

# Filter Pruning via Geometric Median for Deep Convolutional Neural Networks Acceleration

He et al. CVPR 2019

---

Haein Choe

June 14, 2025

Pseudo Lab - On-Device AI: ON THE Air

# Table of contents

1. Introduction

2. Methodology

3. Experiments

4. Conclusion

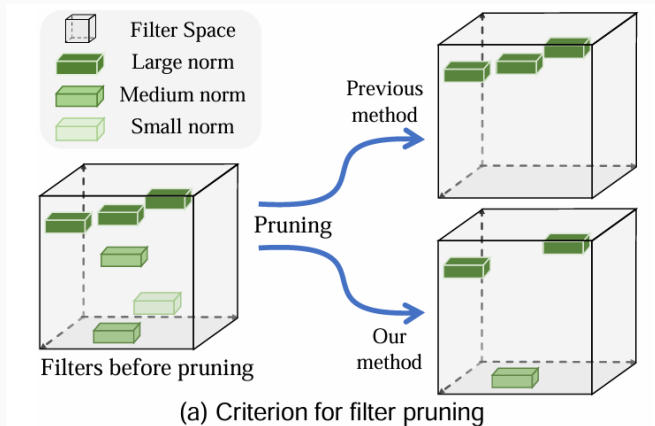
# Introduction

---

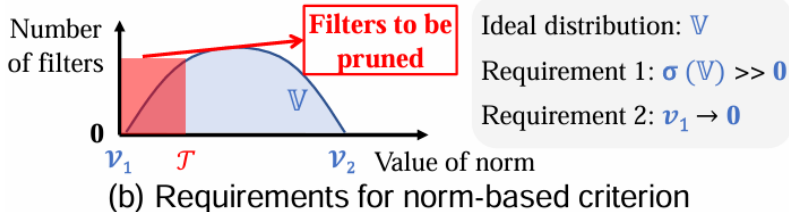
## Introduction

- **filter pruning** directly discards the whole selected filters and leaves a model with regular structures.
- **filter pruning** is more preferred for accelerating the networks and decreasing the model size.
- However, recent practice performs filter pruning by following the “**smaller-norm-less-important**” criterion which relies on two prerequisites that are not always true.
- To solve this problem, the authors propose a novel filter pruning method to compress the CNN model regardless of those two requirements namely **Filter Pruning via Geometric Median (FPGM)**.

# Smaller-norm-less-important criterion



- in the **norm-based criterion**, only the filters with the largest norm are kept.
- the proposed method prunes the filters with **redundant information** in the network.



- Req 1 makes the searching space for  $\mathcal{T}$  wide enough so that separating those filters needed to be pruned would be an easy task.
- Req 2 means that filters with smaller norms are expected to make **small contributions**.
- The **FPGM** chooses filters with **the most replaceable contribution** rather than **relatively less contribution**.

# Methodology

---

## Notation

- $N_i, N_{i+1}$  : the number of input channels and the output channels for the  $i_{th}$  convolution layer.
- $\mathcal{F}_{i,j} \in \mathbb{R}^{N_i \times K \times K}$  :  $j_{th}$  filters of the  $i_{th}$  layer.
- $K$  : kernel size of the network.
- $L$  : the number of layers.
- $\mathbf{W}^{(i)} \in \mathbb{R}^{N_{i+1} \times N_i \times K \times K}, 1 \leq i \leq L$  : model parameter of the  $i_{th}$  layer.



# Filter Pruning via Geometric Median

## Geometric Median

Given a set of  $n$  points  $a^{(1)}, \dots, a^{(n)}$  with each  $a^{(i)} \in \mathbb{R}^d$ , The geometric median is a point  $x^*$  such that minimizes the sum of Euclidean distances to them. i.e,

$$x^* = \arg \min_{x \in \mathbb{R}^d} f(x) \quad \text{where } f(x) \stackrel{\text{def}}{=} \sum_{i \in \{1, \dots, n\}} \|x - a^{(i)}\|_2 \quad (1)$$

- by (1), We use the **geometric median** to get the **common information** of all the filters within the single  $i_{th}$  layer. The geometric median  $x^{GM}$  becomes

$$x^{GM} = \arg \min_{x \in \mathbb{R}^{N_i \times K \times K}} \sum_{j' \in \{1, \dots, N_{i+1}\}} \|x - \mathcal{F}_{i,j'}\|_2 \quad (2)$$

- The filter(s) nearest to the geometric median  $x^{GM}$   $\mathcal{F}_{i,j^*}$  is

$$\mathcal{F}_{i,j^*} = \arg \min_{\mathcal{F}_{i,j'}} \|\mathcal{F}_{i,j'} - x^{GM}\|_2, \quad s.t \quad j' \in \{1, \dots, N_{i+1}\} \quad (3)$$

- As geometric median is a non-trivial problem in computational geometry, the authors do not compute it directly, Instead, they find the filter which **minimizes the summation of the distance with other filters**:

$$\mathcal{F}_{i,x^*} = \arg \min_x \sum_{j' \in \{1, \dots, N_{i+1}\}} \|x - \mathcal{F}_{i,j'}\|_2 \stackrel{\text{def}}{=} \arg \min_x g(x) \quad (4)$$

$s.t \quad x \in \{\mathcal{F}_{i,1}, \dots, \mathcal{F}_{i,N_{i+1}}\}$

- Note that even if  $\mathcal{F}_{i,x^*}$  is not included in the calculation of the geometric median in (4), we obtain the same result.

$$g(x) = \sum_{j' \in [1, N_{i+1}]} \|x - \mathcal{F}_{i,j'}\|_2 \quad s.t \quad x \in \{\mathcal{F}_{i,1}, \dots, \mathcal{F}_{i,N_{i+1}}\} \quad (5)$$

$$= \sum_{j' \in [1, N_{i+1}], \mathcal{F}_{i,j'} \neq x} \|x - \mathcal{F}_{i,j'}\|_2 + \left[ \|x - \mathcal{F}_{i,j'}\|_2 \right]_{\mathcal{F}_{i,j'}=x}^0 \quad (6)$$

$$= g'(x) \quad (7)$$

---

## Algorithm 1 Algorithm Description of FPGM

---

**Input:** training data:  $\mathbf{X}$

- 1: **Given:** pruning rate  $P_i$
- 2: **Initialize:** model parameter  $\mathbf{W} = \{\mathbf{W}^{(i)}, 0 \leq i \leq L\}$
- 3: **for**  $epoch = 1; epoch \leq epoch_{max}; epoch++$  **do**
- 4:     Update the model parameter  $\mathbf{W}$  based on  $\mathbf{X}$
- 5:     **for**  $i = 1; i \leq L; i++$  **do**
- 6:         Find  $N_{i+1}P_i$  filters that satisfy Equation 4
- 7:         Zeroize selected filters
- 8:     **end for**
- 9: **end for**
- 10: Obtain the compact model  $\mathbf{W}^*$  from  $\mathbf{W}$

**Output:** The compact model and its parameters  $\mathbf{W}^*$

---

Using FPGM algorithm, A compact model

$\{\mathbf{W}^{*(i)} \in \mathbb{R}^{N_{i+1}(1-P_i) \times N_i(1-P_{i-1}) \times K \times K}\}$  is obtained.

# Experiments

---

# Experimental Settings

- Dataset
  - CIFAR-10 : 50,000 training images, 10,000 testing images , 10 classes.
  - ILSVRC-2012 : 1.28M training images, 50k validation images, 1,000 classes.
- Models :
  - VGGNet(single-branch network)
  - ResNet(multiple-branch network)
- Pretrained model and scratch model are compared.
- The Pruning rate at each layer is the same.  
( $P_i = P \quad \forall i \in \{1, \dots, L\}$  ).
- The authors use **a mixture of FPGM and previous norm-based method** to show that FPGM could serve as a supplement to pervious methods.

# Single-Branch Network Pruning

Model \ Acc (%)	Baseline	Pruned w.o. FT	FT 40 epochs	FT 160 epochs
PFEC [21]	93.58 ( $\pm 0.03$ )	77.45 ( $\pm 0.03$ )	93.22 ( $\pm 0.03$ )	93.28 ( $\pm 0.07$ )
Ours	93.58 ( $\pm 0.03$ )	<b>80.38</b> ( $\pm 0.03$ )	<b>93.24</b> ( $\pm 0.01$ )	<b>94.00</b> ( $\pm 0.13$ )

Table 3. Pruning pre-trained VGGNet on CIFAR-10. “w.o.” means “without” and “FT” means “fine-tuning” the pruned model.

Model	SA	Baseline	Pruned From Scratch	FLOPs $\downarrow$ (%)
PFEC [21]	Y	93.58 ( $\pm 0.03$ )	93.31 ( $\pm 0.02$ )	34.2
Ours	Y	93.58 ( $\pm 0.03$ )	<b>93.54</b> ( $\pm 0.08$ )	34.2
Ours	N	93.58 ( $\pm 0.03$ )	93.23 ( $\pm 0.13$ )	<b>35.9</b>

Table 4. Pruning scratch VGGNet on CIFAR-10. “SA” means “sensitivity analysis”. Without sensitivity analysis, FPGM can still achieve comparable performances comparing to [21]; after introducing sensitivity analysis, FPGM can surpass [21].

# Multiple-Branch Network Pruning - CIFAR-10

Depth	Method	Fine-tune?	Baseline acc. (%)	Accelerated acc. (%)	Acc. ↓ (%)	FLOPs	FLOPs ↓ (%)
20	SFP [15]	✗	92.20 (±0.18)	90.83 (±0.31)	1.37	2.43E7	42.2
	Ours (FPGM-only 30%)	✗	92.20 (±0.18)	91.09 (±0.10)	1.11	2.43E7	42.2
	Ours (FPGM-only 40%)	✗	92.20 (±0.18)	90.44 (±0.20)	1.76	<b>1.87E7</b>	<b>54.0</b>
	Ours (FPGM-mix 40%)	✗	92.20 (±0.18)	<b>90.62</b> (±0.17)	<b>1.58</b>	<b>1.87E7</b>	<b>54.0</b>
32	MIL [5]	✗	92.33	90.74	1.59	4.70E7	31.2
	SFP [15]	✗	92.63 (±0.70)	92.08 (±0.08)	0.55	4.03E7	41.5
	Ours (FPGM-only 30%)	✗	92.63 (±0.70)	92.31 (±0.30)	0.32	4.03E7	41.5
	Ours (FPGM-only 40%)	✗	92.63 (±0.70)	<b>91.93</b> (±0.03)	<b>0.70</b>	<b>3.23E7</b>	<b>53.2</b>
	Ours (FPGM-mix 40%)	✗	92.63 (±0.70)	91.91 (±0.21)	0.72	<b>3.23E7</b>	<b>53.2</b>
56	PFEC [21]	✗	93.04	91.31	1.75	9.09E7	27.6
	CP [16]	✗	92.80	90.90	1.90	—	50.0
	SFP [15]	✗	93.59 (±0.58)	92.26 (±0.31)	1.33	<b>5.94E7</b>	<b>52.6</b>
	Ours (FPGM-only 40%)	✗	93.59 (±0.58)	<b>92.93</b> (±0.49)	<b>0.66</b>	<b>5.94E7</b>	<b>52.6</b>
	Ours (FPGM-mix 40%)	✗	93.59 (±0.58)	92.89 (±0.32)	0.70	<b>5.94E7</b>	<b>52.6</b>
	PFEC [21]	✓	93.04	93.06	-0.02	9.09E7	27.6
	CP [16]	✓	92.80	91.80	1.00	—	50.0
	Ours (FPGM-only 40%)	✓	93.59 (±0.58)	<b>93.49</b> (±0.13)	<b>0.10</b>	<b>5.94E7</b>	<b>52.6</b>
	Ours (FPGM-mix 40%)	✓	93.59 (±0.58)	93.26 (±0.03)	0.33	<b>5.94E7</b>	<b>52.6</b>
110	MIL [5]	✗	93.63	93.44	0.19	—	34.2
	PFEC [21]	✗	93.53	92.94	0.61	1.55E8	38.6
	SFP [15]	✗	93.68 (±0.32)	93.38 (±0.30)	0.30	1.50E8	40.8
	Ours (FPGM-only 40%)	✗	93.68 (±0.32)	93.73 (±0.23)	-0.05	<b>1.21E8</b>	<b>52.3</b>
	Ours (FPGM-mix 40%)	✗	93.68 (±0.32)	<b>93.85</b> (±0.11)	<b>-0.17</b>	<b>1.21E8</b>	<b>52.3</b>
	PFEC [21]	✓	93.53	93.30	0.20	1.55E8	38.6
	NISP [39]	✓	—	—	0.18	—	43.8
	Ours (FPGM-only 40%)	✓	93.68 (±0.32)	<b>93.74</b> (±0.10)	<b>-0.16</b>	<b>1.21E8</b>	<b>52.3</b>

# Multiple-Branch Network Pruning - ILSVRC-2012

Depth	Method	Fine-tune?	Baseline top-1 acc.(%)	Accelerated top-1 acc.(%)	Baseline top-5 acc.(%)	Accelerated top-5 acc.(%)	Top-1 acc. ↓(%)	Top-5 acc. ↓(%)	FLOPs↓(%)
18	MIL [5]	✗	69.98	66.33	89.24	86.94	3.65	2.30	34.6
	SFP [15]	✗	<b>70.28</b>	67.10	<b>89.63</b>	87.78	3.18	1.85	<b>41.8</b>
	Ours (FPGM-only 30%)	✗	<b>70.28</b>	67.78	<b>89.63</b>	88.01	2.50	1.62	<b>41.8</b>
	Ours (FPGM-mix 30%)	✗	<b>70.28</b>	<b>67.81</b>	<b>89.63</b>	<b>88.11</b>	<b>2.47</b>	<b>1.52</b>	<b>41.8</b>
	Ours (FPGM-only 30%)	✓	<b>70.28</b>	68.34	<b>89.63</b>	<b>88.53</b>	1.94	<b>1.10</b>	<b>41.8</b>
	Ours (FPGM-mix 30%)	✓	<b>70.28</b>	<b>68.41</b>	<b>89.63</b>	88.48	<b>1.87</b>	1.15	<b>41.8</b>
34	SFP [15]	✗	<b>73.92</b>	71.83	<b>91.62</b>	90.33	2.09	1.29	<b>41.1</b>
	Ours (FPGM-only 30%)	✗	<b>73.92</b>	71.79	<b>91.62</b>	<b>90.70</b>	2.13	<b>0.92</b>	<b>41.1</b>
	Ours (FPGM-mix 30%)	✗	<b>73.92</b>	<b>72.11</b>	<b>91.62</b>	90.69	<b>1.81</b>	0.93	<b>41.1</b>
	PFEC [21]	✓	73.23	72.17	–	–	<b>1.06</b>	–	24.2
	Ours (FPGM-only 30%)	✓	<b>73.92</b>	72.54	<b>91.62</b>	<b>91.13</b>	1.38	<b>0.49</b>	<b>41.1</b>
	Ours (FPGM-mix 30%)	✓	<b>73.92</b>	<b>72.63</b>	<b>91.62</b>	91.08	1.29	0.54	<b>41.1</b>
50	SFP [15]	✗	<b>76.15</b>	74.61	<b>92.87</b>	92.06	1.54	0.81	41.8
	Ours (FPGM-only 30%)	✗	<b>76.15</b>	<b>75.03</b>	<b>92.87</b>	<b>92.40</b>	<b>1.12</b>	<b>0.47</b>	42.2
	Ours (FPGM-mix 30%)	✗	<b>76.15</b>	74.94	<b>92.87</b>	92.39	1.21	0.48	42.2
	Ours (FPGM-only 40%)	✗	<b>76.15</b>	74.13	<b>92.87</b>	91.94	2.02	0.93	<b>53.5</b>
	ThiNet [25]	✓	72.88	72.04	91.14	90.67	0.84	0.47	36.7
	SFP [15]	✓	<b>76.15</b>	62.14	<b>92.87</b>	84.60	14.01	8.27	41.8
	NISP [39]	✓	–	–	–	–	0.89	–	44.0
	CP [16]	✓	–	–	92.20	90.80	–	1.40	50.0
	Ours (FPGM-only 30%)	✓	<b>76.15</b>	<b>75.59</b>	<b>92.87</b>	<b>92.63</b>	<b>0.56</b>	0.24	42.2
	Ours (FPGM-mix 30%)	✓	<b>76.15</b>	75.50	<b>92.87</b>	92.63	0.65	<b>0.21</b>	42.2
	Ours (FPGM-only 40%)	✓	<b>76.15</b>	74.83	<b>92.87</b>	92.32	1.32	0.55	<b>53.5</b>
101	Rethinking [38]	✓	<b>77.37</b>	75.27	–	–	2.10	–	<b>47.0</b>
	Ours (FPGM-only 30%)	✓	<b>77.37</b>	<b>77.32</b>	<b>93.56</b>	<b>93.56</b>	<b>0.05</b>	<b>0.00</b>	42.2



## Conclusion

---

- this paper elaborates on the underlying requirements for norm-based criterion and points out their limitations.
- FPGM prunes the most replaceable filters containing redundant information, which can still achieve good performances when norm-based criterion fails;
- FPGM considers the mutual relations between filters. Thanks to this, FPGM achieves the SOTA performance in several benchmarks.