

### 3-3 批量归一化

王中雷

厦门大学王亚南经济研究院和经济学院, 2025

# 内容摘要

1. 简介

2. 前向传播

3. 后向传播

# 简介

1. 被 Ioffe and Szegedy (2015) 提出

- 可用于解决协变量偏差的问题
- 在线性变换后但在激活运算前进行归一化
- 基于两个额外模型参数，调整特征之间的异质性

# 前向传播

1. 对于小批量梯度下降法（批量大小为  $m$ ），第  $l$  层的计算为：

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

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

3. 线性变换后，我们考虑下面的新运算：

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

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

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

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

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

# 前向传播

1. 考虑红色计算的向量化，暂时不考虑蓝色部分的计算细节

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

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

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

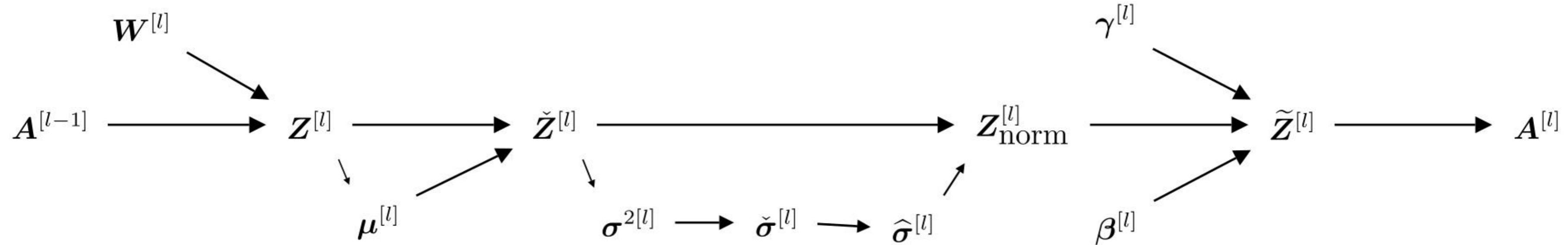
$$\boldsymbol{\sigma}^{2[l]} = m^{-1} \sum_{i=1}^m \check{\mathbf{z}}_i^{[l]} \circ \check{\mathbf{z}}_i^{[l]} \in \mathbb{R}^{d^{[l]} \times 1}$$

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

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

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

2. 注意：偏置项  $\mathbf{b}^{[l]}$  对于批量归一化是无用的



红色计算部分的前向传播过程为：

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

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

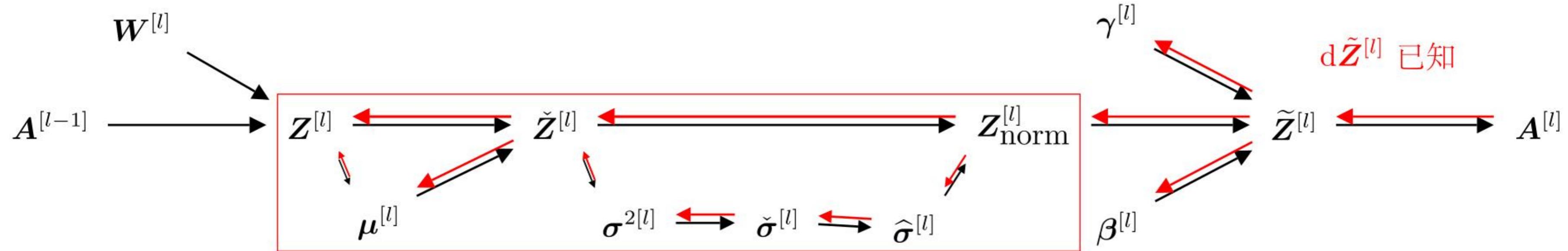
$$\boldsymbol{\mu}^{[l]} = (\mathbf{Z}^{[l]})^T \mathbf{1} \in \mathbb{R}^{d^{[l]} \times 1}$$

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

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

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

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



后向传播：

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

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

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

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

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

# 注意

## 1. 缺点

- 由于我们在激活前使用批量归一化，不同的“输入特征”可能有相同的“激活值”
- 此外，批量归一化对应着更多的模型参数
- 为什么不在激活运算之后进行批量归一化？

# 注意

## 1. 优点：

- 稳定化前向传播
  - ▷ 批量归一化的方差被 $\gamma$  控制
- 可以接受更大的学习率
  - ▷ 批量归一化可使代价函数及其梯度可以更平滑
- 正则化
  - ▷ 批量归一化可被认为将其他训练样本的随机性加入到每个训练样本
  - ▷ 因此，批量归一化可以提升神经网络模型的泛化性