# Batch Normalization: Accelerating Deep Network Training by Reducing Internal Covariate Shift

Sergey Ioffe, Christian Szegedy
Google

Reviewed by Zhao Song

December 17, 2015

# Stochastic gradient descent (SGD)

- Minimize the expected loss over the training set:

$$\hat{\theta} = \arg \min_{\theta} \mathrm{E}_{x \sim D}[\ell(x, \theta)]$$

- Parameters update according to the gradient of mini-batches

$$\theta \leftarrow \theta - \frac{\alpha}{m} \sum_{i=1}^{m} \frac{\partial \ell(x_i, \theta)}{\partial \theta}$$

- Careful tuning of learning rates and initial parameters.

## Internal covariate shift

- Covariate shift: Changes of input distribution to a learning system

$$\ell = F(x, \theta)$$

- Internal covariate shift: Extension to the deep network

$$
\begin{aligned}
\ell &= F_2(F_1(u, \theta_1), \theta_2) \\
&= F_2(x, \theta_2)
\end{aligned}
$$

# Reducing internal covariate shift

- Whitening the inputs to each layer:

$$x \leftarrow \frac{x - \mathrm{E}[x]}{\sqrt{\mathrm{Var}[x]}}$$

  However, gradient descent does not take into account the normalization [Ioffe and Szegedy 2015].

- Mean and variance of an activation depend on model parameters

$$\frac{\partial \, \mathrm{E}[x]}{\partial \theta} \text{ and } \frac{\partial \, \mathrm{Var}[x]}{\partial \theta}$$

# Batch normalizing transform

**Input:** Values of $x$ over a mini-batch: $\mathcal{B} = \{x_{1...m}\}$;
Parameters to be learned: $\gamma, \beta$

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

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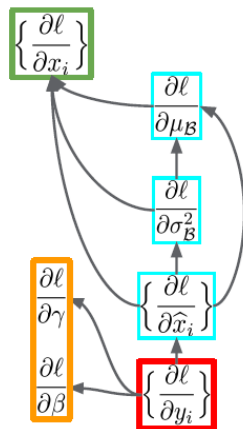
$$\sigma_{\mathcal{B}}^2 \leftarrow \frac{1}{m} \sum_{i=1}^{m} (x_i - \mu_{\mathcal{B}})^2 \qquad \text{// mini-batch variance}$$

$$\widehat{x}_i \leftarrow \frac{x_i - \mu_{\mathcal{B}}}{\sqrt{\sigma_{\mathcal{B}}^2 + \epsilon}} \qquad \text{// normalize}$$

$$y_i \leftarrow \gamma \widehat{x}_i + \beta \equiv \text{BN}_{\gamma,\beta}(x_i) \qquad \text{// scale and shift}$$

**Algorithm 1:** Batch Normalizing Transform, applied to activation $x$ over a mini-batch.

# Backpropagation with batch normalization



$$\frac{\partial \ell}{\partial \widehat{x}_i} = \frac{\partial \ell}{\partial y_i} \cdot \gamma$$

$$\frac{\partial \ell}{\partial \sigma_{\mathcal{B}}^2} = \sum_{i=1}^{m} \frac{\partial \ell}{\partial \widehat{x}_i} \cdot (x_i - \mu_{\mathcal{B}}) \cdot \frac{-1}{2}(\sigma_{\mathcal{B}}^2 + \epsilon)^{-3/2}$$

$$\frac{\partial \ell}{\partial \mu_{\mathcal{B}}} = \sum_{i=1}^{m} \frac{\partial \ell}{\partial \widehat{x}_i} \cdot \frac{-1}{\sqrt{\sigma_{\mathcal{B}}^2 + \epsilon}}$$

$$\frac{\partial \ell}{\partial x_i} = \frac{\partial \ell}{\partial \widehat{x}_i} \cdot \frac{1}{\sqrt{\sigma_{\mathcal{B}}^2 + \epsilon}} + \frac{\partial \ell}{\partial \sigma_{\mathcal{B}}^2} \cdot \frac{2(x_i - \mu_{\mathcal{B}})}{m} + \frac{\partial \ell}{\partial \mu_{\mathcal{B}}} \cdot \frac{1}{m}$$

$$\frac{\partial \ell}{\partial \gamma} = \sum_{i=1}^{m} \frac{\partial \ell}{\partial y_i} \cdot \widehat{x}_i$$

$$\frac{\partial \ell}{\partial \beta} = \sum_{i=1}^{m} \frac{\partial \ell}{\partial y_i}$$

(a)                                                                 (b)

# Training a batch-normalized network

**Input:** Network $N$ with trainable parameters $\Theta$;
  subset of activations $\{x^{(k)}\}_{k=1}^{K}$
**Output:** Batch-normalized network for inference, $N_{\text{BN}}^{\text{inf}}$
1: $N_{\text{BN}}^{\text{tr}} \leftarrow N$  // Training BN network
2: **for** $k = 1 \ldots K$ **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: $N_{\text{BN}}^{\text{inf}} \leftarrow N_{\text{BN}}^{\text{tr}}$  // Inference BN network with frozen
    // parameters
8: **for** $k = 1 \ldots K$ **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:
$$\text{E}[x] \leftarrow \text{E}_{\mathcal{B}}[\mu_{\mathcal{B}}]$$
$$\text{Var}[x] \leftarrow \frac{m}{m-1}\text{E}_{\mathcal{B}}[\sigma_{\mathcal{B}}^2]$$
11:  In $N_{\text{BN}}^{\text{inf}}$, replace the transform $y = \text{BN}_{\gamma, \beta}(x)$ with
$$y = \frac{\gamma}{\sqrt{\text{Var}[x]+\epsilon}} \cdot x + \left(\beta - \frac{\gamma \text{E}[x]}{\sqrt{\text{Var}[x]+\epsilon}}\right)$$
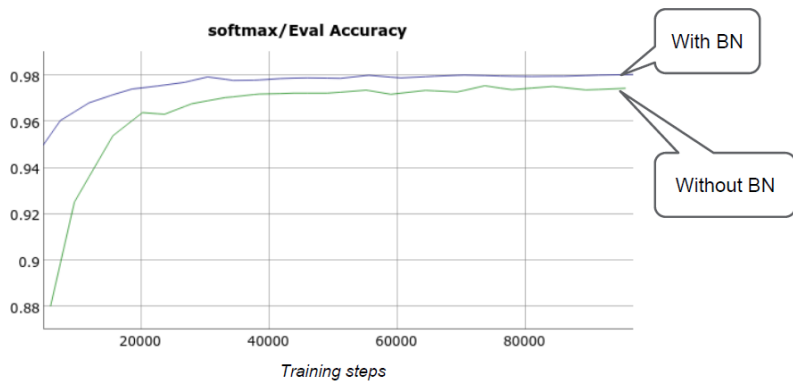12: **end for**

**Algorithm 2:** Training a Batch-Normalized Network
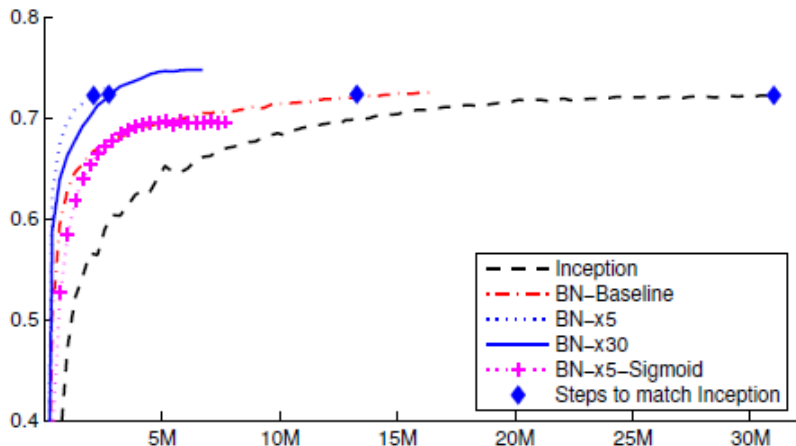
# Experiments

- MNIST
  - 3 fully-connect hidden layers with 100 nodes in each layer.
  - Sigmoid activation function.
  - Mini-batch size to be 60.
- ImageNet
  - The Inception network [Szegedy et al. 2014].
  - SGD with momentum [Sutskever et al. 2013].
  - Mini-batch size to be 32.

# Learning curve on MNIST [Ioffe 2015]



**softmax/Eval Accuracy**

With BN

Without BN

*Training steps*

# Learning curve on ImageNet with single networks

# Classification results on ImageNet

| Model | Resolution | Crops | Models | Top-1 error | Top-5 error |
|---|---|---|---|---|---|
| GoogLeNet ensemble | 224 | 144 | 7 | - | 6.67% |
| Deep Image low-res | 256 | - | 1 | - | 7.96% |
| Deep Image high-res | 512 | - | 1 | 24.88 | 7.42% |
| Deep Image ensemble | up to 512 | - | - | - | 5.98% |
| MSRA multicrop | up to 480 | - | - | - | 5.71% |
| MSRA ensemble | up to 480 | - | - | - | 4.94%* |
| BN-Inception single crop | 224 | 1 | 1 | 25.2% | 7.82% |
| BN-Inception multicrop | 224 | 144 | 1 | 21.99% | 5.82% |
| BN-Inception ensemble | 224 | 144 | 6 | 20.1% | **4.82%*** |

# References I

📄 Sergey Ioffe. *Batch Normalization Presentation*. 2015.

📄 Sergey Ioffe and Christian Szegedy. "Batch normalization: Accelerating deep network training by reducing internal covariate shift." : *ICML*. 2015.

📄 Ilya Sutskever, James Martens, George Dahl, and Geoffrey Hinton. "On the importance of initialization and momentum in deep learning." : *ICML*. 2013.

📄 Christian Szegedy, Wei Liu, Yangqing Jia, Pierre Sermanet, Scott Reed, Dragomir Anguelov, Dumitru Erhan, Vincent Vanhoucke, and Andrew Rabinovich. "Going deeper with convolutions." *arXiv preprint arXiv:1409.4842* (2014).