

XNOR-Net: ImageNet Classification Using Binary Convolutional Neural Networks

2018-11-11

Overview

- 提出 **Binary-Weight-Networks & XNOR-Net**
- **Binary-Weight-Networks:** 二值化filter(weight) ; smaller 32x than an equivalent network with single-precision weight values & 2 speed up
- **XNOR-Net:** 二值化input & weight ; smaller 32x than an equivalent network with single-precision weight values & offering 58x speed up in CPUs

XNOR-Net: ImageNet Classification Using Binary Convolutional Neural Networks

- Mohammad Rastegari, Vicente Ordonez, Joseph Redmon, Ali Farhadi
- Torch 7
- 2016

Binary-Weight-Networks

- 将卷积操作替换为 $\mathbf{I} * \mathbf{W} \approx (\mathbf{I} \oplus \mathbf{B}) \alpha$
- Estimating binary weights \longrightarrow
- Binarization in Forward & Backward

Algorithm 1 Training an L -layers CNN with binary weights:

Input: A minibatch of inputs and targets (\mathbf{I}, \mathbf{Y}) , cost function $C(\mathbf{Y}, \hat{\mathbf{Y}})$, current weight \mathcal{W}^t and current learning rate η^t .

Output: updated weight \mathcal{W}^{t+1} and updated learning rate η^{t+1} .

- 1: Binarizing weight filters:
 - 2: **for** $l = 1$ to L **do**
 - 3: **for** k^{th} filter in l^{th} layer **do**
 - 4: $\mathcal{A}_{lk} = \frac{1}{n} \|\mathcal{W}_{lk}^t\|_{\ell_1}$
 - 5: $\mathcal{B}_{lk} = \text{sign}(\mathcal{W}_{lk}^t)$
 - 6: $\tilde{\mathcal{W}}_{lk} = \mathcal{A}_{lk} \mathcal{B}_{lk}$
 - 7: $\hat{\mathbf{Y}} = \text{BinaryForward}(\mathbf{I}, \mathcal{B}, \mathcal{A})$ // standard forward propagation except that convolutions are computed using equation 1 or 11
 - 8: $\frac{\partial C}{\partial \tilde{\mathcal{W}}} = \text{BinaryBackward}(\frac{\partial C}{\partial \hat{\mathbf{Y}}}, \tilde{\mathcal{W}})$ // standard backward propagation except that gradients are computed using $\tilde{\mathcal{W}}$ instead of \mathcal{W}^t
 - 9: $\mathcal{W}^{t+1} = \text{UpdateParameters}(\mathcal{W}^t, \frac{\partial C}{\partial \tilde{\mathcal{W}}}, \eta^t)$ // Any update rules (e.g., SGD or ADAM)
 - 10: $\eta^{t+1} = \text{UpdateLearningrate}(\eta^t, t)$ // Any learning rate scheduling function
-

$$J(\mathbf{B}, \alpha) = \|\mathbf{W} - \alpha \mathbf{B}\|^2$$

$$\alpha^*, \mathbf{B}^* = \underset{\alpha, \mathbf{B}}{\text{argmin}} J(\mathbf{B}, \alpha) \quad (2)$$

$$J(\mathbf{B}, \alpha) = \alpha^2 \mathbf{B}^T \mathbf{B} - 2\alpha \mathbf{W}^T \mathbf{B} + \mathbf{W}^T \mathbf{W} \quad (3)$$

$$\mathbf{B}^* = \underset{\mathbf{B}}{\text{argmax}} \{\mathbf{W}^T \mathbf{B}\} \quad s.t. \quad \mathbf{B} \in \{+1, -1\}^n \quad (4)$$

$$\alpha^* = \frac{\mathbf{W}^T \mathbf{B}^*}{n} \quad (5)$$

By replacing \mathbf{B}^* with $\text{sign}(\mathbf{W})$

$$\alpha^* = \frac{\mathbf{W}^T \text{sign}(\mathbf{W})}{n} = \frac{\sum |\mathbf{W}_i|}{n} = \frac{1}{n} \|\mathbf{W}\|_{\ell_1} \quad (6)$$

α 的最优解是 \mathbf{W} 的每个元素的绝对值之和的均值

XNOR-Net

- Binary Dot Product(similar as Estimating binary weights)

$$\alpha^*, \mathbf{B}^*, \beta^*, \mathbf{H}^* = \underset{\alpha, \mathbf{B}, \beta, \mathbf{H}}{\operatorname{argmin}} \|\mathbf{X} \odot \mathbf{W} - \beta \alpha \mathbf{H} \odot \mathbf{B}\| \quad (7) \quad \beta \mathbf{H} \text{ 近似表示输入 } \mathbf{X}$$

$$\gamma^*, \mathbf{C}^* = \underset{\gamma, \mathbf{C}}{\operatorname{argmin}} \|\mathbf{Y} - \gamma \mathbf{C}\| \quad (8)$$

$$\mathbf{C}^* = \operatorname{sign}(\mathbf{Y}) = \operatorname{sign}(\mathbf{X}) \odot \operatorname{sign}(\mathbf{W}) = \mathbf{H}^* \odot \mathbf{B}^* \quad (9)$$

Since $|\mathbf{X}_i|, |\mathbf{W}_i|$ are independent, knowing that $\mathbf{Y}_i = \mathbf{X}_i \mathbf{W}_i$ then,
 $\mathbf{E}[|\mathbf{Y}_i|] = \mathbf{E}[|\mathbf{X}_i| |\mathbf{W}_i|] = \mathbf{E}[|\mathbf{X}_i|] \mathbf{E}[|\mathbf{W}_i|]$ therefore,

$$\gamma^* = \frac{\sum |\mathbf{Y}_i|}{n} = \frac{\sum |\mathbf{X}_i| |\mathbf{W}_i|}{n} \approx \left(\frac{1}{n} \|\mathbf{X}\|_{\ell_1} \right) \left(\frac{1}{n} \|\mathbf{W}\|_{\ell_1} \right) = \beta^* \alpha^* \quad (10)$$

- Binary Convolution (使用XNOR替代乘法提速)

$$\mathbf{I} * \mathbf{W} \approx (\operatorname{sign}(\mathbf{I}) \otimes \operatorname{sign}(\mathbf{W})) \odot \mathbf{K} \alpha \quad (11)$$

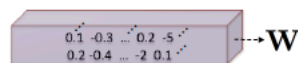
where \otimes indicates a convolutional operation using XNOR and bitcount operations. This

XNOR-Net: ImageNet Classification Using Binary Convolutional Neural Networks

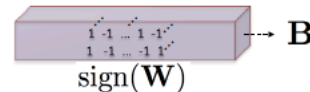


XNOR-Net

(1) Binarizing Weight

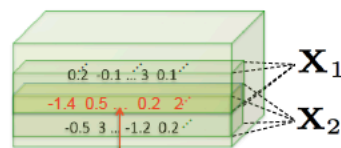


$$\frac{1}{n} \|\mathbf{W}\|_{\ell_1} = \alpha$$



(2) Binarizing Input

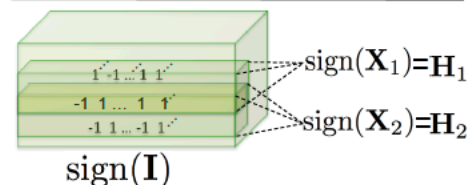
Inefficient



Redundant computations in overlapping areas

$$\frac{1}{n} \|\mathbf{X}_1\|_{\ell_1} = \beta_1$$

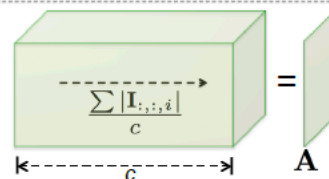
$$\frac{1}{n} \|\mathbf{X}_2\|_{\ell_1} = \beta_2$$



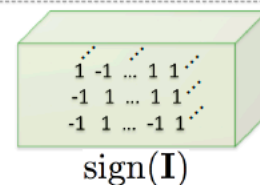
存在重复计算

(3) Binarizing Input

Efficient



$$\mathbf{A} * \mathbf{k} = \mathbf{K}$$

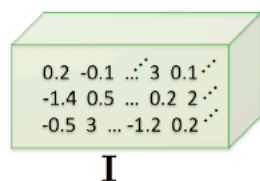


average over absolute values of the elements in the input \mathbf{I}

\mathbf{c} means number of channels

we convolve \mathbf{A} with a 2D filter $\mathbf{k} \in \mathbb{R}^{w \times h}$, $\mathbf{K} = \mathbf{A} * \mathbf{k}$, where $\forall ij \quad k_{ij} = \frac{1}{w \times h}$.

(4) Convolution with XNOR-Bitcount



$$*$$

$$\approx \left[\begin{array}{c} \text{sign}(\mathbf{I}) \\ \text{sign}(\mathbf{W}) \end{array} \right] \odot \odot \alpha$$

Experiments

Classification Accuracy(%)									
Binary-Weight				Binary-Input-Binary-Weight				Full-Precision	
BWN		BC[11]		XNOR-Net		BNN[11]		AlexNet[1]	
Top-1	Top-5	Top-1	Top-5	Top-1	Top-5	Top-1	Top-5	Top-1	Top-5
56.8	79.4	35.4	61.0	44.2	69.2	27.9	50.42	56.6	80.2

Table 1: This table compares the final accuracies (Top1 - Top5) of the full precision network with our binary precision networks; Binary-Weight-Networks(BWN) and XNOR-Networks(XNOR-Net) and the competitor methods; BinaryConnect(BC) and BinaryNet(BNN).

Binary-Weight-Network			XNOR-Network		
Strategy for computing α	top-1	top-5	Block Structure	top-1	top-5
Using equation 6	56.8	79.4	C-B-A-P	30.3	57.5
Using a separate layer	46.2	69.5	B-A-C-P	44.2	69.2

(a)

(b)

Table 3: In this table, we evaluate two key elements of our approach; computing the optimal scaling factors and specifying the right order for layers in a block of CNN with binary input. (a) demonstrates the importance of the scaling factor in training binary-weight-networks and (b) shows that our way of ordering the layers in a block of CNN is crucial for training XNOR-Networks. C,B,A,P stands for Convolutional, BatchNormalization, Active function (here binary activation), and Pooling respectively.

Network Variations	ResNet-18		GoogLenet	
	top-1	top-5	top-1	top-5
Binary-Weight-Network	60.8	83.0	65.5	86.1
XNOR-Network	51.2	73.2	N/A	N/A
Full-Precision-Network	69.3	89.2	71.3	90.0

Table 2: This table compares the final classification accuracy achieved by our binary precision networks with the full precision network in ResNet-18 and GoogLenet architectures.

精度效果

- **Binary-Weight-Networks:** 二值化filter(weight) ; smaller 32x than an equivalent network with single-precision weight values & 2 speed up
- **XNOR-Net:** 二值化input & weight ; smaller 32x than an equivalent network with single-precision weight values & offering 58x speed up in CPUs

压缩效果