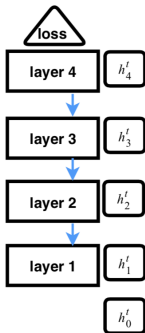


# Training Neural Networks Using Features Replay

**Zhouyuan Huo**<sup>1</sup>, Bin Gu<sup>1,2</sup>, Heng Huang<sup>1,2</sup>

<sup>1</sup>Department of Electrical and Computer Engineering, University of Pittsburgh  
<sup>2</sup> JD.com

November 28, 2018



Method

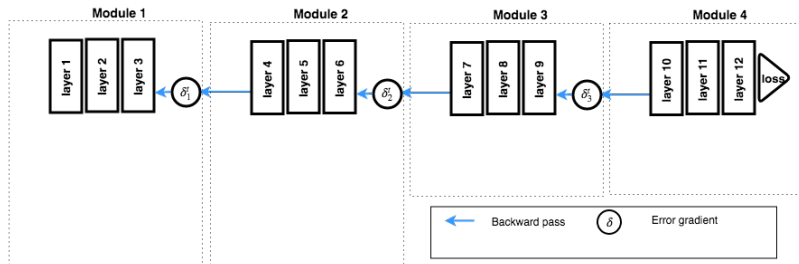
BP

Backpropagation algorithm:

- step 1: Forward pass.
- step 2: Backward pass.

Problem:

- Backward time is about **2** times of forward time.
- Backward locking.
- Backward cannot be parallelized.

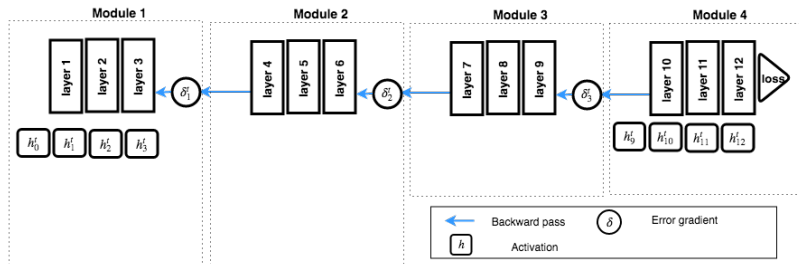


Original formulation:

$$\begin{aligned} \min_w \quad & f(h_L, y) \\ \text{s.t.} \quad & h_l = F_l(h_{l-1}; w_l) \end{aligned}$$

New formulation:

$$\begin{aligned} \min_{w, \delta} \quad & \sum_{k=1}^{K-1} \left\| \delta_k^t - \frac{\partial f_{h_{L_k}^t}(w^t)}{\partial h_{L_k}^t} \right\|_2^2 + f(h_{L_K}^t, y^t) \\ \text{s.t.} \quad & h_{L_k}^t = F_{\mathcal{G}(k)}(h_{L_{k-1}}^t; w_{\mathcal{G}(k)}^t) \end{aligned}$$



Module 1:

$$\min_{w, \delta} \left\| \delta_1^t - \frac{\partial f_{h_{L_1}^t}(w^t)}{\partial h_{L_1}^t} \right\|_2^2$$

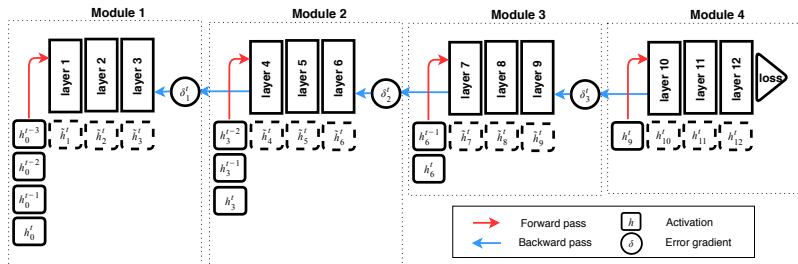
$$s.t. \quad h_{L_1}^t = F_{G(1)}(h_{L_0}^t; w_{G(1)}^t)$$

$$\text{We approximate } \delta_1^t = \frac{\partial f_{h_{L_1}^t}(w^t)}{\partial h_{L_1}^{t-3}}.$$

Module 4:

$$\min_{w, \delta} f(h_{L_4}^t, y^t)$$

$$s.t. \quad h_{L_4}^t = F_{G(4)}(h_{L_3}^t; w_{G(4)}^t)$$



Forward pass:

$$h_{L_k}^t = F_{G(k)}(h_{L_{k-1}}^t; w_{G(k)}^t) \quad (\text{Play})$$

Backward pass:

$$\tilde{h}_{L_k}^t = F_{G(k)}(h_{L_{k-1}}^{t+K-K}; w_{G(k)}^t) \quad (\text{Replay})$$

Apply chain rule using  $\tilde{h}_{L_k}^t$  and  $\delta_k^t$  in each module.

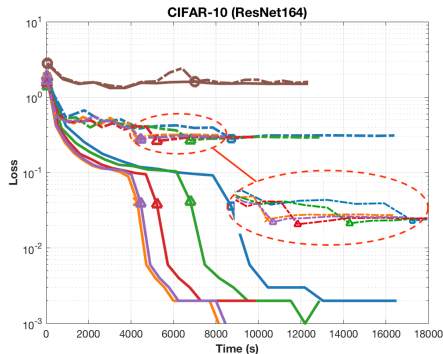
**Assumption 1 (Sufficient direction)** We assume that the expectation of the descent direction  $\mathbb{E} \left[ \sum_{k=1}^K g_{\mathcal{G}(k)}^t \right]$  in Algorithm 1 is a sufficient descent direction of the loss  $f(w^t)$  regarding  $w^t$ . Let  $\nabla f(w^t)$  denote the full gradient of the loss, there exists a constant  $\sigma > 0$  such that,

$$\left\langle \nabla f(w^t), \mathbb{E} \left[ \sum_{k=1}^K g_{\mathcal{G}(k)}^t \right] \right\rangle \geq \sigma \|\nabla f(w^t)\|_2^2. \quad (8)$$

Convergence Guarantee:

$$\frac{1}{\sum_{t=0}^{T-1} \gamma_t} \sum_{t=0}^{T-1} \gamma_t \mathbb{E} \|\nabla f(w^t)\|_2^2 \leq \frac{f(w^0) - f(w^*)}{\sigma \sum_{t=0}^{T-1} \gamma_t} + \frac{LM}{2\sigma} \frac{\sum_{t=0}^{T-1} \gamma_t^2}{\sum_{t=0}^{T-1} \gamma_t}. \quad (1)$$

- Faster Convergence.
- Lower Memory Consumption.
- Better Generalization Error.



# Thanks !

Welcome to poster #12  
Room 210 & 230 AB