

# Normalization techniques in deep learning

Author: Nikita Bashaev, 171

February 28, 2020

# Plan

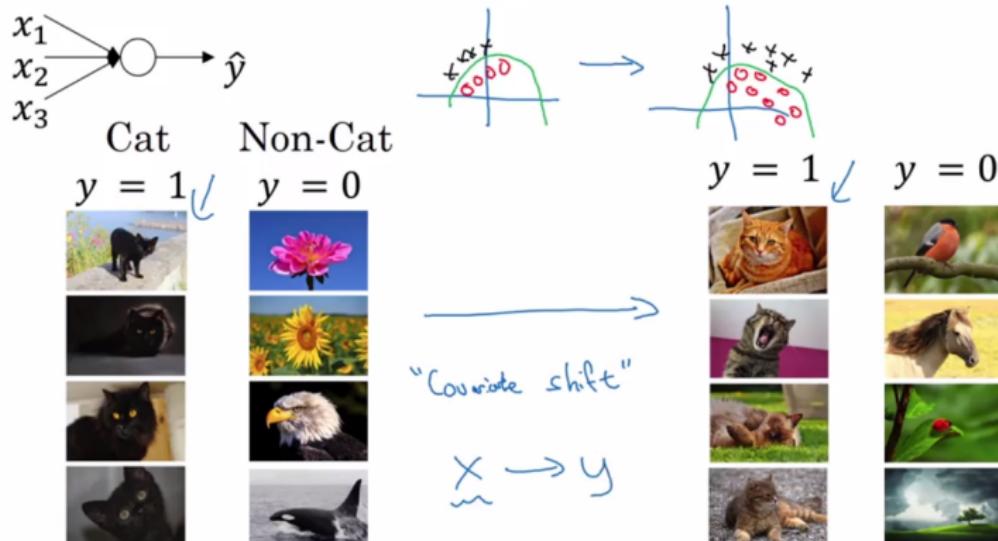
## ① Batch Normalization

- Working principle
- Possible reasons for its phenomenal success
- The demand for other normalization techniques

## ② Layer/Instance/Group Normalization

## ③ GAN's and Spectral Normalization

# Covariate Shift



## BN' Working Principle

**Input:** Values of  $x$  over a mini-batch:  $\mathcal{B} = \{x_1 \dots m\}$ ;

Parameters to be learned:  $\gamma, \beta$

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

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

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

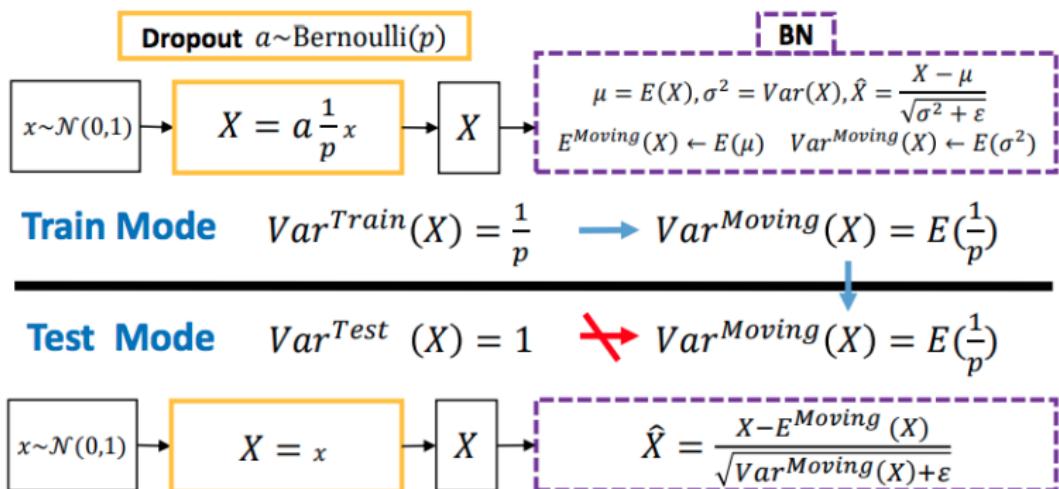
$$\hat{x}_i \leftarrow \frac{x_i - \mu_{\mathcal{B}}}{\sqrt{\sigma_{\mathcal{B}}^2 + \epsilon}} \quad // \text{normalize}$$

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

## BN in CNN

For convolutional layers, we additionally want the normalization to obey the convolutional property – so that different elements of the same feature map, at different locations, are normalized in the same way. To achieve this, we jointly normalize all the activations in a mini-batch, over all locations.

## BN at test time



## BN' effectiveness – covariate shift

The original explanation of "internal covariate shift" is not very precise, so at best it is a heuristic explanation, and at worst it might be completely wrong.

# BN' effectiveness – covariate shift

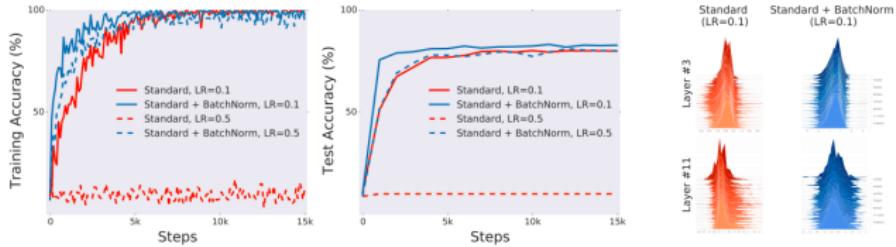


Figure 1: Comparison of (a) training (optimization) and (b) test (generalization) performance of a standard VGG network trained on CIFAR-10 with and without BatchNorm (details in Appendix A). There is a consistent gain in training speed in models with BatchNorm layers. (c) Even though the gap between the performance of the BatchNorm and non-BatchNorm networks is clear, the difference in the evolution of layer input distributions seems to be much less pronounced. (Here, we sampled activations of a given layer and visualized their distribution over training steps.)

# BN' effectiveness – smoothness of loss landscape

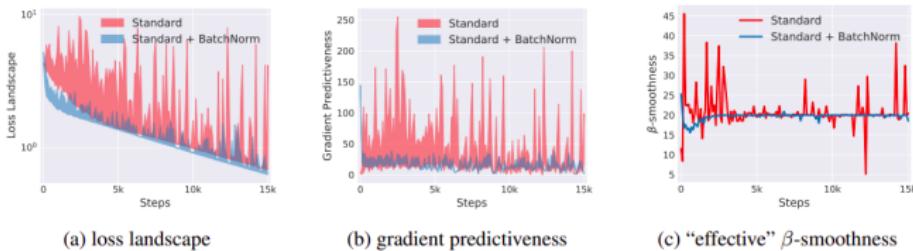


Figure 4: Analysis of the optimization landscape of VGG networks. At a particular training step, we measure the variation (shaded region) in loss (a) and  $\ell_2$  changes in the gradient (b) as we move in the gradient direction. The “effective”  $\beta$ -smoothness (c) refers to the maximum difference (in  $\ell_2$ -norm) in gradient over distance moved in that direction. There is a clear improvement in all of these measures in networks with BatchNorm, indicating a more well-behaved loss landscape. (Here, we cap the maximum distance to be  $\eta = 0.4 \times$  the gradient since for larger steps the standard network just performs worse (see Figure 1). BatchNorm however continues to provide smoothing for even larger distances.) Note that these results are supported by our theoretical findings (Section 4).

## BN' drawbacks

- When batch size is small BN behaves rather randomly
- Statistics become invalid if batch size changes  
(e.g., fine-tuning)
- Batch norm is hard to perform in distributed training
- It is not the best option for RNN's

# Layer/Instance/Group Normalization

$$\mu_i = \frac{1}{m} \sum_{k \in \mathcal{S}_i} x_k, \quad \sigma_i = \sqrt{\frac{1}{m} \sum_{k \in \mathcal{S}_i} (x_k - \mu_i)^2 + \epsilon},$$

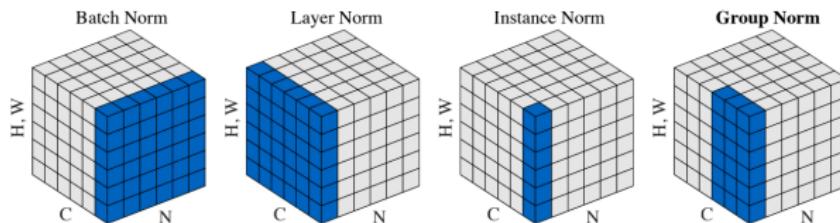
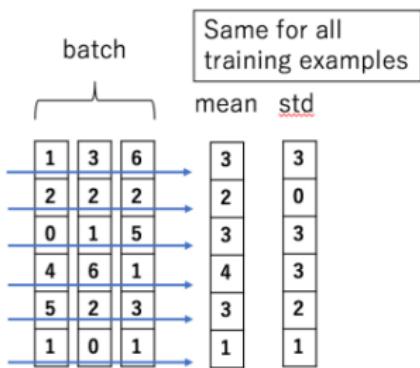


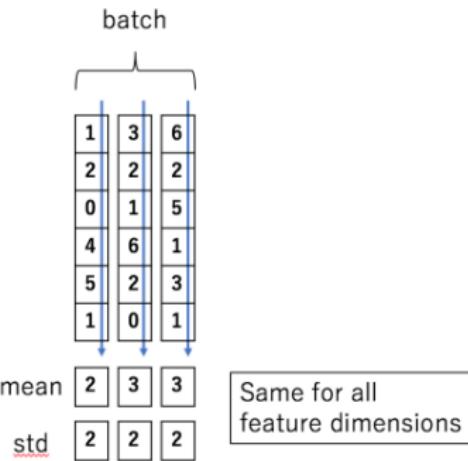
Figure 2. **Normalization methods.** Each subplot shows a feature map tensor, with  $N$  as the batch axis,  $C$  as the channel axis, and  $(H, W)$  as the spatial axes. The pixels in blue are normalized by the same mean and variance, computed by aggregating the values of these pixels.

# Layer/Instance/Group Normalization

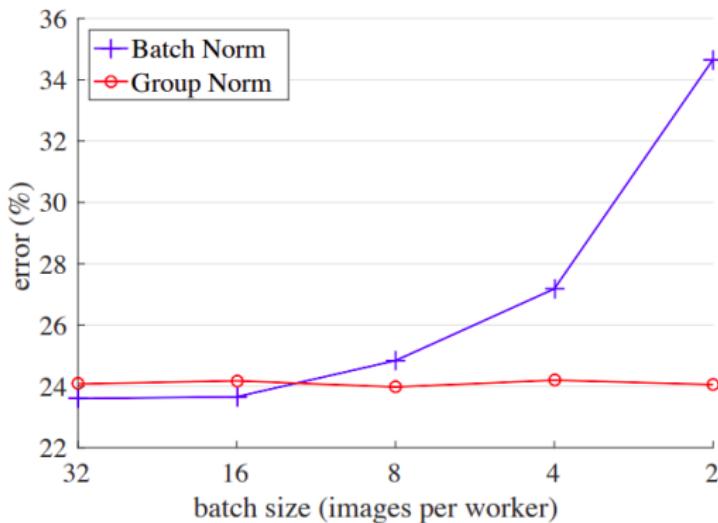
Batch Normalization



Layer Normalization



# Layer/Instance/Group Normalization



**Figure 1. ImageNet classification error vs. batch sizes.** This is a ResNet-50 model trained in the ImageNet training set using 8 workers (GPUs), evaluated in the validation set.

## Layer/Instance/Group Normalization

Applicability:

- Layer Norm – RNN's and Transformers
- Instance Norm – Style Transfer
- Group Norm – CV tasks with small batch size

# Spectral Normalization

Recall that  $f$  is  $L$ -Lipschitz if  $\|f(x_1) - f(x_2)\| \leq L\|x_1 - x_2\|$  for all  $x_1$  and  $x_2$

$$f(\mathbf{x}, \theta) = W^{L+1} a_L(W^L(a_{L-1}(W^{L-1}(\dots a_1(W^1 \mathbf{x}) \dots)))),$$

$$\begin{aligned} \|f\|_{\text{Lip}} &\leq \|(\mathbf{h}_L \mapsto W^{L+1} \mathbf{h}_L)\|_{\text{Lip}} \cdot \|a_L\|_{\text{Lip}} \cdot \|(\mathbf{h}_{L-1} \mapsto W^L \mathbf{h}_{L-1})\|_{\text{Lip}} \\ &\quad \cdots \|a_1\|_{\text{Lip}} \cdot \|(\mathbf{h}_0 \mapsto W^1 \mathbf{h}_0)\|_{\text{Lip}} = \prod_{l=1}^{L+1} \|(\mathbf{h}_{l-1} \mapsto W^l \mathbf{h}_{l-1})\|_{\text{Lip}} = \prod_{l=1}^{L+1} \sigma(W^l). \end{aligned} \quad (7)$$

Our *spectral normalization* normalizes the spectral norm of the weight matrix  $W$  so that it satisfies the Lipschitz constraint  $\sigma(W) = 1$ :

$$\bar{W}_{\text{SN}}(W) := W/\sigma(W). \quad (8)$$

# Spectral Normalization

---

**Algorithm 1** SGD with spectral normalization

---

- Initialize  $\tilde{\mathbf{u}}_l \in \mathcal{R}^{d_l}$  for  $l = 1, \dots, L$  with a random vector (sampled from isotropic distribution).
- For each update and each layer  $l$ :
  1. Apply power iteration method to a unnormalized weight  $W^l$ :

$$\tilde{\mathbf{v}}_l \leftarrow (W^l)^T \tilde{\mathbf{u}}_l / \| (W^l)^T \tilde{\mathbf{u}}_l \|_2 \quad (20)$$

$$\tilde{\mathbf{u}}_l \leftarrow W^l \tilde{\mathbf{v}}_l / \| W^l \tilde{\mathbf{v}}_l \|_2 \quad (21)$$

2. Calculate  $\bar{W}_{\text{SN}}$  with the spectral norm:

$$\bar{W}_{\text{SN}}^l(W^l) = W^l / \sigma(W^l), \text{ where } \sigma(W^l) = \tilde{\mathbf{u}}_l^T W^l \tilde{\mathbf{v}}_l \quad (22)$$

3. Update  $W^l$  with SGD on mini-batch dataset  $\mathcal{D}_M$  with a learning rate  $\alpha$ :

$$W^l \leftarrow W^l - \alpha \nabla_{W^l} \ell(\bar{W}_{\text{SN}}^l(W^l), \mathcal{D}_M) \quad (23)$$

---

# Spectral Normalization

## C.4 GENERATED IMAGES ON CIFAR10 WITH GAN-GP, LAYER NORMALIZATION AND BATCH NORMALIZATION

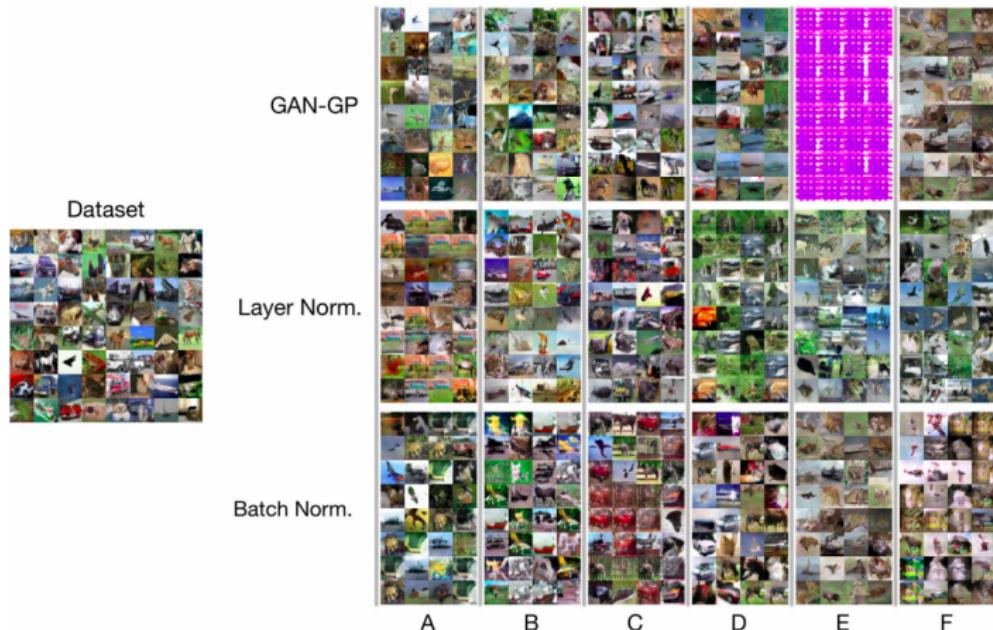


Figure 12: Generated images with GAN-GP, Layer Norm and Batch Norm on CIFAR-10

# Spectral Normalization



## Questions

- ① How does BN Layer work during inference?
- ② Describe the working principle of Group Normalization.
- ③ Why is Spectral Normalization so successful in GAN's training?

## Bibliography

- Batch Normalization: Accelerating Deep Network Training by Reducing Internal Covariate Shift. [arXiv:1502.03167](#)
- How Does Batch Normalization Help Optimization?  
[arXiv:1805.11604](#)
- Group Normalization. [arXiv:1803.08494](#)
- Spectral Normalization for Generative Adversarial Networks.  
[arXiv:1802.05957](#)