

Week 4 Report

Qian Jiang

1. Summarized optimization problems and methods for quantization neural networks

- 1) Training Quantized Nets: A Deeper Understanding

Sum: Problem of low-precision weight: stochastic gradient
→ Small updates of weights → easy to be rounded off by binarization.

This paper did 1) convergence analysis for classical Stochastic Rounding(SR) & Binary Connect(BC), showing they both are capable of solving convex discrete problems
2) explain why BC works better than SR for non-convex problem

- 2) Deep Learning with Limited Numerical Precision

Sum: Idea is to use algorithm-level noise tolerance to leverage hardware requirements → Co-optimized system

This paper did 1) 16-bit fixed point number representation using stochastic rounding 2) Proposed a hardware implementation

- 3) Training and Inference with Integers in Deep Neural Networks

Sum: Problem: how to quantize both operands and operations & how many bits are needed for SGD

This paper proposed WAGE method to discretize both training and inference: 1) For operands, linear mapping & orientation-preserved shifting; For operations, replace batch normalization by constant scaling factor and simplify unnecessary techniques 2) For inference, accumulate-compare cycles; For training, low-bitwidth multiply-accumulate cycles

(Detailed notes attached below)

2. Read code of [HWGQ](#)

Set up environments and learn caffe

Training Quantized Nets: A Deeper Understanding

convex optimization where all the constraints are convex

Key: training methods for QNN

- accuracy guarantee for methods under convexity assumptions
- explore algorithm behavior for non-convex problem

functions 1. Intro

(eg 

Problem of low-precision weights:

stochastic gradient method \rightarrow small updates of w

\Downarrow
easy to be rounded off
by binarization/
discretization

this paper

① convergence analysis: classical stochastic rounding

\Downarrow
Binary Connect

Capable of solving convex discrete problems
(depends on quantization level)

② why algorithm "maintain floating-point (like BC) work well"

\Downarrow
2. fully quantized (like SR)
stall before training is complete

long-term behavior of BC has important property
needed for non-convex optimization
which SR lacks.

2. Background & Related Work

backprop =
back
propagation

{ quantizing pre-trained models with / without retraining
training quantized model from scratch
↳ what this paper focus on

3. Training QNN

$$\min F(w) := \frac{1}{m} \sum_{i=1}^m f_i(w) \quad (1)$$

$$GSD: w^{t+1} = w^t - \alpha_t \nabla \tilde{f}(w^t) \quad (2)$$

$$\hat{w} = Q(w)$$

⇒ ① Deterministic Rounding:

$$Q_d(w) = \text{sign}(w) \cdot \Delta \cdot \left[\frac{|w|}{\Delta} + \frac{1}{2} \right] \quad (3)$$

$$W_b^{t+1} = Q_d(W_b^t - \alpha_t \nabla \tilde{f}(W_b^t)) \quad (4)$$

② Stochastic Rounding: random number

$$Q_s(w) = \Delta \cdot \begin{cases} \left\lfloor \frac{|w|}{\Delta} \right\rfloor + 1, & \text{for } p \leq \frac{w}{\Delta} - \left\lfloor \frac{w}{\Delta} \right\rfloor \\ \left\lfloor \frac{w}{\Delta} \right\rfloor, & \text{otherwise} \end{cases} \quad (5)$$

$$\underset{\text{expectation}}{\downarrow} E[Q_s(w)] = w$$

$$W_b^{t+1} = Q_s(W_b^t - \alpha_t \nabla \tilde{f}(W_b^t)) \quad (6)$$

③ Binary Connect:

full-precision buffer W_r

$$W_r^{t+1} = W_r^t - \alpha_t \nabla \tilde{f}(Q(W_r^t)) \quad (7)$$

quantizes weights just before the gradient computation

Deep learning with Numerical Precision

Key: ① 1b-bit fixed point number representation.

low-precision

using stochastic rounding

: less than

\Rightarrow little degradation in classification acc

8-bits

1. Intro

② Hardware.

idea: algorithm-level noise tolerance

↑ leverage
hardware requirements

\Rightarrow co-optimized system

allowing low-level hardware components

to perform approximate, non-deterministic
computations

Step 1: low-precision fixed-point arithmetic

2: stochastic rounding scheme

3: test on MNIST, CIFAR-10

4: propose a hardware accelerator on FPGA

2. Related Works

① implementing just the feed-forward.

② training using different number representations
(most recent 32-bit, this paper 1b-bit)

•

3. Limited Precision Arithmetic

$[Q_I \cdot Q_F]$
integer \rightarrow fraction

$$IL + FL = WL$$

E: smallest positive
number can be
represented
 $= 2^{-FL}$

① Rounding Modes

1. Round to nearest

$$= \begin{cases} \lfloor x \rfloor & , \text{ if } \lfloor x \rfloor \leq x \leq \lfloor x \rfloor + \frac{\epsilon}{2} \\ \lfloor x \rfloor + \epsilon & , \text{ if } \lfloor x \rfloor + \frac{\epsilon}{2} < x \leq \lfloor x \rfloor + \epsilon \end{cases}$$

2. Stochastic rounding

$$= \begin{cases} \lfloor x \rfloor & , \text{ w.p. } (1 - \frac{x - \lfloor x \rfloor}{\epsilon}) \\ \lfloor x \rfloor + \epsilon & , \text{ w.p. } (\frac{x - \lfloor x \rfloor}{\epsilon}) \end{cases}$$

$$(E(\text{Round}(x, \langle IL, FL \rangle))) = x$$

$$\text{Convert}(x, \langle IL, FL \rangle) = \begin{cases} -2^{IL-1} & , \text{ if } x \leq 2^{IL-1} \\ 2^{IL-1} - 2^{-FL} & , \text{ if } x \geq 2^{IL-1} - 2^{-FL} \\ \text{Round}(x, \langle IL, FL \rangle), \text{ otherwise} \end{cases}$$

② Multiply and accumulate (MACC) operation

$$c_0 = a_k \cdot b \quad \text{step 1: } z = \sum_{i=1}^d a_i b_i \\ [\tilde{IL}, \tilde{IF}] \quad [IL, IF] \quad (L \text{ of } z \leq \log_2 d + 2WL)$$

step 2: Convert

$$c_0 = \text{Convert}(z, \langle \tilde{IL}, \tilde{IF} \rangle)$$

Pros:

1. mimic the hardware behavior

2. reduce hardware overhead.

3. efficiently simulate fixed-point computations

using CPU/GPUs and vendor-supplied BLAS

4. Training Deep Networks

constrained parameters: $W^L, B^L, Y^L, S^L, \Delta W^L, \Delta B^L$

unchanged: W initialization, learning rate ...

Training and Inference with Integers in Deep NNs.

- Key:
1. WAGE : discretize both training and inference
 2. replace batch normalization by constant scaling layer

training
↓

inference 1. Intro

① how to quantize all the operands and operations

② how many bits or states are needed for SGD

WAGE : weight. Activation. Gradient. Error

step 1: linear mapping. orientation-preserved shifting
are applied for operands → ternary weights

8-bit integers
for A · G

2: batch normalization replaced by constant scaling factor
for operations.

3. Simplify other techniques (L2 regularization...)

Streamline: inference → accumulate - compare
training → low-bitwidth multiply - accumulate

2. Related Work

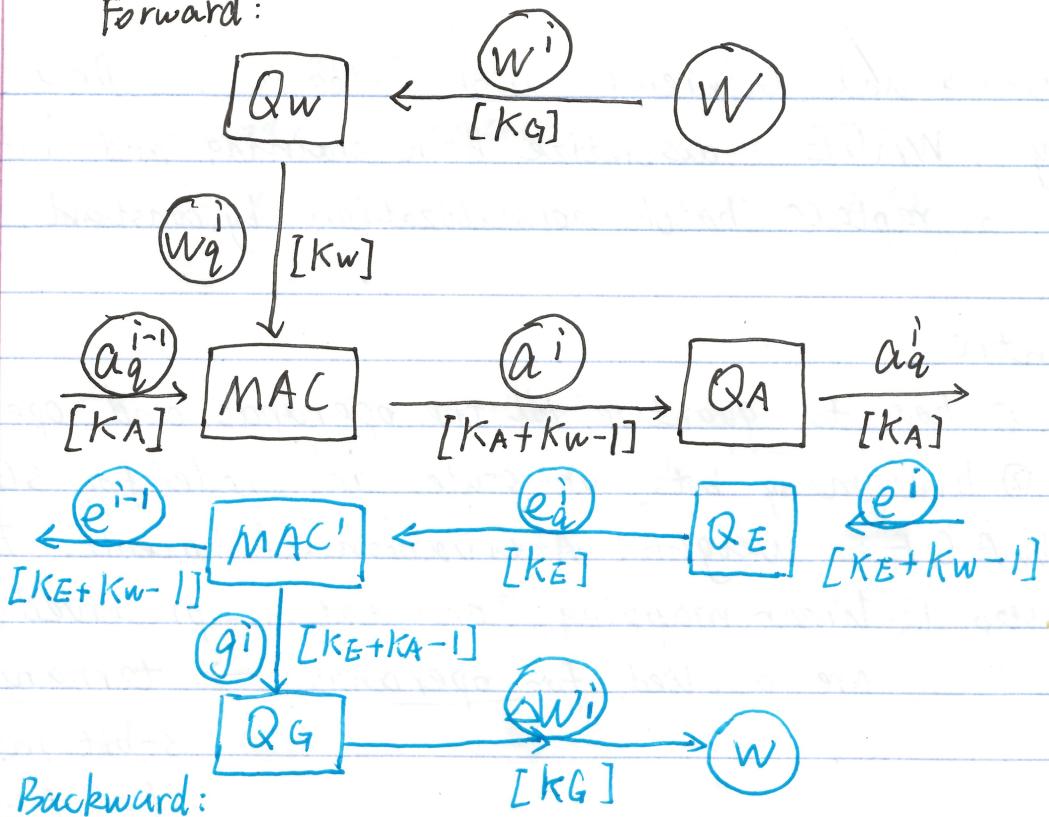
① Weight and activation: BC · BNN

② Gradient computation and accumulation: DoReFa-Net

3. WAGE Quantization

$$e^i = \frac{\partial L}{\partial a^i} \cdot g^i = \frac{\partial L}{\partial w^i} \quad (L: \text{Loss function})$$

Forward:



Backward:

operator \mathcal{Q} : reduce precision

① Linear-mapping & Stochastic Rounding

$$\mathcal{D}(k) = 2^{1-k}, k \in N^+$$

$$\mathcal{Q}(x, k) = \text{clip}\left\{\frac{x}{\mathcal{D}(k)}, \text{round}\left[\frac{x}{\mathcal{D}(k)}\right], -1 + \mathcal{D}(k), 1 - \mathcal{D}(k)\right\}$$

$$\text{Shift}(x) = 2^{\text{round}(\log_2 x)}$$

② Weight Initialization

$$W \sim U(-L, +L), L = \max\{\sqrt{b/n_{in}}, L_{min}\}, \underline{L_{min}}$$

$$(L_{min} = \beta \sigma)$$

③ Quantization Details

$$W_q = Q_W(W) = \mathcal{Q}(W, k_w) \quad (\text{scaling factor } \alpha = \max\{\text{shift} \frac{L_{min}}{L}, 1\})$$

$$a_q = Q_A(a) = \mathcal{Q}\left(\frac{a}{\alpha}, k_a\right) \quad e_q = Q_E(e) = \mathcal{Q}\left(\frac{e}{\text{shift}(\max\{|e|\})}, k_e\right)$$

$$g_s = \eta \cdot g / \text{shift}(\max\{|g|\}) \quad \Delta W = W_{t+1} - \dots$$