

# Week 3 Report

## Qian Jiang

1. Backpropagation of quantization neural network  
(Detailed notes attached below)
2. Deep Learning with Limited Numerical Precision [2015 ICML]

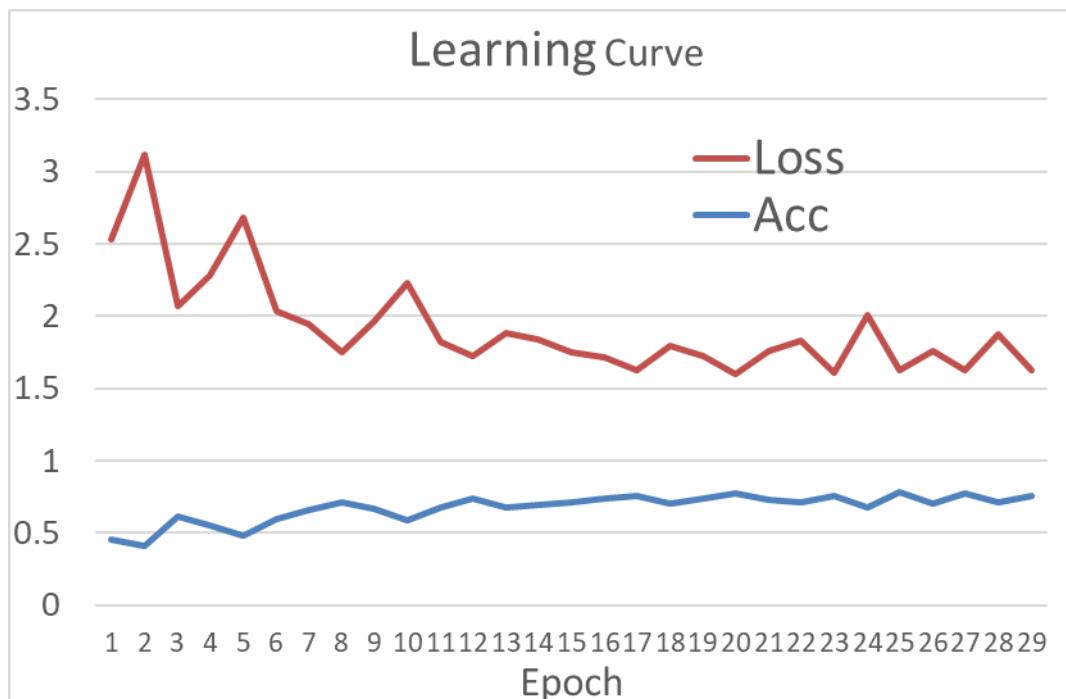
Summary:

- 1) Problem: Large model size& computational cost  
→ Solu: HWGQ in feed forward
- 2) Problem: Step-wise response leads to weak gradient when backward. → Solu: Using continuous approximation of operators( ReLU or tanh)
- 3) Problem: mismatch in forward and backward → Solu: Clipped ReLU & Log-tailed ReLU when backward

(Detailed notes attached below)

3. Learning Curve of Cifar\_10

(Currently unable to run on server, so for now only 30 epochs)



## Backpropagation of XNOR

epoch:

one forward pass

+ one back pass

batch-size:

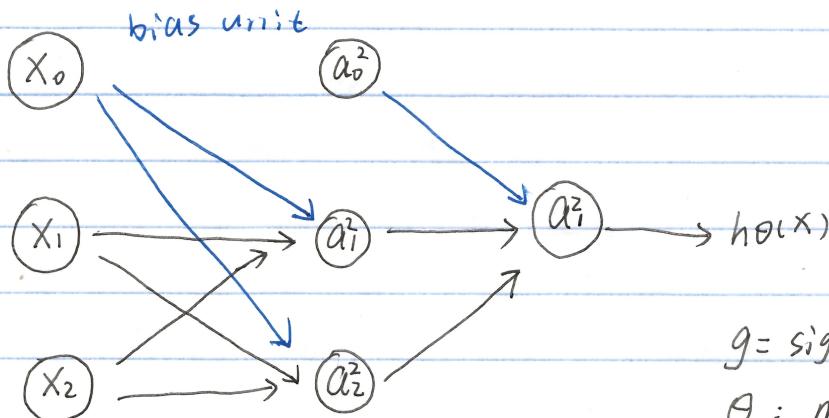
number of

examples in  
one pass

iteration:

number of

passes



$$g = \text{sigmoid} = \frac{1}{1+e^{-z}}$$

$\theta$ : matrix of weights

$$a_0^2 = g(\theta_{00}^{(1)} x_0 + \theta_{01}^{(1)} x_1 + \theta_{02}^{(1)} x_2) = g(\theta_0^T x) = g(z_0^2)$$

$$a_1^2 = g(\theta_{10}^{(1)} x_0 + \theta_{11}^{(1)} x_1 + \theta_{12}^{(1)} x_2) = g(\theta_1^T x) = g(z_1^2)$$

$$a_2^2 = g(\theta_{20}^{(1)} x_0 + \theta_{21}^{(1)} x_1 + \theta_{22}^{(1)} x_2) = g(\theta_2^T x) = g(z_2^2)$$

$$h_\theta(x) = a_1^3 = g(\theta_{10}^{(2)} a_0^2 + \theta_{11}^{(2)} a_1^2 + \theta_{12}^{(2)} a_2^2)$$

$$g = \frac{1}{1+e^{-z}} \quad g' = \frac{e^{-z}}{(1+e^{-z})^2} = g(1-g)$$

Forward Propagation:

$$x = \begin{pmatrix} x_0 \\ x_1 \\ x_2 \end{pmatrix} \quad z^2 = \begin{pmatrix} z_0^2 \\ z_1^2 \\ z_2^2 \end{pmatrix}$$

$$z^2 = \theta^{(1)} x = \theta^{(1)} \cdot a^{(1)} \quad a^2 = g(z^2)$$

$$z^3 = \theta^{(2)} a^2 \quad h_\theta(x) = g(z^3)$$

Back Propagation:

$$\delta_j^3 = a_j^3 - y_j \Rightarrow \delta^3 = a^3 - y$$

$$\delta^2 = (\theta^2)^T \cdot \delta^3 g'(z^2)$$

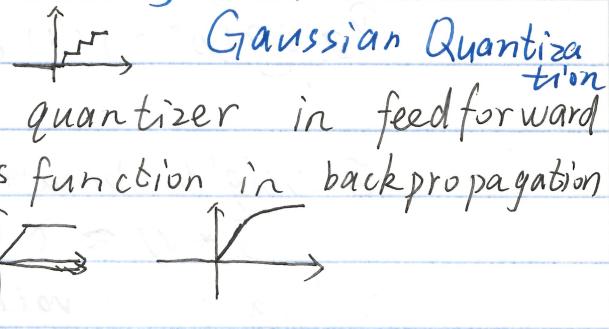
$$= (\theta^2)^T \cdot \delta^3 \cdot \cancel{z^2} \cdot (1 - \cancel{a^2})$$

$$\frac{\partial J(\theta)}{\partial \theta_{ij}^L} = a_j^L \delta_i^{L+1}$$

# Deep Learning with Low Precision by Half-wave

piecewise  
function  
= hybrid  
function

- Key:
- ① half-wave Gaussian quantizer in feedforward
  - ② piece-wise continuous function in backpropagation



hyperbolic

- tangent

1.

Intro

tanh

Problem 1. large model size



2. Large computational cost

↓ Soln:

compressed models : [quantization, low-rank matrix factorization, pruning, architecture design]

binarization of activations:

problem: step-wise response  $\rightarrow$  weak gradient

↓ Soln: using continuous approximation of operators

$\hookrightarrow$  problem: mismatch

↓

Soln: view quantization operator to be 2 functions.

ReLu  $\gg$  tanh: stronger gradient magnitude

2. linearization, gradient clipping, gradient suppression

Related work:

1. redundancy of weights ... (see above [ ] )

2. weight binarization/quantization

3. activation quantization { speed up ↑  
memory saving }

### 3. Binary Networks (floating-point)

①  $Z = g(W^T X)$  (1) complexity: 1. memory to store  $W$   
 2. dot-product  $W^T X$

② weight binarization

$$I * W \approx \alpha(I \oplus B) \quad (2)$$

③ binary activation quantization.

$$Z = \text{sign}(x) = \begin{cases} +1 & \text{if } x \geq 0 \\ -1 & \text{otherwise} \end{cases} \quad (3) \Rightarrow \begin{array}{l} \text{simplify computation} \\ \text{increase difficulty of learning} \end{array}$$

$$\frac{\partial C}{\partial W} = \frac{\partial C}{\partial Z} g'(W^T X) \quad (4)$$

↓ if  $g = \text{sign}$ ,  $g'(\dots) = 0$  everywhere

$$\text{soln: } g = \text{hard tanh} \quad g'(c) = \begin{cases} 1 & |x| \leq 1 \\ 0 & \text{otherwise} \end{cases} \quad (5)$$

↓ problem: gradient mismatch

### 4. HWGQ

$$\text{① ReLU: } g(x) = \max(0, x) \quad (6)$$

$$\text{② Quantizer: } Q(x) = q_i, \text{if } x \in [t_i, t_{i+1}] \quad (7)$$

$$\text{Quantization step: } q_{i+1} - q_i = \Delta, \forall i \quad (8)$$

$$Q^*(x) = \arg \min_Q \mathbb{E}_x [L(Q(x) - x)^2]$$

$$= \arg \min_Q \int_Q p(x) (Q(x) - x)^2 dx \quad (9)$$

Lloyd's  
algorithm

dot-products tend to be Gaussian distribution

ReLU is half wave.  $\rightarrow$  optimal parameters for Gaussian distribution

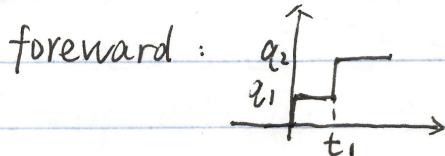
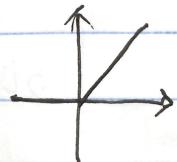
$$\Rightarrow Q(x) = \begin{cases} q_i & \text{if } x \in [t_i, t_{i+1}] \\ 0 & x \leq 0 \end{cases} \quad (10)$$

batch normalization: force layer to have zero mean, unit variance

③ Backward.

Vanilla ReLu:  $\tilde{Q}'(x) = \begin{cases} 1, & \text{if } x > 0 \\ 0, & \text{otherwise} \end{cases}$

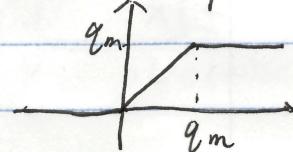
(11)



$\Rightarrow$  mismatch (large on the tail)

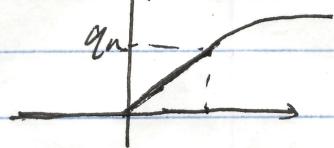
$\downarrow$   
soln: ① Clipped ReLU

$$Q' = \begin{cases} 0 & \\ 1 & \\ 0 & \end{cases} \quad \tilde{Q}_c(x) = \begin{cases} q_m, & x > q_m \\ x, & x \in [0, q_m] \\ 0, & \text{otherwise} \end{cases} \quad (\text{con: loss of info})$$



② Log-tailed ReLU

$$Q' = \begin{cases} \frac{1}{x-\tau} & \\ 1 & \\ 0 & \end{cases} \quad \tilde{Q}_l(x) = \begin{cases} q_m + \log(x-\tau), & x > q_m \\ x, & x \in [0, q_m] \\ 0, & \text{otherwise} \end{cases}$$



## 5. Experimental Results

- ① Full-precision Activation Comparison
- ② Low-bit Activation Quantization Results
- ③ Backward Approximation Comparison
- ④ Bit-width Impact
- ⑤ Comparison with the state-of-the-art