

Review of Batch Normalization

Guannan Hu

April 13, 2018

Batch Normalization allows us to **use much higher learning rates and be less careful about initialization**, and it also acts as a **regularizer**, in some cases eliminating the need for Dropout.

Algorithm 1 Batch Normalizing Transform, applied to activation x over a mini-batch

Input: Values of x over a mini-batch: $\mathcal{B} = x_{1\dots m}$; Parameters to be learned: γ, β

Output: $y_i = \text{BN}_{\gamma, \beta}(x_i)$

$$\mu_{\mathcal{B}} \leftarrow \frac{1}{m} \sum_{i=1}^m x_i \quad \triangleright \text{mini-batch mean}$$

$$\sigma_{\mathcal{B}}^2 \leftarrow \frac{1}{m} \sum_{i=1}^m (x_i - \mu_{\mathcal{B}})^2 \quad \triangleright \text{mini-batch variance}$$

$$\hat{x}_i = \frac{x_i - \mu_{\mathcal{B}}}{\sqrt{\sigma_{\mathcal{B}}^2 + \epsilon}} \quad \triangleright \text{normalize}$$

$$y_i \leftarrow \gamma \hat{x}_i + \beta \equiv \text{BN}_{\gamma, \beta}(x_i) \quad \triangleright \text{Scale and shift}$$

During training we need to backpropagate the gradient of loss l through this transformation, as well as compute the gradients with respect to the parameters of the BN transform. We use chain rule, as follows:

$$\begin{aligned} \frac{\partial l}{\partial \hat{x}_i} &= \frac{\partial l}{\partial y_i} \cdot \gamma \\ \frac{\partial l}{\partial \sigma_{\mathcal{B}}^2} &= \sum_{i=1}^m \frac{\partial l}{\partial \hat{x}_i} \cdot (x_i - \mu_{\mathcal{B}}) \cdot \frac{-1}{2} (\sigma_{\mathcal{B}}^2 + \epsilon)^{-\frac{3}{2}} \\ \frac{\partial l}{\partial \mu_{\mathcal{B}}} &= \left(\sum_{i=1}^m \frac{\partial l}{\partial \hat{x}_i} \cdot \frac{-1}{\sqrt{\sigma_{\mathcal{B}}^2 + \epsilon}} \right) + \frac{\partial l}{\partial \sigma_{\mathcal{B}}^2} \cdot \frac{\sum_{i=1}^m -2(x_i - \mu_{\mathcal{B}})}{m} \\ \frac{\partial l}{\partial x_i} &= \frac{\partial l}{\partial \hat{x}_i} \cdot \frac{1}{\sqrt{\sigma_{\mathcal{B}}^2 + \epsilon}} + \frac{\partial l}{\partial \sigma_{\mathcal{B}}^2} \cdot \frac{2(x_i - \mu_{\mathcal{B}})}{m} + \frac{\partial l}{\partial \mu_{\mathcal{B}}} \cdot \frac{1}{m} \\ \frac{\partial l}{\partial \gamma} &= \sum_{i=1}^m \frac{\partial l}{\partial y_i} \cdot \hat{x}_i \\ \frac{\partial l}{\partial \beta} &= \sum_{i=1}^m \frac{\partial l}{\partial y_i} \end{aligned}$$

Algorithm 2 Training a Batch-Normalized Network

Input: Network N with trainable paramters Θ ;

Output: Batch-normalized network for ingerence, $N_{\text{BN}}^{\text{inf}}$

- 1: $N_{\text{BN}}^{\text{tr}} \leftarrow N$ \triangleright Training BN network
- 2: **for** $k = 1 \dots K$ **do**
- 3: Add transformation $y(k) = \text{BN}_{\gamma^{(k)}, \beta^{(k)}}(x^{(k)})$ to $N_{\text{BN}}^{\text{tr}}$ (Alg.1)
- 4: Modify each layer in $N_{\text{BN}}^{\text{tr}}$ with input $x^{(k)}$ to take $y^{(k)}$ instead
- 5: **end for**
- 6: Train $N_{\text{BN}}^{\text{tr}}$ to optimize the parameters $\Theta \cup \{\gamma^{(k)}, \beta^{(k)}\}_{k=1}^K$
- 7: $N_{\text{BN}}^{\text{inf}} \leftarrow N_{\text{BN}}^{\text{tr}}$ \triangleright Inference BN network with frozen paramters
- 8: **for** $k = 1 \dots K$ **do**
- 9: //For clarity, $x \equiv x^{(k)}, \gamma \equiv \gamma^{(k)}, \mu_{\mathcal{B}} \equiv \mu_{\mathcal{B}}^{(k)}$, etc.
- 10: Process multiple training mini-batches \mathcal{B} , each of size m , and average over them:

$$\begin{aligned} \text{E}[x] &\leftarrow \text{E}_{\mathcal{B}}[\mu_{\mathcal{B}}] \\ \text{Var}[x] &\leftarrow \frac{m}{m-1} \text{E}_{\mathcal{B}}[\sigma_{\mathcal{B}}^2] \end{aligned}$$

- 11: In $N_{\text{BN}}^{\text{inf}}$, replace the transform $y = \text{BN}_{\gamma, \beta}(x)$ with $y = \frac{\gamma}{\sqrt{\text{Var}[x] + \epsilon}} \cdot x + (\beta - \frac{\gamma \text{E}(x)}{\sqrt{\text{Var}[x] + \epsilon}})$
 - 12: **end for**
-