

# Loss-Aware Binarization Of Deep Networks

Credit: Lu Hou, Quanming Yao, James T. Kwok

Hong Kong University of Science and Technology

Presenter: Ryan McCampbell

<https://qdata.github.io/deep2Read>

# Outline

- 1 Introduction
- 2 Prior Work
- 3 Loss-Aware Binarization
- 4 Results
- 5 Conclusion

# Outline

- 1 Introduction
- 2 Prior Work
- 3 Loss-Aware Binarization
- 4 Results
- 5 Conclusion

# Introduction

- Deep neural networks are expensive in space and time
  - Training
  - Inference
  - Weight storage
- Not good for embedded systems and mobile devices
- How can we make them more efficient?

- Train neural network then compress
- Train and compress network simultaneously
  - Tensor decomposition
  - Parameter quantization
  - Binarization

# Binarization

- Store only one bit for each weight value
- Significantly reduces storage
- Eliminates most multiplications in forward pass

# Outline

- 1 Introduction
- 2 **Prior Work**
- 3 Loss-Aware Binarization
- 4 Results
- 5 Conclusion

- BinaryConnect - Transform each weight to  $\pm 1$  with the sign function

$$\text{Binarize}(w_j^t) = \text{sign}(w_j^t)$$

- Binary Weight Network - also learn scaling parameter

$$\text{Binarize}(w_j^t) = \alpha_j^t \mathbf{b}_j^t$$
$$\alpha_j^t = \frac{\|w_j^t\|_1}{n_j}, \quad \mathbf{b}_j^t = \text{sign}(w_j^t)$$



# Prior Work

- Use binarized weights for forward and backward pass
- Use full-precision weights in update, then re-binarize

- Use binarized weights for forward and backward pass
- Use full-precision weights in update, then re-binarize
- **Problem:** This doesn't consider effect of binarization on loss

# Outline

- 1 Introduction
- 2 Prior Work
- 3 Loss-Aware Binarization**
- 4 Results
- 5 Conclusion

# Proximal Newton Algorithm

- Proximal Newton algorithm: composite optimization problems where  $g$  may not be smooth

$$\min_x f(x) + g(x)$$

- Use second-order information: approximate Hessian  $\mathbf{H}$  of  $f$

$$x_{t+1} = \operatorname{argmin}_x \nabla f(x_t)^T (x - x_t) + (x - x_t)^T \mathbf{H} (x - x_t) + g(x)$$

# Loss-Aware Binarization

- Minimize  $\ell(\hat{\mathbf{w}})$  given constraints

$$\hat{\mathbf{w}}_l = \alpha_l \mathbf{b}_l, \alpha_l > 0, b_l \in \{\pm 1\}_l^n$$

- Define  $I_C(\hat{\mathbf{w}})$  as indicator function, 1 if constraints hold and  $\infty$  otherwise.

$$\min \ell(\hat{\mathbf{w}}) + I_C(\hat{\mathbf{w}})$$

- Use proximal Newton method:

$$\min_{\hat{\mathbf{w}}^t} \nabla \ell(\hat{\mathbf{w}}^{t-1})^T (\hat{\mathbf{w}}^t - \hat{\mathbf{w}}^{t-1}) + (\hat{\mathbf{w}}^t - \hat{\mathbf{w}}^{t-1})^T \mathbf{D}^{t-1} (\hat{\mathbf{w}}^t - \hat{\mathbf{w}}^{t-1})$$

given the constraints, where  $\mathbf{D}$  is an approximate diagonal Hessian matrix.

# Solution

- Let  $\mathbf{d}_l^{t-1} = \text{diag}(\mathbf{D}_l^{t-1})$

$$w_l^t = \hat{w}_l^{t-1} - \nabla_l \ell(\hat{w}^{t-1}) / \mathbf{d}_l^{t-1}$$

$$\alpha_l^t = \frac{\|\mathbf{d}_l^{t-1} \circ w_l^t\|_1}{\|\mathbf{d}_l^{t-1}\|_1}, \quad \mathbf{b}_l^t = \text{sign}(w_l^t)$$

- Each iteration: gradient descent with adaptive learning rate  $1/\mathbf{d}_l^{t-1}$
- Project to binary

# Algorithm

**Input:** Minibatch  $\{(\mathbf{x}_0^t, \mathbf{y}^t)\}$ , current full-precision weights  $\{\mathbf{w}_l^t\}$ , first moment  $\{\mathbf{m}_l^{t-1}\}$ , second moment  $\{\mathbf{v}_l^{t-1}\}$ , and learning rate  $\eta^t$ .

- 1: **Forward Propagation**
- 2: **for**  $l = 1$  to  $L$  **do**
- 3:  $\alpha_l^t = \frac{\|\mathbf{d}_l^{t-1} \odot \mathbf{w}_l^t\|_1}{\|\mathbf{d}_l^{t-1}\|_1}$ ;
- 4:  $\mathbf{b}_l^t = \text{sign}(\mathbf{w}_l^t)$ ;
- 5: rescale the layer- $l$  input:  $\tilde{\mathbf{x}}_{l-1}^t = \alpha_l^t \mathbf{x}_{l-1}^t$ ;
- 6: compute  $\mathbf{z}_l^t$  with input  $\tilde{\mathbf{x}}_{l-1}^t$  and binary weight  $\mathbf{b}_l^t$ ;
- 7: apply batch-normalization and nonlinear activation to  $\mathbf{z}_l^t$  to obtain  $\mathbf{x}_l^t$ ;
- 8: **end for**
- 9: compute the loss  $\ell$  using  $\mathbf{x}_L^t$  and  $\mathbf{y}^t$ ;
- 10: **Backward Propagation**
- 11: initialize output layer's activation's gradient  $\frac{\partial \ell}{\partial \mathbf{x}_L^t}$ ;
- 12: **for**  $l = L$  to 2 **do**
- 13: compute  $\frac{\partial \ell}{\partial \mathbf{x}_{l-1}^t}$  using  $\frac{\partial \ell}{\partial \mathbf{x}_l^t}$ ,  $\alpha_l^t$  and  $\mathbf{b}_l^t$ ;
- 14: **end for**

# Algorithm

- 15: **Update parameters using Adam**
- 16: **for**  $l = 1$  to  $L$  **do**
- 17:   compute gradients  $\nabla_l \ell(\hat{\mathbf{w}}^t)$  using  $\frac{\partial \ell}{\partial \mathbf{x}_i^t}$  and  $\mathbf{x}_{l-1}^t$ ;
- 18:   update first moment  $\mathbf{m}_l^t = \beta_1 \mathbf{m}_l^{t-1} + (1 - \beta_1) \nabla_l \ell(\hat{\mathbf{w}}^t)$ ;
- 19:   update second moment  $\mathbf{v}_l^t = \beta_2 \mathbf{v}_l^{t-1} + (1 - \beta_2) (\nabla_l \ell(\hat{\mathbf{w}}^t) \odot \nabla_l \ell(\hat{\mathbf{w}}^t))$ ;
- 20:   compute unbiased first moment  $\hat{\mathbf{m}}_l^t = \mathbf{m}_l^t / (1 - \beta_1^t)$ ;
- 21:   compute unbiased second moment  $\hat{\mathbf{v}}_l^t = \mathbf{v}_l^t / (1 - \beta_2^t)$ ;
- 22:   compute current curvature matrix  $\mathbf{d}_l^t = \frac{1}{\eta^t} (\epsilon \mathbf{1} + \sqrt{\hat{\mathbf{v}}_l^t})$ ;
- 23:   update full-precision weights  $\mathbf{w}_l^{t+1} = \mathbf{w}_l^t - \hat{\mathbf{m}}_l^t \odot \mathbf{d}_l^t$ ;
- 24:   update learning rate  $\eta^{t+1} = \text{UpdateRule}(\eta^t, t + 1)$ ;
- 25: **end for**



# Further Optimization

- An additional optimization can be obtained by binarizing the activation function
- This reduces additions and multiplications to XNOR bitwise operations

# Outline

- 1 Introduction
- 2 Prior Work
- 3 Loss-Aware Binarization
- 4 Results**
- 5 Conclusion

# Results

Table 1: Test error rates (%) for feedforward neural network models.

		<i>MNIST</i>	<i>CIFAR-10</i>	<i>SVHN</i>
(no binarization)	full-precision	1.190	11.900	2.277
(binarize weights)	BinaryConnect	1.280	<b>9.860</b>	2.450
	BWN	1.310	10.510	2.535
	LAB	<b>1.180</b>	10.500	<b>2.354</b>
(binarize weights and activations)	BNN	1.470	12.870	3.500
	XNOR	1.530	12.620	3.435
	LAB2	<b>1.380</b>	<b>12.280</b>	<b>3.362</b>

# Results

Table 2: Test error rates (%) on *SVHN*, for CNNs with different numbers of filters. Number in brackets is the difference between the errors of the binarized scheme and the full-precision network.

	K = 16	K = 32	K = 64	K = 128
full-precision	2.738	2.585	2.277	2.146
BinaryConnect	3.200 (0.462)	2.777 (0.192)	2.450 (0.173)	2.315 (0.169)
BWN	3.119 (0.461)	2.743 (0.158)	2.535 (0.258)	2.319 (0.173)
LAB	<b>3.050</b> (0.312)	<b>2.742</b> (0.157)	<b>2.354</b> (0.077)	<b>2.200</b> (0.054)

Table 3: Testing cross-entropy values of LSTM.

		<i>War and Peace</i>	<i>Linux Kernel</i>
(no binarization)	full-precision	1.268	1.329
	BinaryConnect	2.942	3.532
(binarize weights)	BWN	1.313	1.307
	LAB	<b>1.291</b>	<b>1.305</b>
	BNN	3.050	3.624
(binarize weights and activations)	XNOR	1.424	1.426
	LAB2	<b>1.376</b>	<b>1.409</b>

# Outline

- 1 Introduction
- 2 Prior Work
- 3 Loss-Aware Binarization
- 4 Results
- 5 Conclusion**

# Conclusion

- This method has similar performance to the full-precision networks for many tasks
- Enables efficient storage and evaluation of deep neural networks
- However it is only slightly better in some cases than BWN
- Binarizing activations reduces performance