

3-3 Batch Normalization

Zhonglei Wang

WISE and SOE, XMU, 2025

Contents

1. Introduction

2. Forward propagation

3. Backpropagation

Introduction

1. Proposed by Ioffe and Szegedy (2015)

- Tries to solve an interval covariate shift problem
- Normalize values **after** linear transformation but **before** activation for each neuron
- Introduce two more parameters to allow for heterogeneity

Forward propagation

1. For a mini-batch with m training examples, the forward propagation for the l th layer is

$$\mathbf{Z}^{[l]} = \mathbf{A}^{[l-1]}(\mathbf{W}^{[l]})^T + (\mathbf{b}^{[l]})^T \in \mathbb{R}^{m \times d^{[l]}}; \quad \mathbf{A}^{[l]} = \sigma^{[l]}(\mathbf{Z}^{[l]}) \in \mathbb{R}^{m \times d^{[l]}}$$

2. Denote $\mathbf{Z}^{[l]} = (\mathbf{z}_1^{[l]}, \dots, \mathbf{z}_m^{[l]})^T$

3. After linear transformation, we consider the following **new** calculations:

$$\boldsymbol{\mu}^{[l]} = m^{-1}(\mathbf{Z}^{[l]})^T \mathbf{1} \in \mathbb{R}^{d^{[l]} \times 1}$$

$$\boldsymbol{\sigma}^{2[l]} = m^{-1} \sum_{i=1}^m \left(\mathbf{z}_i^{[l]} - \boldsymbol{\mu}^{[l]} \right) \circ \left(\mathbf{z}_i^{[l]} - \boldsymbol{\mu}^{[l]} \right) \in \mathbb{R}^{d^{[l]} \times 1}$$

$$\tilde{\mathbf{z}}_{i,norm}^{[l]} = \frac{\mathbf{z}_i^{[l]} - \boldsymbol{\mu}^{[l]}}{\sqrt{\boldsymbol{\sigma}^{2[l]} + \epsilon}} \in \mathbb{R}^{d^{[l]} \times 1}$$

$$\tilde{\mathbf{z}}_i^{[l]} = \gamma^{[l]} \circ \tilde{\mathbf{z}}_{i,norm}^{[l]} + \boldsymbol{\beta}^{[l]} \in \mathbb{R}^{d^{[l]} \times 1}$$

$$\tilde{\mathbf{Z}}^{[l]} = (\tilde{\mathbf{z}}_1^{[l]}, \dots, \tilde{\mathbf{z}}_m^{[l]})^T \in \mathbb{R}^{m \times d^{[l]}}$$

Forward propagation

1. Vectorization for the **red** calculations, but leave the **blue** parts alone

$$\mathbf{Z}^{[l]} = \mathbf{A}^{[l-1]}(\mathbf{W}^{[l]})^T \in \mathbb{R}^{m \times d^{[l]}}$$

$$\boldsymbol{\mu}^{[l]} = m^{-1}(\mathbf{Z}^{[l]})^T \mathbf{1} \in \mathbb{R}^{d^{[l]} \times 1}$$

$$\check{\mathbf{Z}}^{[l]} = \mathbf{Z}^{[l]} - \mathbf{1}(\boldsymbol{\mu}^{[l]})^T \in \mathbb{R}^{m \times d^{[l]}}$$

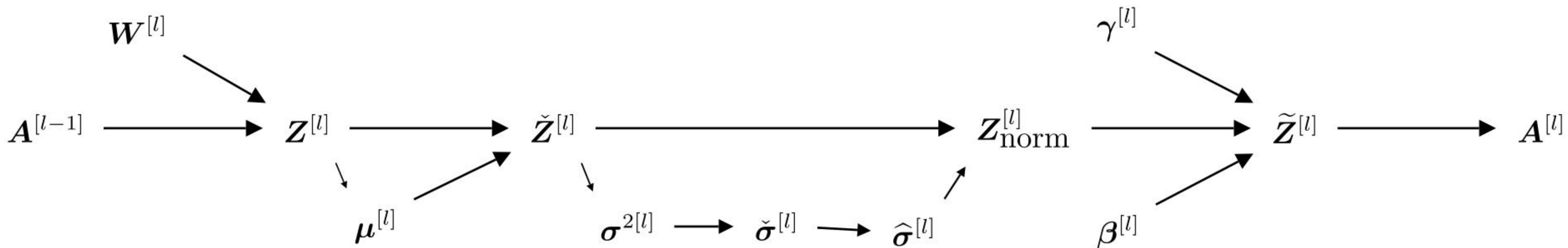
$$\boldsymbol{\sigma}^{2[l]} = m^{-1} \sum_{i=1}^m \check{\mathbf{z}}_i^{[l]} \circ \check{\mathbf{z}}_i^{[l]} \in \mathbb{R}^{d^{[l]} \times 1}$$

$$\check{\boldsymbol{\sigma}}^{[l]} = \sqrt{\boldsymbol{\sigma}^{2[l]} + \epsilon} \in \mathbb{R}^{d^{[l]} \times 1}$$

$$\hat{\boldsymbol{\sigma}}^{[l]} = (\check{\boldsymbol{\sigma}}^{[l]})^{-1} \in \mathbb{R}^{d^{[l]} \times 1}$$

$$\mathbf{Z}_{\text{norm}}^{[l]} = \check{\mathbf{Z}}^{[l]} \circ \{\mathbf{1}(\hat{\boldsymbol{\sigma}}^{[l]})^T\} \in \mathbb{R}^{m \times d^{[l]}}$$

2. Notice that the previous bias term $\mathbf{b}^{[l]}$ is useless for batch normalization



Forward propagation for the **red** parts:

$$\mathbf{Z}^{[l]} = \mathbf{A}^{[l-1]}(\mathbf{W}^{[l]})^T \in \mathbb{R}^{m \times d^{[l]}}$$

$$\check{\sigma}^{[l]} = \sqrt{\sigma^2[l] + \epsilon} \in \mathbb{R}^{d^{[l]} \times 1}$$

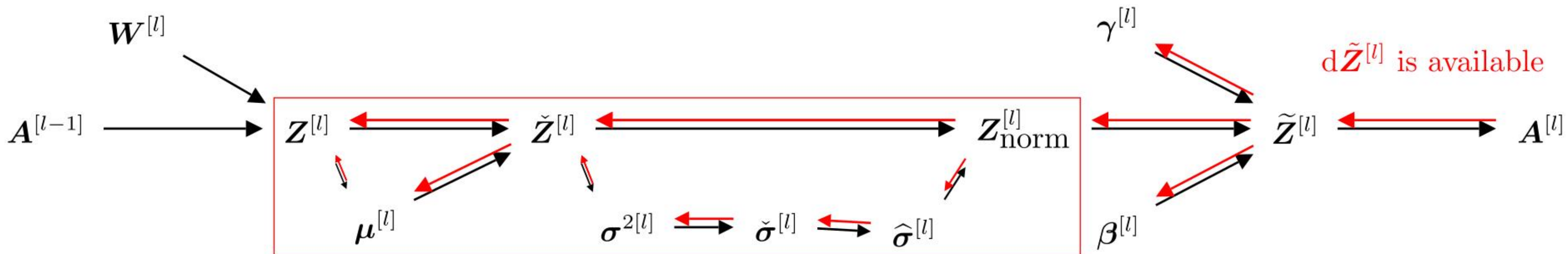
$$\mu^{[l]} = (\mathbf{Z}^{[l]})^T \mathbf{1} \in \mathbb{R}^{d^{[l]} \times 1}$$

$$\hat{\sigma}^{[l]} = (\check{\sigma}^{[l]})^{-1} \in \mathbb{R}^{d^{[l]} \times 1}$$

$$\check{\mathbf{Z}}^{[l]} = \mathbf{Z}^{[l]} - \mathbf{1}(\mu^{[l]})^T \in \mathbb{R}^{m \times d^{[l]}}$$

$$\mathbf{Z}_{\text{norm}}^{[l]} = \check{\mathbf{Z}}^{[l]} \circ \{\mathbf{1}(\hat{\sigma}^{[l]})^T\} \in \mathbb{R}^{m \times d^{[l]}}$$

$$\sigma^2[l] = m^{-1} \sum_{i=1}^m \check{\mathbf{z}}_i^{[l]} \circ \check{\mathbf{z}}_i^{[l]} \in \mathbb{R}^{d^{[l]} \times 1}$$



Backpropagation:

$$d\mathbf{Z}^{[l]} = d\mathbf{Z}_1^{[l]} + d\mathbf{Z}_2^{[l]} \quad d\boldsymbol{\beta}^{[l]} = \left(d\tilde{\mathbf{Z}}^{[l]}\right)^T \mathbf{1} \quad d\boldsymbol{\gamma}^{[l]} = \left(d\tilde{\mathbf{Z}}^{[l]} \circ \mathbf{Z}_{\text{norm}}^{[l]}\right)^T \mathbf{1}$$

$$d\check{\mathbf{Z}}^{[l]} = d\check{\mathbf{Z}}_1^{[l]} + d\check{\mathbf{Z}}_2^{[l]} \quad d\check{\boldsymbol{\sigma}}^{[l]} = -d\hat{\boldsymbol{\sigma}}^{[l]} \circ \hat{\boldsymbol{\sigma}}^{[l]} \circ \hat{\boldsymbol{\sigma}}^{[l]} \quad d\mathbf{Z}_{\text{norm}}^{[l]} = d\tilde{\mathbf{Z}}^{[l]} \circ \left\{ \mathbf{1} (\boldsymbol{\gamma}^{[l]})^T \right\}$$

$$d\mathbf{Z}_1^{[l]} = d\check{\mathbf{Z}}^{[l]} \quad d\boldsymbol{\sigma}^{2[l]} = d\check{\boldsymbol{\sigma}}^{[l]} \circ \hat{\boldsymbol{\sigma}}^{[l]} / 2 \quad d\check{\mathbf{Z}}_1^{[l]} = d\mathbf{Z}_{\text{norm}}^{[l]} \circ \left\{ \mathbf{1} (\hat{\boldsymbol{\sigma}}^{[l]})^T \right\}$$

$$d\boldsymbol{\mu}^{[l]} = -\left(d\check{\mathbf{Z}}^{[l]}\right)^T \mathbf{1} \quad d\check{\mathbf{Z}}_2^{[l]} = 2m^{-1} \check{\mathbf{Z}}^{[l]} \circ \left\{ \mathbf{1} (d\boldsymbol{\sigma}^{2[l]})^T \right\} \quad d\hat{\boldsymbol{\sigma}}^{[l]} = \left(d\mathbf{Z}_{\text{norm}}^{[l]} \circ \check{\mathbf{Z}}^{[l]}\right)^T \mathbf{1}$$

$$d\mathbf{Z}_2^{[l]} = m^{-1} \mathbf{1} (d\boldsymbol{\mu}^{[l]})^T$$

Remarks

1. Disadvantage

- Since we normalize “inputs” before activation, many different “inputs” may result in same “activations”
- Besides, batch normalization also introduces more model parameters

Remarks

1. Advantage

- Stabilize forward propagation
- Stabilize forward propagation
 - ▷ The variance can be controlled by γ 's
- Higher learning rates
 - ▷ Batch normalization can make loss and its gradients more smooth
- Regularization
 - ▷ We injects noises from other training examples through mean and variance
 - ▷ Thus, batch normalization may improve the generality of the network