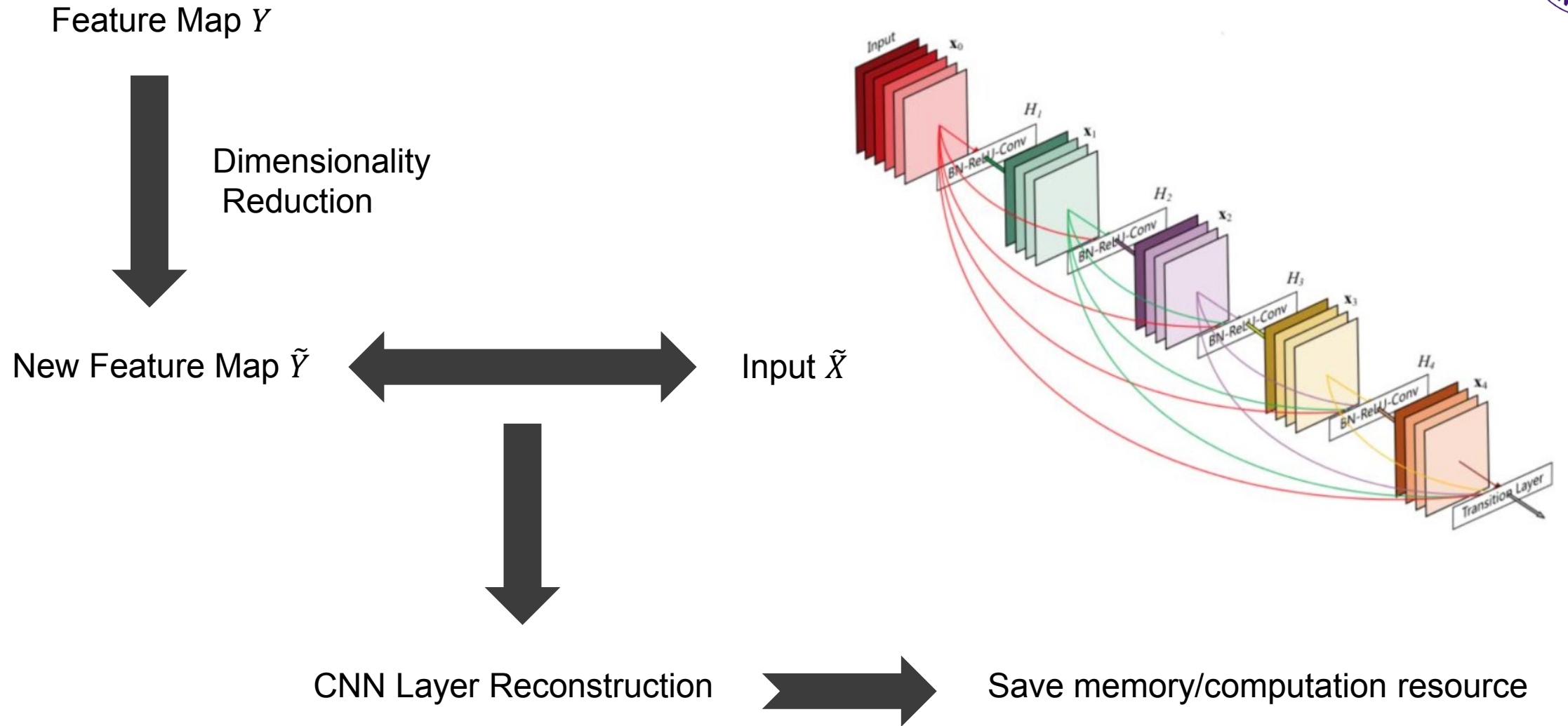


Beyond Filters: Compact Feature Map for Portable Deep Model

ICML2017



Framework



Motivation

$$X \in \mathbb{R}^{H \times W \times c} \xrightarrow{F \in \mathbb{R}^{s_1 \times s_2 \times c \times d}} Y \in \mathbb{R}^{H' \times W' \times d}$$

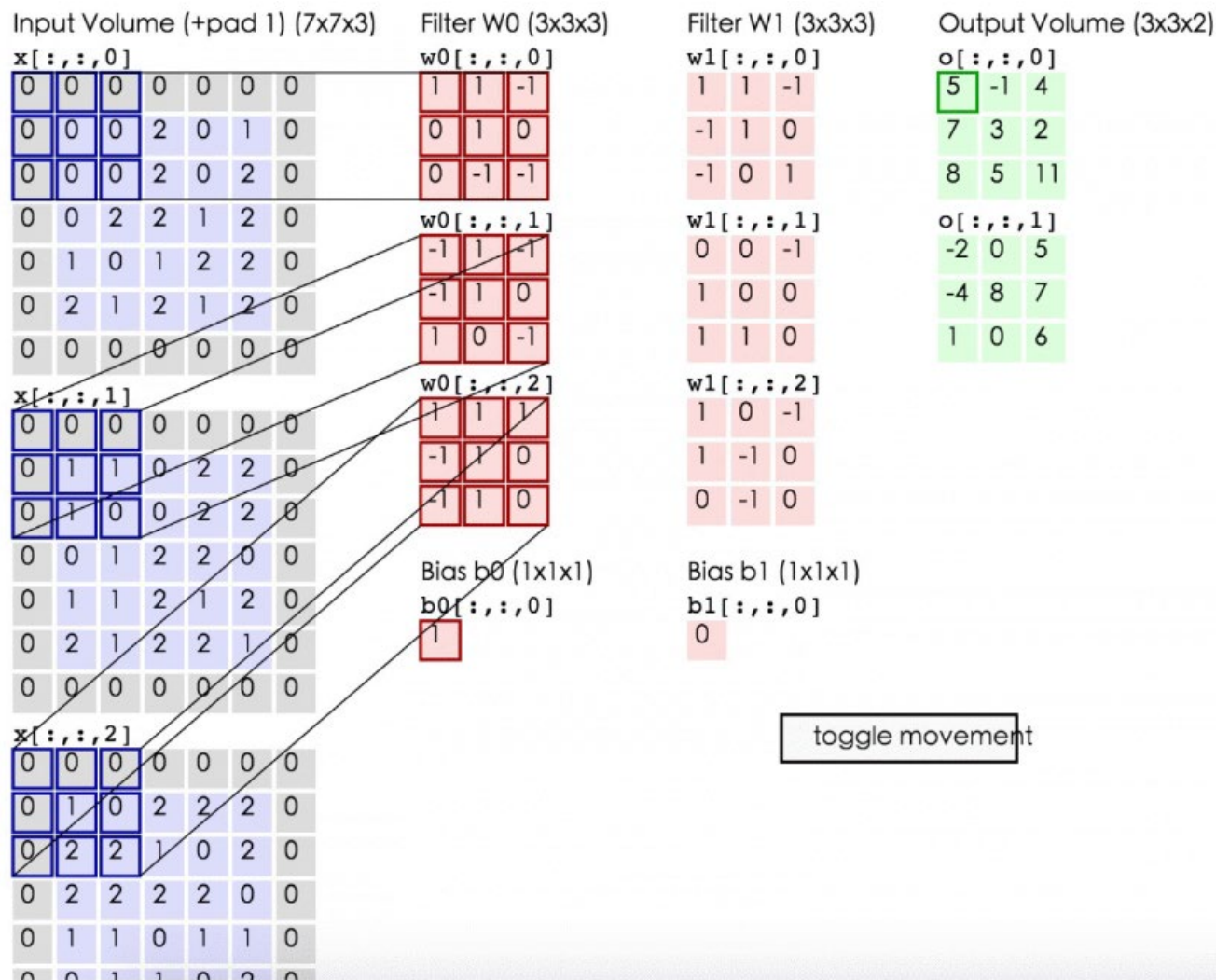
d is the number of convolution filters in this layer and d is much larger than $s = s_1 \times s_2$.

1	1	1	0	0
0	1	1	1	0
0	0	1	1	1
0	0	1	1	0
0	1	1	0	0

Image

4		

Convolved Feature



toggle movement

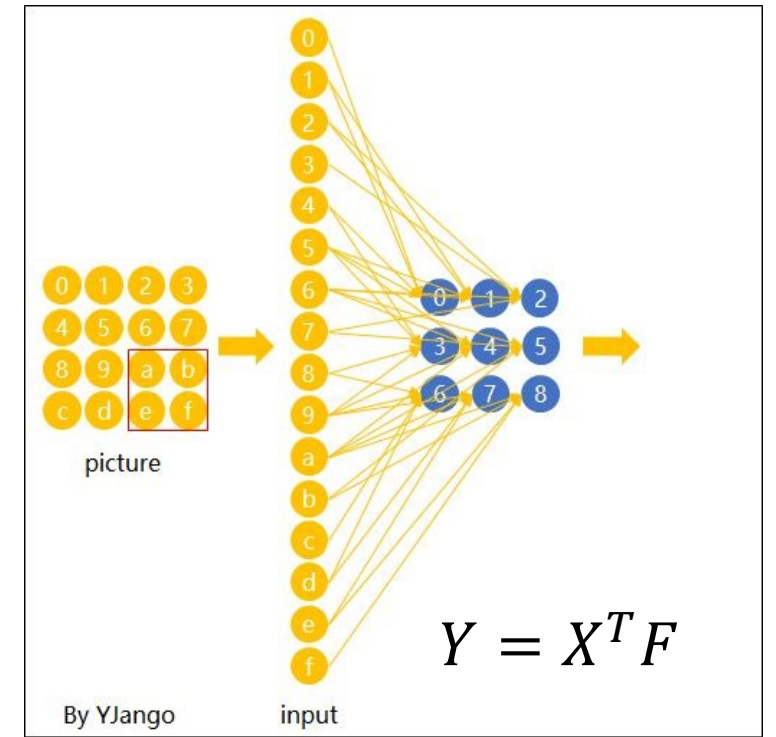
$$X \in \mathbb{R}^{H \times W \times c}$$

$$X = [\text{vec}(X_1), \dots, \text{vec}(X_c)] \in \mathbb{R}^{HW \times c}$$

$$Y \in \mathbb{R}^{H' \times W' \times d}$$

$$Y = [\text{vec}(Y_1), \dots, \text{vec}(Y_d)] \in \mathbb{R}^{H'W' \times d}$$

Model



Feature maps from different filters should be independent of each other as far as possible

$$\theta(Y) = \|Y^T Y \circ (1 - I)\|_F^2$$

$$\tilde{Y} = \phi(Y) = YP^T \quad P \in \mathbb{R}^{d \times d} \text{ is the projection matrix}$$



$$\min_{P, C} \|PY^TYP^T - C\|_F^2, s. t. C = \text{diag}(c)$$



An optimal projection should thus not only remove redundancy between feature maps, but also preserve the discriminability of these features.

$$\min_{P, C} \|PY^TYP^T - C\|_F^2 + \lambda \|D(\tilde{Y}) - D(Y)\|$$
$$s. t. \tilde{Y} = YP^T, C = \text{diag}(c)$$

D_{ij} is the Euclidean distance between the i -th column and the j -th column of Y

Distances between feature maps can be easily preserved if P is orthogonal, i.e. $PP^T = I$

$$\|y_1P^T\|_2 = \|y_1\|_2 \quad \|y_1P^T - y_2P^T\|_2^2 = \|y_1 - y_2\|_2^2$$



$$\min_{P,c} \|PY^TYP^T - C\|_F^2 \quad s.t. \quad C = \text{diag}(c), PP^T = I$$



The orthogonal transformation P learned is able to extract the intrinsic representation and preserve distances between feature maps, but the dimensionality has not been reduced since P is a square matrix.

Propose to use a sparse matrix to approximate the representation generated by P .

$$\min_{\tilde{Y}} \|\tilde{Y} - YP^T\|_F^2 + \lambda \|\tilde{Y}\|_{2,1} \quad \|\tilde{Y}\|_{2,1} = \sum_i \|\tilde{Y}^i\| \quad \tilde{Y}^i \text{ is the } i\text{-th column in } \tilde{Y}$$

$$\tilde{Y}^i = \begin{cases} \frac{\|u_i\| - \lambda}{\|u_i\|} u_i & \text{if } \lambda < \|u_i\| \\ 0 & \text{otherwise} \end{cases} \quad u_i = YP_i^T$$

Robust Subspace Segmentation by Low-Rank Representation (ICML2010)



$$\min_{P,c} \|PY^TYP^T - C\|_F^2 + \beta \|c\|_1$$

$$s.t. C = \text{diag}(c), PP^T = I$$

Some small valued columns in $\tilde{Y} = YP^T$ will be discarded.

$$\min_{P,c} \|PY^TYP^T - C\|_F^2 + \alpha \|PP^T - I\|_F^2 + \beta \|c\|_1$$

$$s.t. P = \text{circ}(p), C = \text{diag}(c)$$

$$p = (p_1, \dots, p_d)^T$$

$$\text{circ}(p) := \begin{bmatrix} p_1 & p_d & \dots & p_3 & p_2 \\ p_2 & p_1 & p_d & & p_3 \\ \vdots & p_2 & p_1 & \ddots & \vdots \\ p_{d-1} & & \ddots & \ddots & p_d \\ p_d & p_{d-1} & \dots & p_2 & p_1 \end{bmatrix}$$

..., circulant matrices have complex internal structures and strong diversity thus can be utilized for approximating huge matrices

An exploration of parameter redundancy in deep networks with circulant projections(CVPR2015)



$$\min_{P,c} \|PY^TYP^T - C\|_F^2 + \alpha \|PP^T - I\|_F^2 + \beta \|c\|_1$$
$$s.t. P = \text{circ}(p), C = \text{diag}(c)$$

Solve c for a fixed projection P

$$\min_{P,c} \|PY^TYP^T - C\|_F^2 + \beta \|c\|_1$$
$$s.t. C = \text{diag}(c)$$

$$\tilde{c} = \text{diag}(PY^TYP^T)$$

$$\min_c \|\tilde{c} - c\|_2^2 + \beta \|c\|_1 \quad c = \text{sign}(\tilde{c}) \circ \max\{|\tilde{c}| - \beta/2, 0\}$$

机器学习11.4

Solve P for a fixed projection c

循环矩阵

$$\mathbf{x} = [x_0, x_1, x_2, \dots, x_{N-1}]$$

$$\mathbf{X} = C(\mathbf{x}) = \mathbf{F} \cdot \text{diag}(\mathcal{F}(\mathbf{x})) \cdot \mathbf{F}^H$$

$$\mathbf{X} = C(\mathbf{x}) = \begin{bmatrix} x_0 & x_1 & x_2 & \cdots & x_{N-1} \\ x_{N-1} & x_0 & x_1 & \cdots & x_{N-2} \\ x_{N-2} & x_{N-1} & x_0 & \cdots & x_{N-3} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ x_1 & x_2 & x_3 & \cdots & x_0 \end{bmatrix}$$

$$\hat{x} = \mathcal{F}(x) = Fx$$

$$w = e^{-j2\pi/N}$$

$$\mathbf{F}^H \mathbf{F} = \mathbf{F} \mathbf{F}^H = \mathbf{I}$$

$$\mathbf{F} = \frac{1}{\sqrt{n}} \begin{bmatrix} 1 & 1 & 1 & \cdots & 1 \\ 1 & w & w^2 & \cdots & w^{N-1} \\ 1 & w^2 & w^4 & \cdots & w^{2(N-1)} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & w^{N-1} & w^{2(N-1)} & \cdots & w^{(N-1)^2} \end{bmatrix}$$

循环矩阵

$$\mathbf{A} + \mathbf{B} = C(\mathbf{a} + \mathbf{b})$$

$$\mathbf{A} \cdot \mathbf{B} = C(\mathcal{F}^{-1}(\mathcal{F}(\mathbf{a}) \odot \mathcal{F}(\mathbf{b})))$$

$$\mathbf{X} = C(\mathbf{x}) = \begin{bmatrix} x_0 & x_1 & x_2 & \cdots & x_{N-1} \\ x_{N-1} & x_0 & x_1 & \cdots & x_{N-2} \\ x_{N-2} & x_{N-1} & x_0 & \cdots & x_{N-3} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ x_1 & x_2 & x_3 & \cdots & x_0 \end{bmatrix}$$

$$\mathbf{X}^T = \mathbf{F} \cdot \text{diag}(\mathcal{F}^*(\mathbf{x})) \cdot \mathbf{F}^H$$

$$\mathbf{X}^{-1} = \mathbf{F} \cdot \text{diag}(\mathcal{F}(\mathbf{x}))^{-1} \cdot \mathbf{F}^H$$

$$\mathcal{F}(\mathbf{X}\mathbf{y}) = \mathcal{F}^*(\mathbf{x}) \odot \mathcal{F}(\mathbf{y})$$



$$\min_P \|PY^TYP^T - C\|_F^2 + \alpha \|PP^T - I\|_F^2$$

$$s.t. P = \text{circ}(p)$$

Since P and its relevant terms can be efficiently calculated in the frequency domain, transfer the problem into the Fourier domain for having a higher training speed.

$$\min_P L(P) = L_1(P) + \alpha L_2(P), s.t. P = \text{circ}(p)$$

矩阵分析与应用 (第2版)
3. 矩阵微分

$$\frac{\partial L_1(P)}{\partial P} = 4HPPT^T HP - 4HPC = \frac{4}{d^2} HS^H \hat{P} \hat{P}^T SHS^H \hat{P} S - \frac{4}{d} HS^H \hat{P} S^H C \quad \hat{P} = \text{diag}(\hat{p})$$

$$\frac{\partial L_1(P)}{\partial \hat{P}} = \frac{1}{d} S \frac{\partial L_1(P)}{\partial P} S^H = \frac{4}{d^2} (\hat{H} \hat{P} \hat{P}^H \hat{H} \hat{P} - \hat{H} \hat{P} \hat{C}) \quad \longrightarrow \quad \frac{\partial L_1(P)}{\partial \hat{p}} = \text{diag} \left(\frac{\partial L_1(P)}{\partial \hat{P}} \right)$$

$$\hat{H} = SHS^H \quad \hat{C} = SCS^H$$

$$\min_{P,c} \|PY^TYP^T - C\|_F^2 + \alpha \|PP^T - I\|_F^2 + \beta \|c\|_1$$

$$s. t. P = \text{circ}(p), C = \text{diag}(c)$$

$$L_2(P) = \|PP^T - I\|_F^2 = \left\| \frac{1}{d} S^H \hat{P}^H \hat{P} S - \frac{1}{d} S^H I S \right\|_F^2$$

$$= \text{Tr} \left[\frac{1}{d} S^H (\hat{P}^H \hat{P} - I)^H (\hat{P}^H \hat{P} - I) S \right]$$

$$= \text{Tr} \left[(\hat{P}^H \hat{P} - I)^H (\hat{P}^H \hat{P} - I) \right] = \|\hat{p}^H \circ \hat{p} - 1\|_2^2$$

$$\frac{\partial L_2(P)}{\partial \hat{p}} = 4(\hat{p} \hat{p}^H \circ I) \hat{p} - 4\hat{p}$$

$$\hat{p} = \hat{p} - \eta \left(\frac{\partial L_1(P)}{\partial \hat{p}} + \alpha \frac{\partial L_2(P)}{\partial \hat{p}} \right)$$

$$P = \frac{1}{d} \text{diag}(\bar{c}) S^H \hat{P} S$$

$$\bar{c}_i = 0, \text{ if } c_i = 0 \text{ and } \bar{c}_i = 1 \text{ otherwise}$$

$$\tilde{Y} = Y P_{\tilde{d}}^T = X^T F P_{\tilde{d}}^T = X^T \tilde{F}$$

CNN Layer Reconstruction

$$\min_{\tilde{F}} \|\tilde{Y} - \tilde{X}^T \tilde{F}\|_F^2 + \gamma \|\tilde{F}\|_F^2$$

$$\tilde{F} = (\tilde{X}\tilde{X}^T + \gamma I)^{-1} \tilde{X}\tilde{Y}$$

上一层的一个输出 map Mj

1	1	1	1
0	0	1	1
0	1	1	0
0	1	1	0

卷积核 Kj

1	1
0	1

sigmoid (sum(Mij*Kij)+bias)

=

卷积层的输出 map M2j

2	3	3
1	2	2
2	3	1

$$L(\tilde{F}) = \text{Tr}(\tilde{F}^T \tilde{X}^T \tilde{X} \tilde{F}) - 2\text{Tr}(\tilde{F}^T \tilde{X}^T \tilde{Y}) + \gamma \text{Tr}(\tilde{F}^T \tilde{F})$$

$$\tilde{F} = \tilde{F} - \eta \frac{\partial L(\tilde{F})}{\partial \tilde{F}}$$

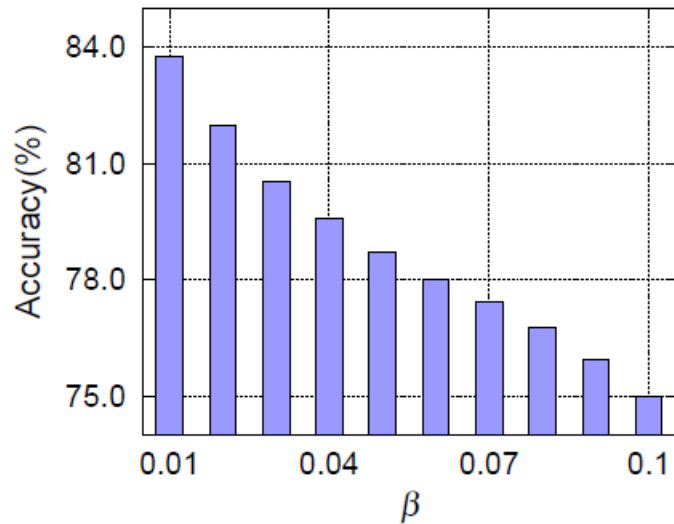
Experiments

Speed-up ratio

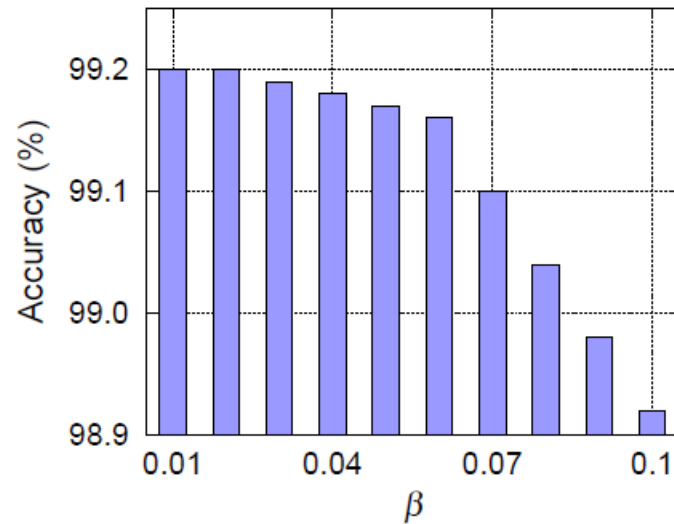
$$r_s = \frac{s_i^2 d_{i-1} d_i H'_i W'_i}{s_i^2 \tilde{d}_{i-1} \tilde{d}_i H'_i W'_i} = \frac{d_{i-1} d_i}{\tilde{d}_{i-1} \tilde{d}_i}$$

Compression ratio

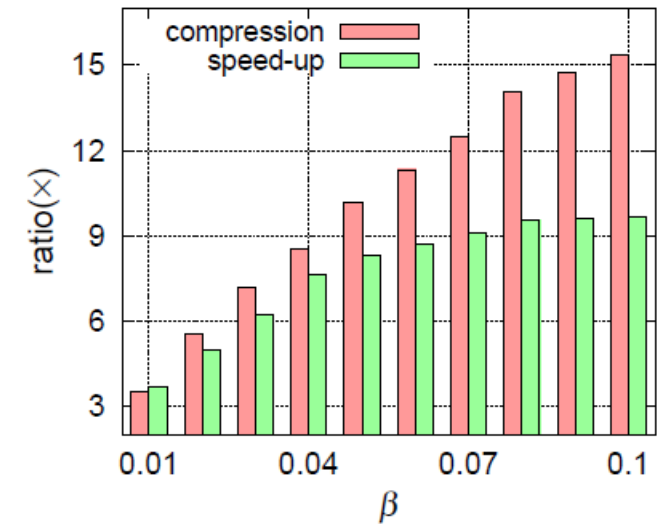
$$r_c = \frac{s_i^2 d_{i-1} d_i + d_i H'_i W'_i}{s_i^2 \tilde{d}_{i-1} \tilde{d}_i + \tilde{d}_i H'_i W'_i}$$



(a) Accuracy of the reconstructed model.



(b) Model accuracy after fine-tuning.



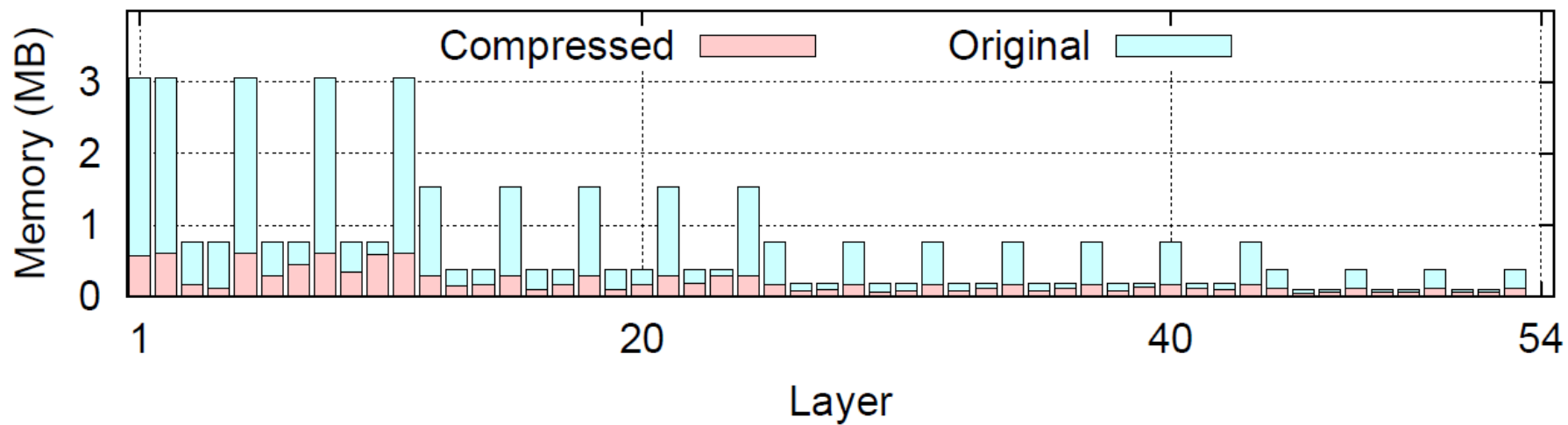
(c) Compression and speed-up ratios.

Table 1. Compression statistics for AlexNet.

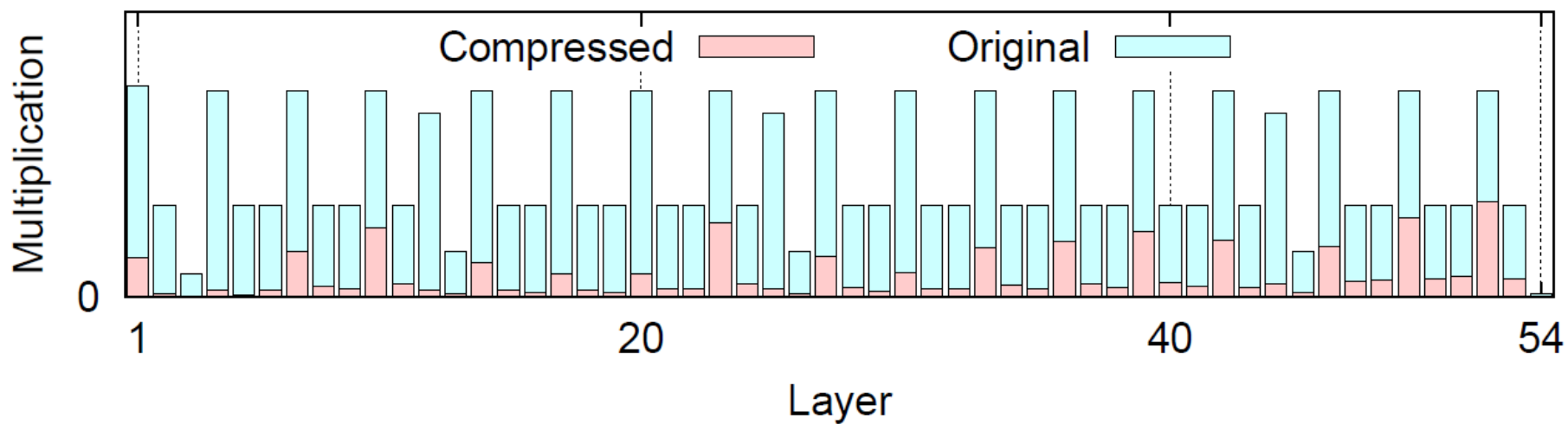
Layer	Memory	r_c	Mult.	r_s
conv1	1.24MB	$2.4\times$	1.05×10^8	$2.4\times$
conv2	1.88MB	$10.0\times$	2.23×10^8	$15.4\times$
conv3	3.62MB	$15.5\times$	1.49×10^8	$21.2\times$
conv4	2.78MB	$3.0\times$	1.12×10^8	$3.6\times$
conv5	1.85MB	$1.2\times$	0.74×10^8	$1.2\times$
fc6	144MB	$3.5\times$	0.37×10^8	$3.5\times$
fc7	64.02MB	$72.0\times$	0.16×10^8	$72.1\times$
fc8	15.63MB	$23.4\times$	0.04×10^8	$23.5\times$
Total	235.03MB	$5.06\times$	7.24×10^8	$4.31\times$

Table 2. Compression statistics for VGG-16 Net.

Layer	Memory	r_c	Mult.	r_s
conv1_1	12.26MB	$10.7\times$	0.11×10^9	$10.7\times$
conv1_2	12.39MB	$1.8\times$	2.41×10^9	$19.0\times$
conv2_1	6.41MB	$3.3\times$	1.20×10^9	$5.7\times$
conv2_2	6.69MB	$4.0\times$	2.41×10^9	$12.0\times$
conv3_1	4.19MB	$7.2\times$	1.20×10^9	$21.9\times$
conv3_2	5.31MB	$2.7\times$	2.41×10^9	$10.2\times$
conv3_3	5.31MB	$2.9\times$	2.41×10^9	$4.2\times$
conv4_1	6.03MB	$41.3\times$	1.20×10^9	$56.2\times$
conv4_2	10.53MB	$16.5\times$	2.41×10^9	$70.1\times$
conv4_3	10.53MB	$3.9\times$	2.41×10^9	$5.1\times$
conv5_1	9.38MB	$7.8\times$	0.60×10^9	$8.0\times$
conv5_2	9.38MB	$18.5\times$	0.60×10^9	$21.3\times$
conv5_3	9.38MB	$23.4\times$	0.60×10^9	$26.7\times$
fc6	392MB	$8.6\times$	0.41×10^9	$8.6\times$
fc7	64MB	$3.3\times$	0.16×10^8	$3.3\times$
fc8	15.62MB	$2.2\times$	0.41×10^7	$2.2\times$
Total	579.46MB	$6.19\times$	2.04×10^{10}	$9.63\times$



(b) Compression ratio of feature maps of convolutional layers (r_{c_2}).



(c) Speed-up ratio of all convolutional layers (r_s).



Input: A pre-trained convolutional neural network \mathcal{N} learned through a dataset \mathcal{X} with k convolutional layers: $\mathcal{L}_1, \dots, \mathcal{L}_k$, weight parameter γ , learning rate η .

1: Calculate feature maps of each layer by using the original network: $\{\mathbf{Y}_1, \dots, \mathbf{Y}_k\} \leftarrow \mathcal{N}(\mathcal{X})$;

2: **for** $i = 1$ to $k - 1$ **do**

3: Learn the projection P_i by solving Fcn. 9;

4: Calculate new feature maps: $\tilde{\mathbf{Y}}_i \leftarrow \mathbf{Y}_i P_i^T$;

5: **end for**

6: Keep feature maps of the k -th layer: $\tilde{\mathbf{Y}}_k \leftarrow \mathbf{Y}_k$;

7: Construct a new network $\tilde{\mathcal{N}}$ according to $\{\tilde{\mathbf{Y}}_1, \dots, \tilde{\mathbf{Y}}_k\}$ and initialize convolution filters $\{\tilde{\mathbf{F}}_1, \dots, \tilde{\mathbf{F}}_k\}$ by random values from the standard normal distribution;

8: **repeat**

9: Randomly select a batch \mathcal{X}_j from \mathcal{X} ;

10: **for** $i = 1$ to k **do**

11: Generate input data $\tilde{\mathbf{X}}_i$ of \mathcal{L}_i exploiting $\tilde{\mathcal{N}}$;

12: Estimate the new filter matrix (Fcn. 20):

$$\tilde{\mathbf{F}}_i \leftarrow \tilde{\mathbf{F}}_i - \eta \partial L(\tilde{\mathbf{F}}_i) / \partial \tilde{\mathbf{F}}_i;$$

13: Convert $\tilde{\mathbf{F}}_i$ into filter data and fill it in $\tilde{\mathcal{N}}$;

14: **end for**

15: **until** Convergence;

Output: The new convolutional neural network $\tilde{\mathcal{N}}$.

It is worth mentioning that **input data of the first layer of the original network N and feature map of the last layer** (closely related to the number of classification labels) are **kept unchanged**.

As for other intermediate convolutional layers and fully connected layers, we can **generate compact feature maps \tilde{Y}** from the original feature maps Y .

Then, calculate the input data \tilde{X} using the compressed network \tilde{N} and estimate the filter matrix \tilde{N} .