



Training Binary Neural Networks

11月8日
周争光



二值化

□ 二值化

- 浮点数 $\rightarrow \{-1, 1\}$

□ 权值二值化

- 卷积乘法 \rightarrow 加法与减法
- 减少 ~ 32 倍参数量
- ~ 2 倍理论加速

□ 权值二值化+激活值二值化

- 卷积乘法 \rightarrow 位运算 (XNOR, bitcount)
- 减少 ~ 32 倍参数量
- ~ 58 倍理论加速



BinaryConnect

□ 权值二值化

- 确定方式

- 随机方式

$$w_b = \begin{cases} +1 & \text{if } w \geq 0, \\ -1 & \text{otherwise.} \end{cases}$$

$$w_b = \begin{cases} +1 & \text{with probability } p = \sigma(w), \\ -1 & \text{with probability } 1 - p. \end{cases}$$

$$\sigma(x) = \text{clip}\left(\frac{x+1}{2}, 0, 1\right) = \max\left(0, \min\left(1, \frac{x+1}{2}\right)\right)$$

□ 反向传播与梯度计算

- 保留 w 浮点数值

- 梯度累加，将超出 $[-1, 1]$ 的部分裁剪

$$w_t \leftarrow \text{clip}\left(w_{t-1} - \eta \frac{\partial C}{\partial w_b}\right)$$



BinaryConnect

□ 训练算法

Algorithm 1 SGD training with BinaryConnect. C is the cost function for minibatch and the functions $\text{binarize}(w)$ and $\text{clip}(w)$ specify how to binarize and clip weights. L is the number of layers.

Require: a minibatch of (inputs, targets), previous parameters w_{t-1} (weights) and b_{t-1} (biases), and learning rate η .

Ensure: updated parameters w_t and b_t .

1. Forward propagation:

$w_b \leftarrow \text{binarize}(w_{t-1})$

For $k = 1$ to L , compute a_k knowing a_{k-1} , w_b and b_{t-1}

2. Backward propagation:

Initialize output layer's activations gradient $\frac{\partial C}{\partial a_L}$

For $k = L$ to 2, compute $\frac{\partial C}{\partial a_{k-1}}$ knowing $\frac{\partial C}{\partial a_k}$ and w_b

3. Parameter update:

Compute $\frac{\partial C}{\partial w_b}$ and $\frac{\partial C}{\partial b_{t-1}}$ knowing $\frac{\partial C}{\partial a_k}$ and a_{k-1}

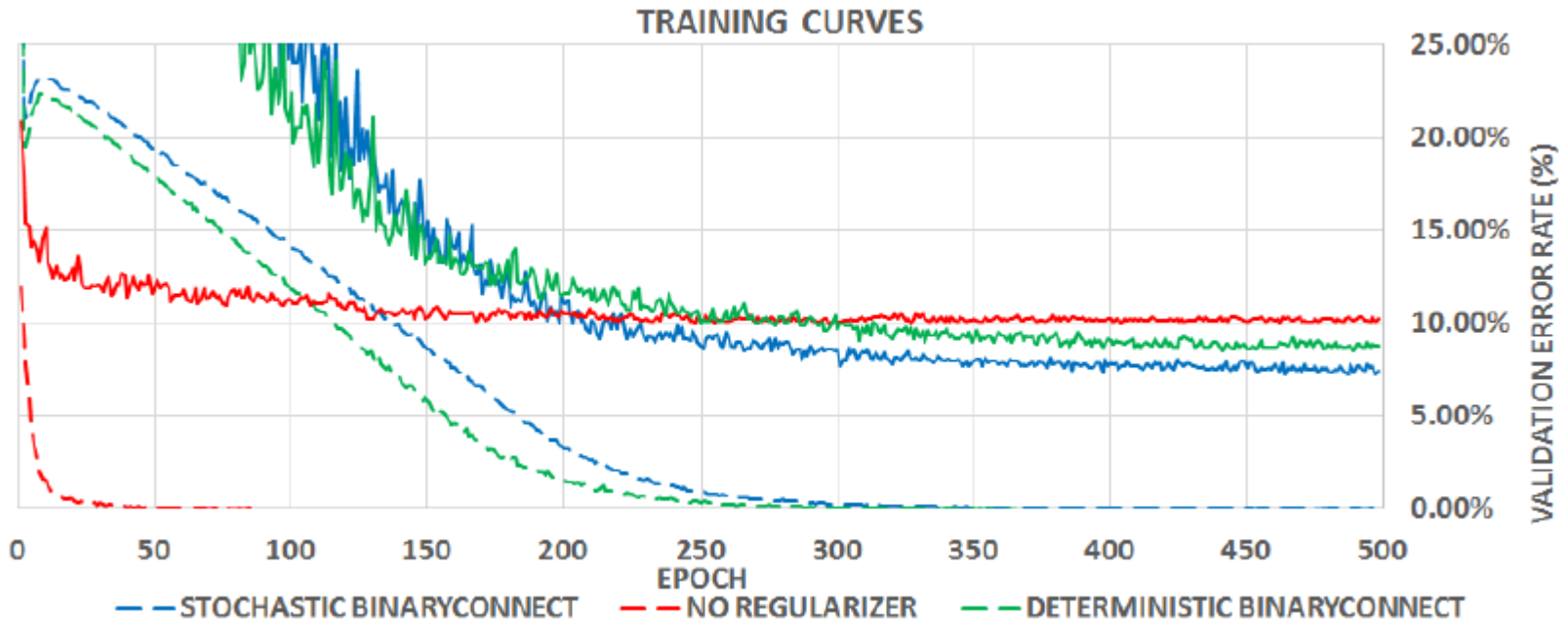
$w_t \leftarrow \text{clip}(w_{t-1} - \eta \frac{\partial C}{\partial w_b})$

$b_t \leftarrow b_{t-1} - \eta \frac{\partial C}{\partial b_{t-1}}$



BinaryConnect

□ 实验结果



- 二值化相当于给权值添加了噪声，具有正则化作用，可以防止模型过拟合



Binarized Neural Networks (BNN)

□ 量化函数（权值+激活值）

$$x^b = \text{Sign}(x) = \begin{cases} +1 & \text{if } x \geq 0, \\ -1 & \text{otherwise,} \end{cases}$$

- 有确定和随机方式两种
- 但随机形式需要生成随机数耗时，一般都采用确定形式的

□ 对sign求导（“straight through estimator”）

$$q = \text{Sign}(r),$$

$$g_r = g_q \mathbf{1}_{|r| \leq 1}.$$

$$\text{Htanh}(x) = \text{Clip}(x, -1, 1)$$



BNN

```
{1. Computing the parameters gradients:}
{1.1. Forward propagation:}
for  $k = 1$  to  $L$  do
   $W_k^b \leftarrow \text{Binarize}(W_k)$ 
   $s_k \leftarrow a_{k-1}^b W_k^b$ 
   $a_k \leftarrow \text{BatchNorm}(s_k, \theta_k)$ 
  if  $k < L$  then
     $a_k^b \leftarrow \text{Binarize}(a_k)$ 
  end if
end for
{1.2. Backward propagation:}
{Please note that the gradients are not binary.}
Compute  $g_{a_L} = \frac{\partial C}{\partial a_L}$  knowing  $a_L$  and  $a^*$ 
for  $k = L$  to 1 do
  if  $k < L$  then
     $g_{a_k} \leftarrow g_{a_k^b} \circ 1_{|a_k| \leq 1}$ 
  end if
   $(g_{s_k}, g_{\theta_k}) \leftarrow \text{BackBatchNorm}(g_{a_k}, s_k, \theta_k)$ 
   $g_{a_{k-1}^b} \leftarrow g_{s_k} W_k^b$ 
   $g_{W_k^b} \leftarrow g_{s_k} a_{k-1}^b$ 
end for
{2. Accumulating the parameters gradients:}
for  $k = 1$  to  $L$  do
   $\theta_k^{t+1} \leftarrow \text{Update}(\theta_k, \eta, g_{\theta_k})$ 
   $W_k^{t+1} \leftarrow \text{Clip}(\text{Update}(W_k, \gamma_k \eta, g_{W_k^b}), -1, 1)$ 
   $\eta^{t+1} \leftarrow \lambda \eta$ 
end for
```

- 训练
- 运行

Algorithm 5 Running a BNN. L is the number of layers.

Require: a vector of 8-bit inputs a_0 , the binary weights W^b , and the BatchNorm parameters θ .

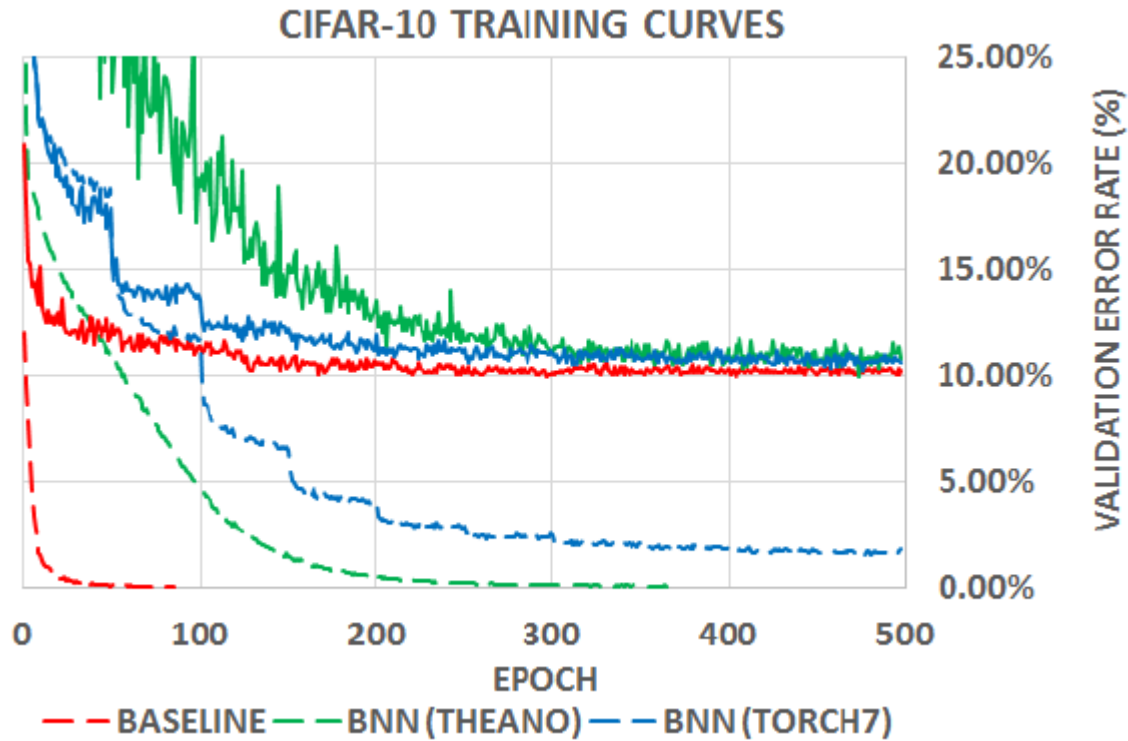
Ensure: the MLP output a_L .

```
{1. First layer:}
 $a_1 \leftarrow 0$ 
for  $n = 1$  to 8 do
   $a_1 \leftarrow a_1 + 2^{n-1} \times \text{XnorDotProduct}(a_0^n, W_1^b)$ 
end for
 $a_1^b \leftarrow \text{Sign}(\text{BatchNorm}(a_1, \theta_1))$ 
{2. Remaining hidden layers:}
for  $k = 2$  to  $L - 1$  do
   $a_k \leftarrow \text{XnorDotProduct}(a_{k-1}^b, W_k^b)$ 
   $a_k^b \leftarrow \text{Sign}(\text{BatchNorm}(a_k, \theta_k))$ 
end for
{3. Output layer:}
 $a_L \leftarrow \text{XnorDotProduct}(a_{L-1}^b, W_L^b)$ 
 $a_L \leftarrow \text{BatchNorm}(a_L, \theta_L)$ 
```



BNN

□ 实验结果





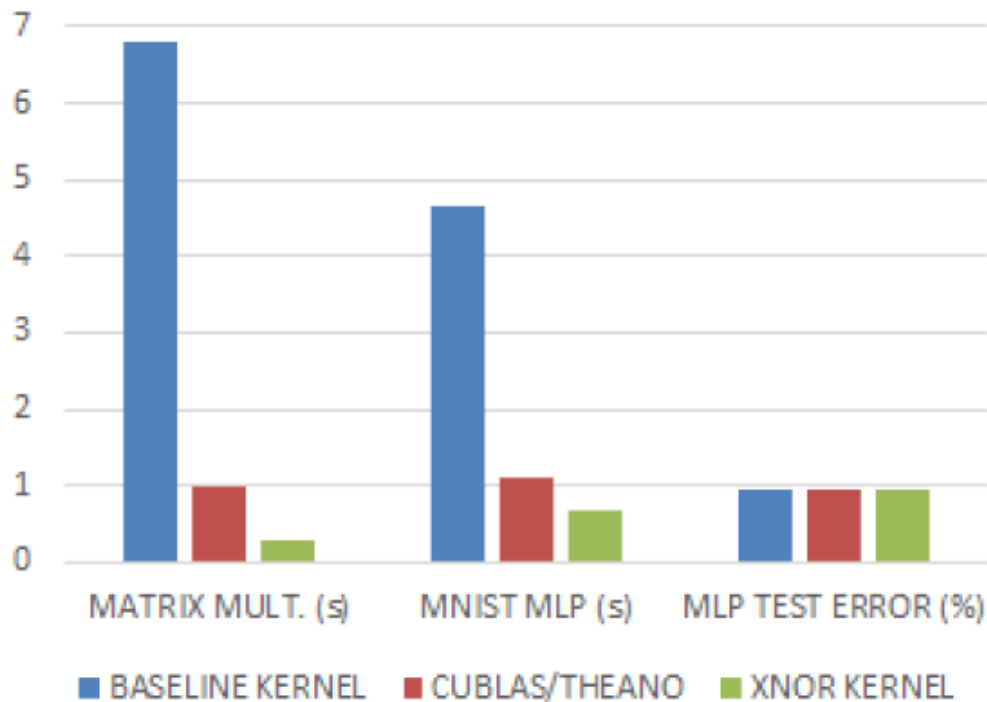
BNN

□ 实现优化

$$a_1 + = \text{popcount}(\text{xnor}(a_0^{32b}, w_1^{32b})),$$

- 位操作使用SIMD指令加速，将32二值数合并并在一个32位寄存器中，从而32个时钟现在可用1+4+1=6个时钟完成

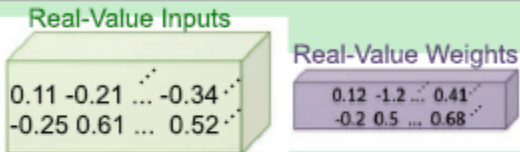
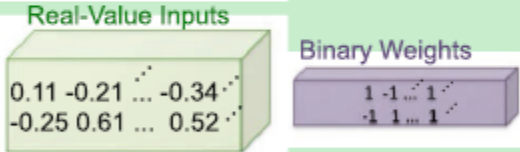

GPU KERNELS' EXECUTION TIMES



XNOR-Net

□ 提出了两种二值网络

- Binary-Weight-Networks (BWN): 只二值化权值
- XNOR-Networks (XNOR-Net): 二值化权值和激活值
- 由于引入缩放因子, BWN和XNOR-Net不是纯粹的二值网络

	Network Variations	Operations used in Convolution	Memory Saving (Inference)	Computation Saving (Inference)	Accuracy on ImageNet (AlexNet)
Standard Convolution	 <p>Real-Value Inputs</p> <p>Real-Value Weights</p>	+ , - , ×	1x	1x	%56.7
Binary Weight	 <p>Real-Value Inputs</p> <p>Binary Weights</p>	+ , -	~32x	~2x	%56.8
BinaryWeight Binary Input (XNOR-Net)	 <p>Binary Inputs</p> <p>Binary Weights</p>	XNOR , bitcount	~32x	~58x	%44.2



BWN

□ 近似

$$\mathbf{W} \approx \alpha \mathbf{B}, \quad \mathbf{I} * \mathbf{W} \approx (\mathbf{I} \oplus \mathbf{B}) \alpha, \quad \mathbf{B} \in \{+1, -1\}^{c \times w \times h}$$

□ 优化

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

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

□ 求解

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

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



BWN

□ 训练

$$\frac{\partial \text{sign}}{\partial r} = r \mathbf{1}_{|r| \leq 1}$$

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: $\widetilde{\mathcal{W}}_{lk} = \mathcal{A}_{lk} \mathcal{B}_{lk}$
- 7: $\hat{\mathbf{Y}} = \mathbf{BinaryForward}(\mathbf{I}, \mathcal{B}, \mathcal{A})$ // standard forward propagation except that convolutions are computed using Eq. 1 or 11
- 8: $\frac{\partial C}{\partial \widetilde{\mathcal{W}}} = \mathbf{BinaryBackward}(\frac{\partial C}{\partial \hat{\mathbf{Y}}}, \widetilde{\mathcal{W}})$ // standard backward propagation except that gradients are computed using $\widetilde{\mathcal{W}}$ instead of \mathcal{W}^t
- 9: $\mathcal{W}^{t+1} = \mathbf{UpdateParameters}(\mathcal{W}^t, \frac{\partial C}{\partial \widetilde{\mathcal{W}}}, \eta_t)$ // Any update rules (*e.g.*, SGD or ADAM)
- 10: $\eta^{t+1} = \mathbf{UpdateLearningrate}(\eta^t, t)$ // Any learning rate scheduling function



XNOR-Net

□ 优化

$$\alpha^*, \mathbf{B}^*, \beta^*, \mathbf{H}^* = \operatorname{argmin}_{\alpha, \mathbf{B}, \beta, \mathbf{H}} \|\mathbf{X} \odot \mathbf{W} - \beta \alpha \mathbf{H} \odot \mathbf{B}\|$$

$$\gamma^*, \mathbf{C}^* = \operatorname{argmin}_{\gamma, \mathbf{C}} \|\mathbf{Y} - \gamma \mathbf{C}\|$$

□ 求解同BWN

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

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

□ 卷积实现

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

XNOR-Net

卷积实现

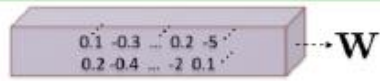
$$A = \frac{\sum |I_{:, :, i}|}{c}$$

$$K = A * k,$$

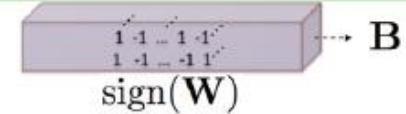
$$\forall ij \ k_{ij} = \frac{1}{w \times h}.$$

$$I * W \approx (\text{sign}(I) \circledast \text{sign}(W)) \odot K \alpha$$

(1) Binarizing Weight

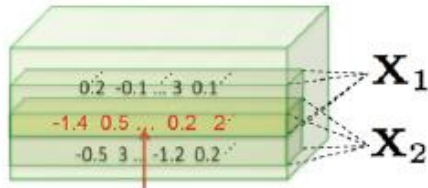


$$\frac{1}{n} \|W\|_{\ell_1} = \alpha$$



(2) Binarizing Input

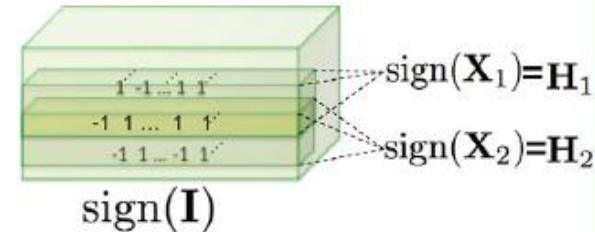
Inefficient



Redundant computations in overlapping areas

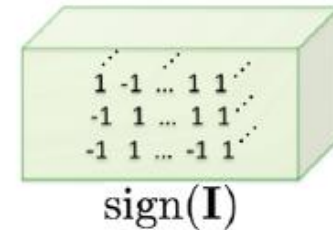
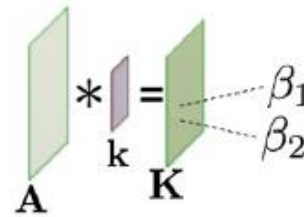
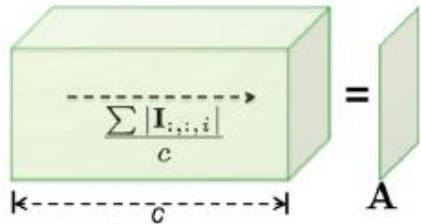
$$\frac{1}{n} \|X_1\|_{\ell_1} = \beta_1$$

$$\frac{1}{n} \|X_2\|_{\ell_1} = \beta_2$$

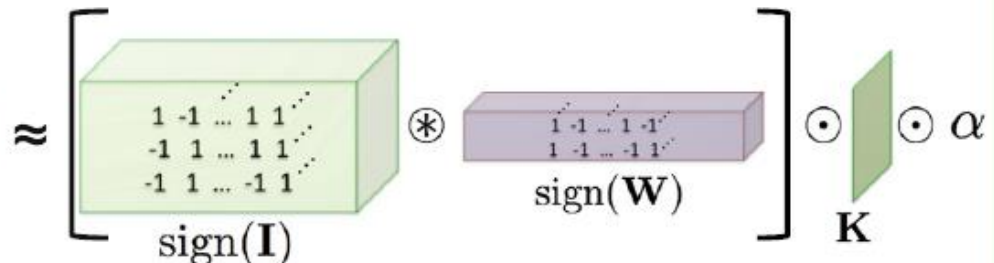
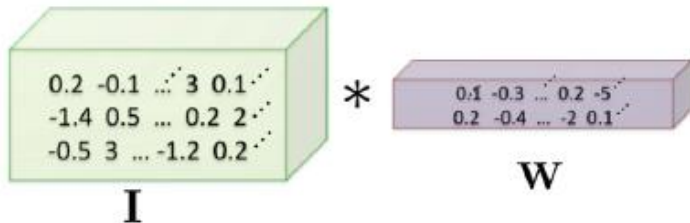


(3) Binarizing Input

Efficient

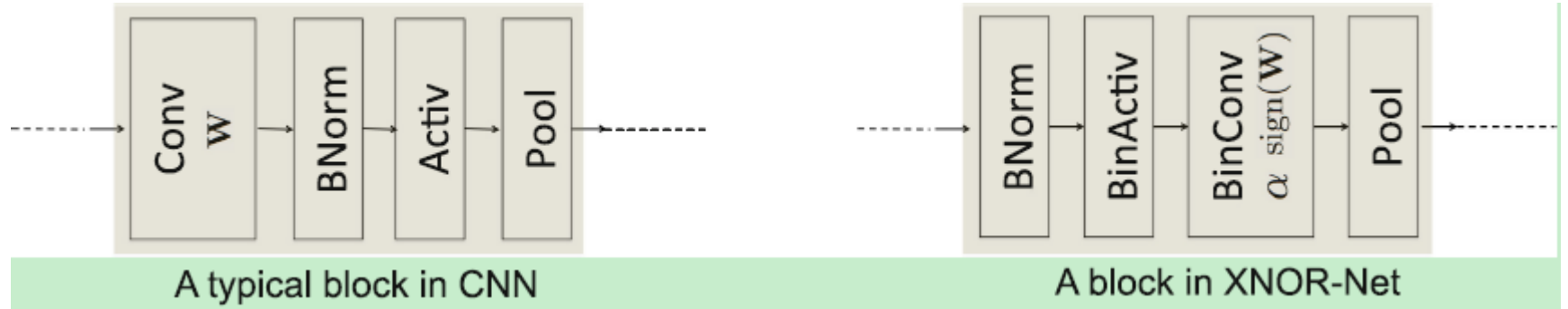


(4) Convolution with XNOR-Bitcount



XNOR-Net

□ 网络结构



□ 性能分析

- 一个卷积含有 $cN_{\mathbf{W}}N_{\mathbf{I}}$ 个二值运算和 $N_{\mathbf{I}}$ 个非二值运算
- 一个CPU时钟可运行64个二值运算，因此计算加速为：

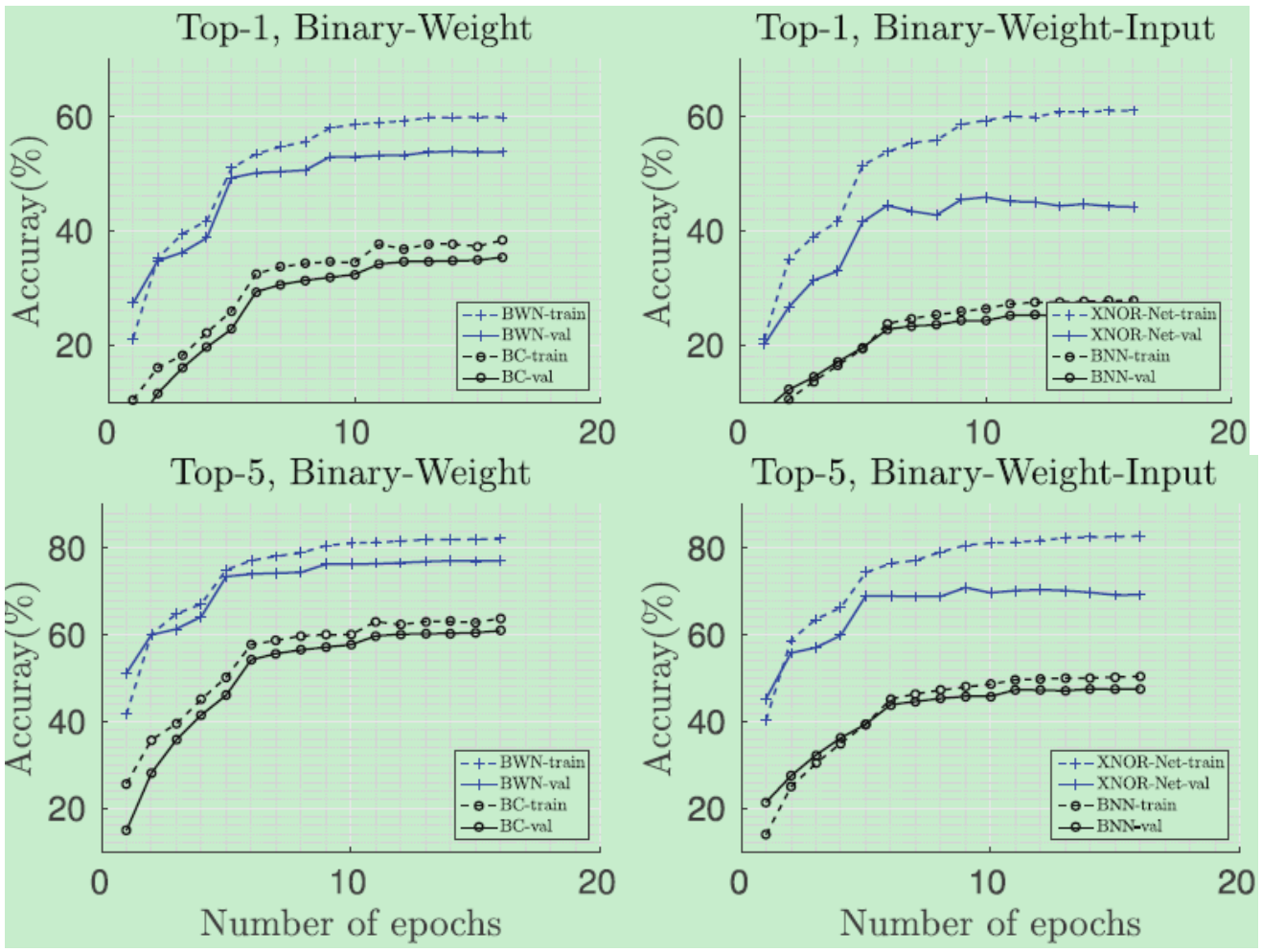
$$S = \frac{cN_{\mathbf{W}}N_{\mathbf{I}}}{\frac{1}{64}cN_{\mathbf{W}}N_{\mathbf{I}} + N_{\mathbf{I}}} = \frac{64cN_{\mathbf{W}}}{cN_{\mathbf{W}} + 64}$$

- 获得62.27倍理论加速但实际只有58倍加速

$$c = 256, n_{\mathbf{I}} = 14^2 \text{ and } n_{\mathbf{W}} = 3^2$$

实验结果

与BC和BNN方法对比





Ternary weight networks

□ 三值网络 (-1,0,1)

□ 优化

$$\begin{cases} \alpha^*, \mathbf{W}^{t*} = \arg \min_{\alpha, \mathbf{W}^t} J(\alpha, \mathbf{W}^t) = \|\mathbf{W} - \alpha \mathbf{W}^t\|_2^2 \\ \text{s.t.} \quad \alpha \geq 0, \mathbf{W}_i^t \in \{-1, 0, 1\}, i = 1, 2, \dots, n. \end{cases}$$

□ 不好求解，近似

$$\mathbf{W}_i^t = f_t(\mathbf{W}_i | \Delta) = \begin{cases} +1, & \text{if } \mathbf{W}_i > \Delta \\ 0, & \text{if } |\mathbf{W}_i| \leq \Delta \\ -1, & \text{if } \mathbf{W}_i < -\Delta \end{cases}$$

$$\alpha^*, \Delta^* = \arg \min_{\alpha \geq 0, \Delta > 0} (|\mathbf{I}_\Delta| \alpha^2 - 2 \left(\sum_{i \in \mathbf{I}_\Delta} |\mathbf{W}_i| \right) \alpha + c_\Delta)$$

$$\alpha_\Delta^* = \frac{1}{|\mathbf{I}_\Delta|} \sum_{i \in \mathbf{I}_\Delta} |\mathbf{W}_i|.$$

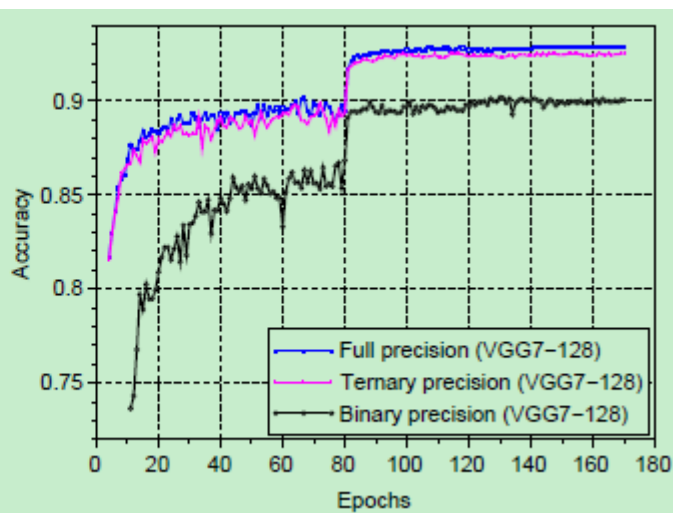
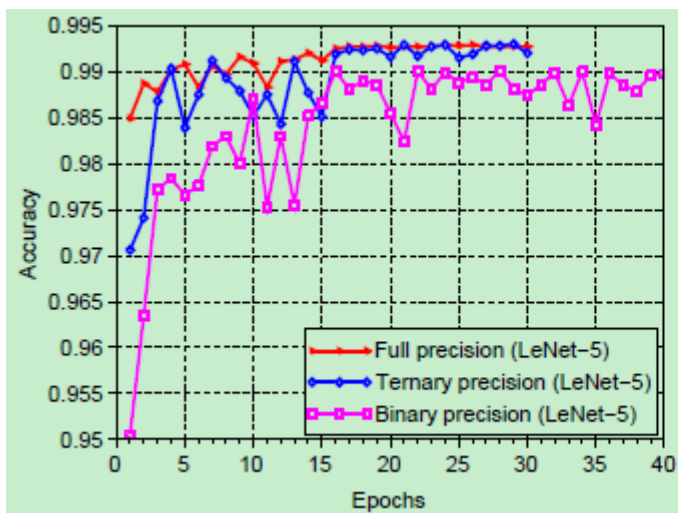
$$\Delta^* = \arg \max_{\Delta > 0} \frac{1}{|\mathbf{I}_\Delta|} \left(\sum_{i \in \mathbf{I}_\Delta} |\mathbf{W}_i| \right)^2$$

$$\Delta^* \approx 0.7 \cdot \mathbb{E}(|\mathbf{W}|) \approx \frac{0.7}{n} \sum_{i=1}^n |\mathbf{W}_i|$$

TWN

□ 实验结果

	MNIST	CIFAR-10	ImageNet (top-1)	ImageNet (top-5)
TWNs	99.35	92.56	61.8 / 65.3	84.2 / 86.2
BPWNs	99.05	90.18	57.5 / 61.6	81.2 / 83.9
FPWNs	99.41	92.88	65.4 / 67.6	86.76 / 88.0
BinaryConnect	98.82	91.73	-	-
Binarized Neural Networks	88.6	89.85	-	-
Binary Weight Networks	-	-	60.8	83.0
XNOR-Net	-	-	51.2	73.2





TTQ

- Learn the ternary assignments and ternary values.

- 三值量化
$$w_l^t = \begin{cases} W_l^p & : \tilde{w}_l > \Delta_l \\ 0 & : |\tilde{w}_l| \leq \Delta_l \\ -W_l^n & : \tilde{w}_l < -\Delta_l \end{cases}$$

- 求导：权值和缩放因子

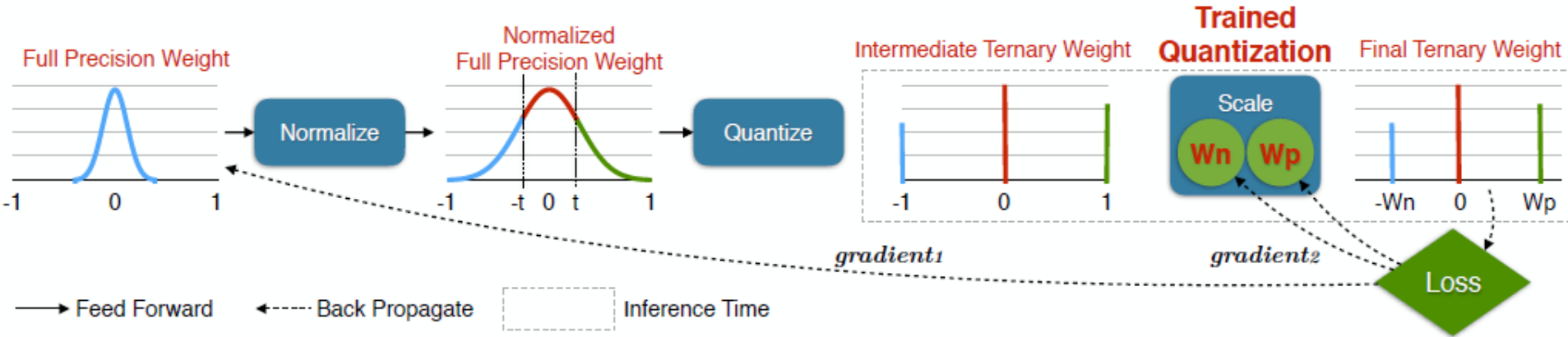
$$\frac{\partial L}{\partial W_l^p} = \sum_{i \in I_l^p} \frac{\partial L}{\partial w_l^t(i)}, \quad \frac{\partial L}{\partial W_l^n} = \sum_{i \in I_l^n} \frac{\partial L}{\partial w_l^t(i)}$$

$$\frac{\partial L}{\partial \tilde{w}_l} = \begin{cases} W_l^p \times \frac{\partial L}{\partial w_l^t} & : \tilde{w}_l > \Delta_l \\ 1 \times \frac{\partial L}{\partial w_l^t} & : |\tilde{w}_l| \leq \Delta_l \\ W_l^n \times \frac{\partial L}{\partial w_l^t} & : \tilde{w}_l < -\Delta_l \end{cases}$$

- 阈值选择

- 依据该层最大权值 $\Delta_l = t \times \max(|\tilde{w}|)$
- 保持一个固定的稀疏度r

TTQ



TTQ



Error	Full precision	1-bit (DoReFa)	2-bit (TWN)	2-bit (Ours)
Top1	42.8%	46.1%	45.5%	42.5%
Top5	19.7%	23.7%	23.2%	20.3%

Table 2: Top1 and Top5 error rate of AlexNet on ImageNet

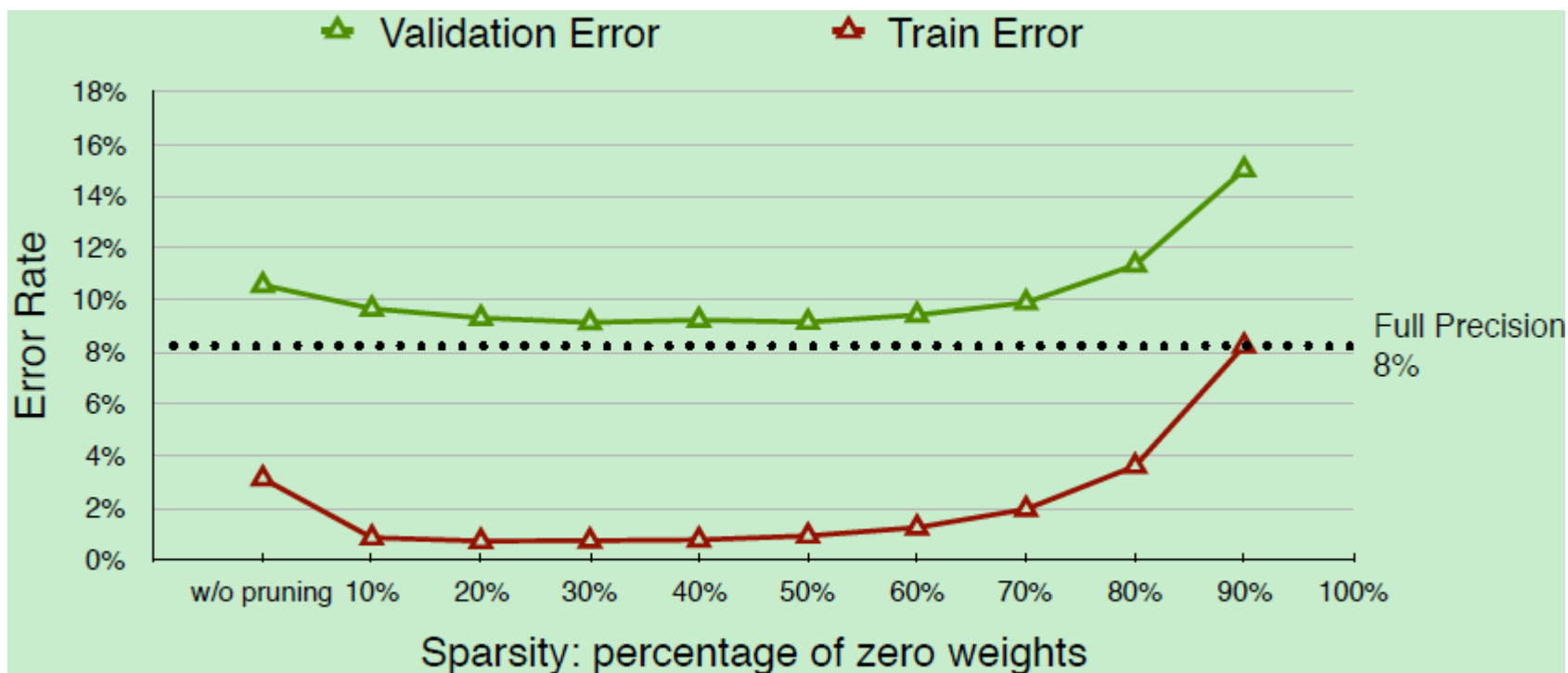


Figure 5: Accuracy v.s. Sparsity on ResNet-20



DoReFa-Net

- 在BNN和XNOR-Net中，梯度保持浮点型
- DoReFa-Net
 - 对权值、激活值、梯度均做量化
 - 量化到任意比特，

$$\text{Forward: } r_o = \frac{1}{2^k - 1} \text{round}((2^k - 1)r_i)$$

$$\text{Backward: } \frac{\partial c}{\partial r_i} = \frac{\partial c}{\partial r_o}$$

- STE直通估计：BNN，XNOR-Net

$$\text{Forward: } r_o = \text{sign}(r_i)$$

$$\text{Backward: } \frac{\partial c}{\partial r_i} = \frac{\partial c}{\partial r_o} \mathbb{I}_{|r_i| \leq 1}$$

$$\text{Forward: } r_o = \text{sign}(r_i) \times \mathbf{E}_F(|r_i|)$$

$$\text{Backward: } \frac{\partial c}{\partial r_i} = \frac{\partial c}{\partial r_o}$$



DoReFa-Net

- 二值化：一层只分配一个缩放因子

Forward: $r_o = \text{sign}(r_i) \times \mathbf{E}(|r_i|)$

Backward: $\frac{\partial c}{\partial r_i} = \frac{\partial c}{\partial r_o}$

- k比特量化

Forward: $r_o = f_{\omega}^k(r_i) = 2 \text{quantize}_k\left(\frac{\tanh(r_i)}{2 \max(|\tanh(r_i)|)} + \frac{1}{2}\right) - 1.$

Backward: $\frac{\partial c}{\partial r_i} = \frac{\partial r_o}{\partial r_i} \frac{\partial c}{\partial r_o}$ 4

- 激活值量化

1. $h(x) = \frac{\tanh(x)+1}{2}$

2. $h(x) = \text{clip}(x, 0, 1)$

3. $h(x) = \min(1, |x|)$

$f_{\alpha}^k(r) = \text{quantize}_k(r).$



DoReFa-Net

□ 量化梯度（随机方式）

$$f_{\gamma}^k(dr) = 2 \max_0(|dr|) \left[\text{quantize}_k \left[\frac{dr}{2 \max_0(|dr|)} + \frac{1}{2} + N(k) \right] - \frac{1}{2} \right].$$

$$N(k) = \frac{\sigma}{2^k - 1} \text{ where } \sigma \sim \text{Uniform}(-0.5, 0.5).$$

□ 合并非线性函数和量化函数

```
1: for  $k = 1$  to  $L$  do  
2:    $W_k^b \leftarrow f_{\omega}^W(W_k)$   
3:    $\tilde{a}_k \leftarrow \text{forward}(a_{k-1}^b, W_k^b)$   
4:    $a_k \leftarrow h(\tilde{a}_k)$   
5:   if  $k < L$  then  
6:      $a_k^b \leftarrow f_{\alpha}^A(a_k)$   
7:   end if
```

$$a_k^b = f_{\alpha}(h(a_k))$$



DoReFa-Net

Algorithm 1 Training a L -layer DoReFa-Net with W -bit weights and A -bit activations using G -bit gradients. Weights, activations and gradients are quantized according to [Eqn. 9](#), [Eqn. 11](#), [Eqn. 12](#), respectively.

Require: a minibatch of inputs and targets (a_0, a^*) , previous weights W , learning rate η

Ensure: updated weights W^{t+1}

```
{1. Computing the parameter gradients:}
{1.1 Forward propagation:}
1: for  $k = 1$  to  $L$  do
2:    $W_k^b \leftarrow f_\omega^W(W_k)$ 
3:    $\bar{a}_k \leftarrow \text{forward}(a_{k-1}^b, W_k^b)$ 
4:    $a_k \leftarrow h(\bar{a}_k)$ 
5:   if  $k < L$  then
6:      $a_k^b \leftarrow f_\alpha^A(a_k)$ 
7:   end if
8:   Optionally apply pooling
9: end for
{1.2 Backward propagation:}
Compute  $g_{a_L} = \frac{\partial C}{\partial a_L}$  knowing  $a_L$  and  $a^*$ .
10: for  $k = L$  to 1 do
11:   Back-propagate  $g_{a_k}$  through activation function  $h$ 
12:    $g_{a_k}^b \leftarrow f_\gamma^G(g_{a_k})$ 
13:    $g_{a_{k-1}} \leftarrow \text{backward\_input}(g_{a_k}^b, W_k^b)$ 
14:    $g_{W_k^b} \leftarrow \text{backward\_weight}(g_{a_k}^b, a_{k-1}^b)$ 
15:   Back-propagate gradients through pooling layer if there is one
16: end for
{2. Accumulating the parameters gradients:}
17: for  $k = 1$  to  $L$  do
18:    $g_{W_k} = g_{W_k^b} \frac{\partial W_k^b}{\partial W_k}$ 
19:    $W_k^{t+1} \leftarrow \text{Update}(W_k, g_{W_k}, \eta)$ 
20: end for
```



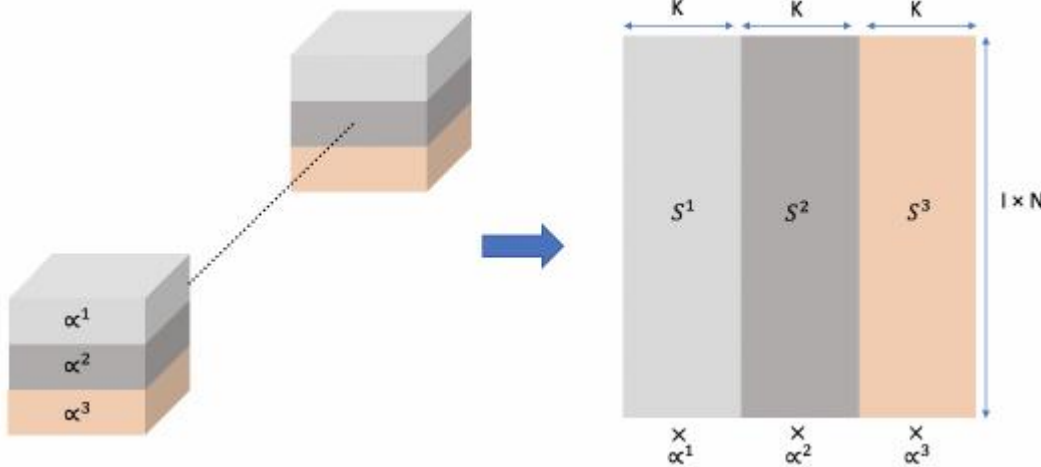
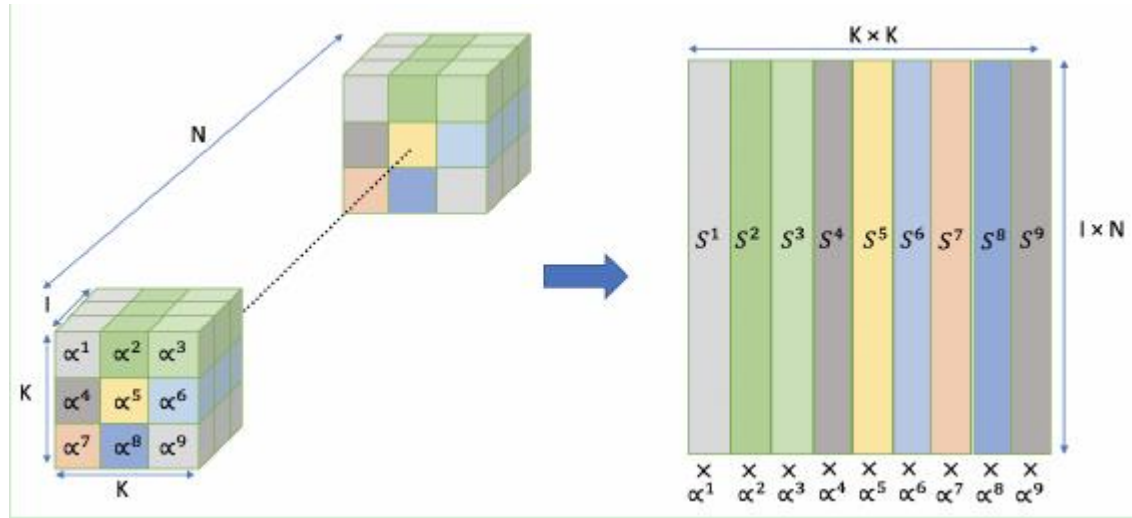
DoReFa-Net

W	A	G	Training Complexity	Inference Complexity	Storage Relative Size	AlexNet Accuracy
1	1	6	7	1	1	0.395
1	1	8	9	1	1	0.395
1	1	32	-	1	1	0.279 (BNN)
1	1	32	-	1	1	0.442 (XNOR-Net)
1	1	32	-	1	1	0.401
1	1	32	-	1	1	0.436 (initialized)
1	2	6	8	2	1	0.461
1	2	8	10	2	1	0.463
1	2	32	-	2	1	0.477
1	2	32	-	2	1	0.498 (initialized)
1	3	6	9	3	1	0.471
1	3	32	-	3	1	0.484
1	4	6	-	4	1	0.482
1	4	32	-	4	1	0.503
1	4	32	-	4	1	0.530 (initialized)
8	8	8	-	-	8	0.530
32	32	32	-	-	32	0.559

SYQ

□ 不一样的缩放因子

- Pixel-wise
- Row-wise





SYQ

□ 权值量化

$$\hat{w} = \text{sign}(w)$$

$$w_q = \text{diag}(\alpha)\hat{w}^T, \text{diag}(\alpha) \in \mathbb{R}^{d^2 \times d^2}$$

$$Q_l = \text{sign}(W_l) \odot M_l$$

$$M_{l_{i,j}} = \begin{cases} 1 & \text{if } |W_{l_{i,j}}| \geq \eta \\ 0 & \text{if } -\eta < W_{l_{i,j}} < \eta \end{cases}$$

□ 激活值量化同DoReFa-Net

$$\hat{a} = \text{clip}(a, 0, 1)$$

$$a_q = \frac{1}{2^k - 1} \text{round}((2^k - 1)\hat{a})$$

□ BP阶段

$$\frac{\partial E}{\partial w^i} = \frac{\partial E}{\partial w_q^i} \frac{\partial w_q^i}{\partial w^i} = \alpha^i \frac{\partial E}{\partial w_q^i}$$

$$\frac{\partial E}{\partial a} = \frac{\partial E}{\partial a_q} \frac{\partial \hat{a}}{\partial a} \quad \frac{\partial E}{\partial \alpha^i} = \sum_{j \in S^i} \frac{\partial E}{\partial w^j}$$

初始化

$$\alpha_{l_0}^i = \frac{\sum_{j \in S_l^i} |W_{l_{i,j}}|}{I \times N}$$



SYQ

Table 3. Comparison to previously published AlexNet results

Model	Weights	Act.	Top-1	Top-5
DoReFa-Net [33]	1	2	49.8	-
QNN [15]	1	2	51.0	73.7
HWGQ [2]	1	2	52.7	76.3
SYQ	1	2	55.4	78.6
DoReFa-Net [33]	1	4	53.0	-
SYQ	1	4	56.2	79.4
BWN [24]	1	32	56.8	79.4
SYQ	1	8	56.6	79.4
SYQ	2	2	55.8	79.2
FGQ [21]	2	8	49.04	-
TTQ [34]	2	32	57.5	79.7
SYQ	2	8	58.1	80.8



TBN

□ Ternary-Binary Network

	Methods	Inputs	Weights	MACs	Binary operations	Speedup	Operations
	Full-precision	\mathbb{R}	\mathbb{R}	$n \times m \times q$	0	$1 \times$	+, x
Quantize Weights	TTQ [60]	\mathbb{R}	$\{-\alpha^n, 0 + \alpha^p\}$	$n \times m \times q$	0	$\sim 2 \times$	+,-
	TWN [33]	\mathbb{R}	$\{-\alpha, 0, -\alpha\}$	$n \times m \times q$	0	$\sim 2 \times$	+,-
	BWN [42]	\mathbb{R}	$\{-\alpha, +\alpha\}$	$n \times m \times q$	0	$\sim 2 \times$	+,-
	BC [6]	\mathbb{R}	$\{-1, +1\}$	$n \times m \times q$	0	$\sim 2 \times$	+,-
Quantize Inputs and Weights	TNN [1]	$\{-1, 0, 1\}$	$\{-1, 0, 1\}$	0	$8 \times n \times m \times q$	$15 \times$	AND, bitcount
	GXNOR [9]	$\{-1, 0, 1\}$	$\{-1, 0, 1\}$	0	$5 \times n \times m \times q$	$15 \times$	AND, bitcount
	BNN [7]	$\{-1, +1\}$	$\{-1, +1\}$	0	$2 \times n \times m \times q$	$64 \times$	XOR, bitcount
	XNOR [42]	$\{-\beta, +\beta\}$	$\{-\alpha, +\alpha\}$	$2 \times n \times m$	$2 \times n \times m \times q$	$58 \times$	XOR, bitcount
	HORQ [34]	$\{-\beta, +\beta\} \times 2$	$\{-\alpha, +\alpha\}$	$4 \times n \times m$	$4 \times n \times m \times q$	$29 \times$	XOR, bitcount
	DoReFa* [59]	$\{0, 1\} \times 2$	$\{0, 1\}$	0	$4 \times n \times m \times q$	$30 \times$	AND, bitcount
	TBN	$\{-1, 0, +1\}$	$\{-\alpha, +\alpha\}$	$n \times m$	$3 \times n \times m \times q$	$40 \times$	AND, XOR, bitcount

*We adopt DoReFa Network with 1-bit weight, 2-bit activation.

TBN

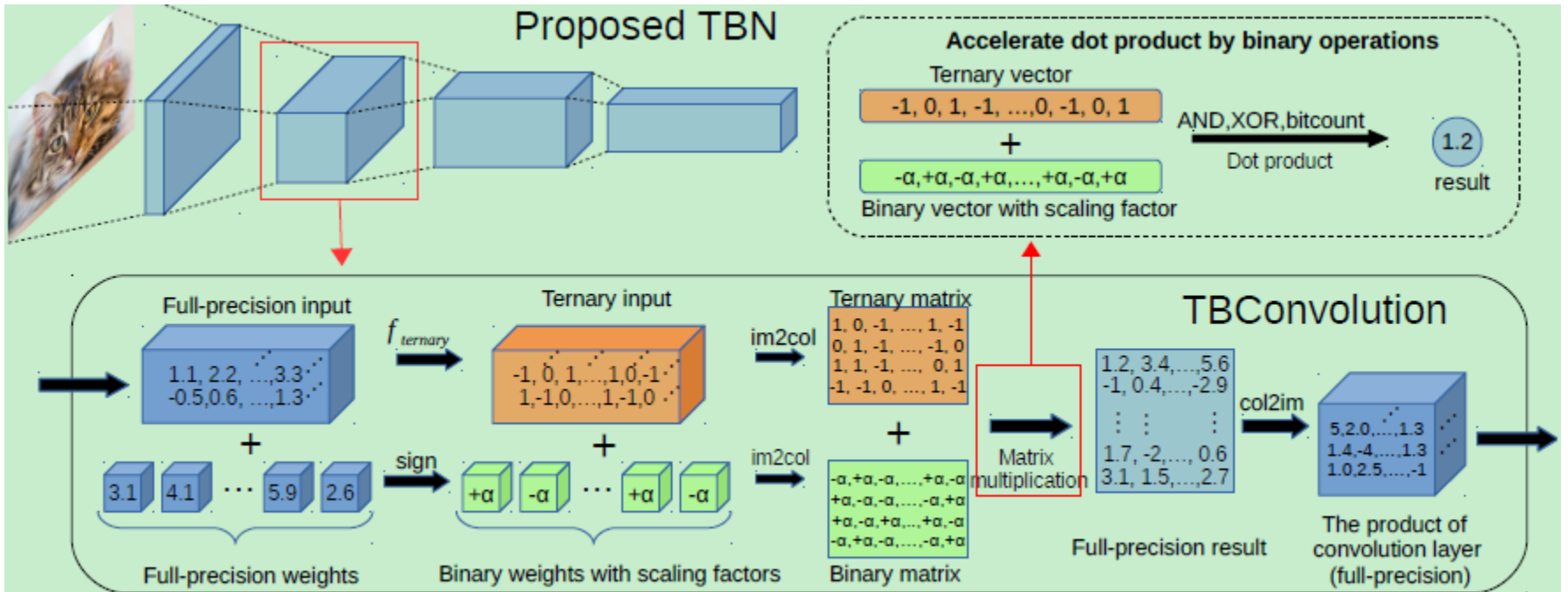
- 权值量化同XNOR-Net:
- 激活值量化像TWN:
- 求导同上:

$$\mathbf{B} = \text{sign}(\mathbf{W}), \quad \alpha = \frac{1}{c \times h \times w} \|\mathbf{W}\|_1.$$

$$\mathbf{T}_i = f_{\text{ternary}}(\mathbf{I}_i, \Delta) = \begin{cases} +1, & \mathbf{I}_i > \Delta; \\ 0, & |\mathbf{I}_i| \leq \Delta; \\ -1, & \mathbf{I}_i < -\Delta; \end{cases}$$

$$\frac{\partial \text{sign}}{\partial r} = \frac{\partial f_{\text{ternary}}}{\partial r} = \mathbf{1}_{|r| < 1} = \begin{cases} 1, & |r| < 1 \\ 0, & \text{otherwise} \end{cases}$$

$$\Delta = \delta \times \mathbf{E}(|\mathbf{I}|) \approx \frac{\delta}{c \times h_{in} \times w_{in}} \|\mathbf{I}\|_1$$





TBN

□ 卷积为矩阵乘 $C = \text{mat2ten}(\widetilde{\mathbf{W}}\widetilde{\mathbf{I}})$, $\widetilde{\mathbf{W}} = \text{ten2mat}(\mathbf{W})$, $\widetilde{\mathbf{I}} = \text{ten2mat}(\mathbf{I})$
($q = c \times h \times w$), $\widetilde{\mathbf{W}} \in \mathbb{R}^{n \times q}$, $\widetilde{\mathbf{I}} \in \mathbb{R}^{q \times m}$ ($m = h_{out} \times w_{out}$)

$$\widetilde{\mathbf{W}}_i \approx \alpha b, b = \text{ten2mat}(\mathbf{B}) \in \{-1, 1\}^q \quad t \in \{-1, 0, +1\}^q$$

□ 三值-二值向量点乘加速:

$$C_{ij} = \alpha(c_t - 2 \times \text{bitcount}((b \text{ XOR } t') \text{ AND } t'')),$$

□ 其中

$$t'_i = \begin{cases} 1, & t_i = 1 \\ -1, & \text{otherwise} \end{cases}, t''_i = \begin{cases} 0, & t_i = 0 \\ 1, & \text{otherwise} \end{cases}, i = 1, \dots, q$$

$$t_i = t'_i \times t''_i, c_t = \text{bitcount}(t'') = \|t\|_1$$

- XOR、AND均为逻辑运算
- 1在 b, t', t'' 中视为逻辑真, 0和-1则为逻辑假

TBN

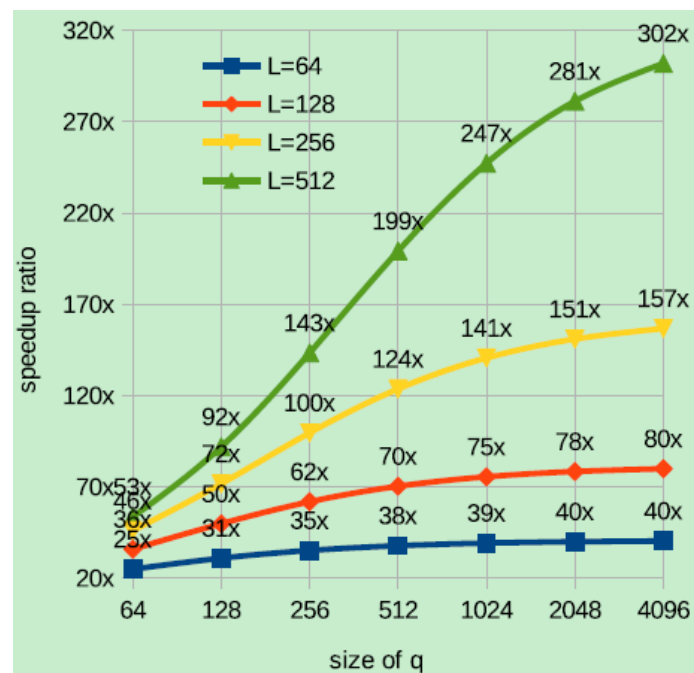
效率分析

- 原浮点运算 $\tilde{C} = \tilde{W}\tilde{I}$, $n \times m \times q$
- TBN运算: $n \times m$ 次浮点 + $n \times m \times q$ AND, XOR and bitcount
- 一个时钟执行L=64位二值运算, 设

$$\gamma = \frac{\text{average time required by MAC}}{\text{average time required by } L\text{-bits binary operation}}$$

- 加速比为

$$S = \frac{\gamma n m q}{\gamma n m + 3 n m \lceil \frac{q}{L} \rceil} = \frac{\gamma q}{\gamma + 3 \lceil \frac{q}{L} \rceil}$$





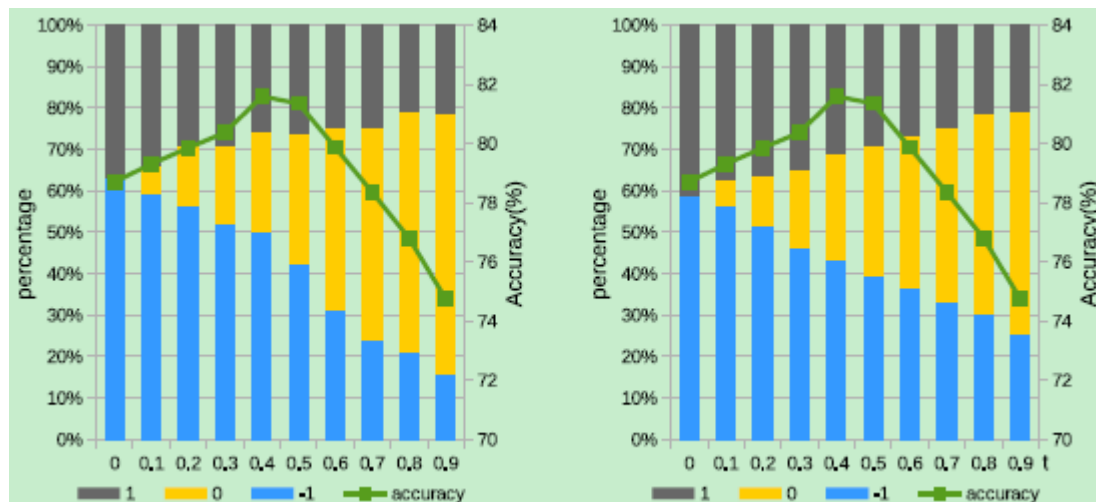
TBN

Dataset		MNIST	CIFAR-10	SVHN	ImageNet	ImageNet	ImageNet
Models		LeNet-5	VGG-7	VGG-7	AlexNet	ResNet-18	ResNet-34
	Full-precision	99.48	92.88	97.68	57.2/80.2	69.3/89.2	73.3/91.4
Quantize Weights	BC [6]	98.82	91.73	97.85	35.5/61.0	-	-
	BWN [42]	99.38	92.58	97.46	56.8/79.4	60.8/83.0	-
	TWN [33]	99.38	92.56	-	54.5/76.8	65.3/86.2	-
	TTQ [60]	-	-	-	57.5/79.7	66.6/87.2	-
Other Methods	FFN [53]	-	-	-	55.5/79.0	-	-
	LCNN-fast [3]	-	-	-	44.3/68.7	51.8/76.8	-
	LCNN-accurate [3]	-	-	-	55.1/78.1	62.2/84.6	-
	LBCNN [24]	99.51	92.66	94.50	54.9/-	-	-
Quantize Inputs and Weights	TNN [1]	98.33	87.89	97.27	-	-	-
	GXNOR [9]	99.32	92.50	97.37	-	-	-
	BNN [7]	98.60	89.85	97.47	27.9/50.42	-	-
	DoReFa-Net* [59]	-	-	97.6	47.7/-	-	-
	BinaryNet [51]	-	-	-	46.6/71.1	-	-
	HORQ [34]	99.38	91.18	97.41	-	55.9/78.9	-
	XNOR-Network [42]	99.21	90.02	96.96	44.2/69.2	51.2/73.2	55.9/79.1
	TBN	99.38	90.85	97.27	49.7/74.2	55.6/79.0	58.2/81.0

*We adopt DoReFa-Net with 1-bit weight, 2-bit activation and 32-bit gradient for fair comparison.

TBN

稀疏度 $\delta = 0.4$



检测

method	full-precision VGG-16	full-precision ResNet-34	XNOR-Networks ResNet-34	TBN ResNet-34
Faster R-CNN	73.2	75.6	54.7	59.0
SSD 300	74.3	75.5	55.1	59.5



谢谢!