

Xavier and Kaiming initialization

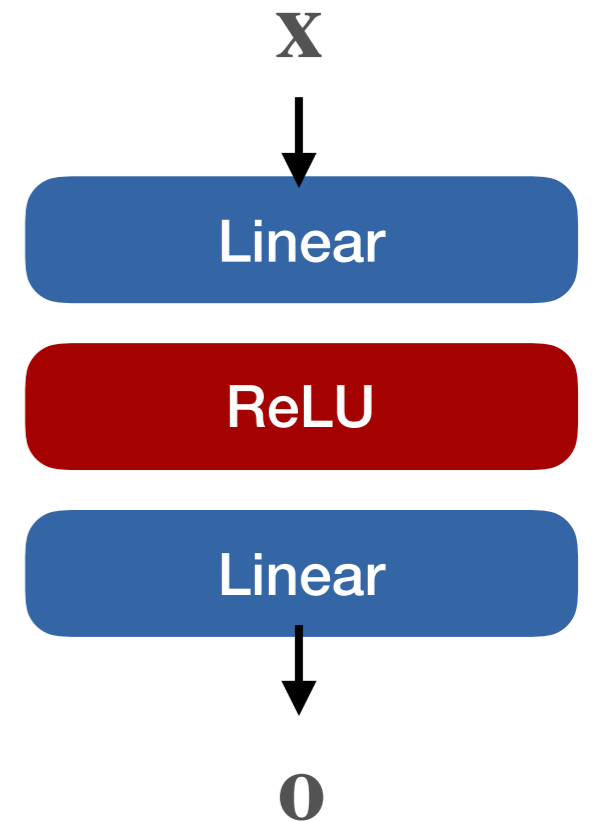
© 2019 Philipp Krähenbühl and Chao-Yuan Wu

Xavier and Kaiming initialization

- Strategy to set variance σ^2 of Normal initialization
- All activations are of similar scale

$$\mathbf{W}_1 \sim \mathcal{N}(\mu_1, \sigma_1^2 \mathbf{I})$$

$$\mathbf{W}_3 \sim \mathcal{N}(\mu_3, \sigma_3^2 \mathbf{I})$$



Random matrix multiplication

$$\mathbf{a}^\top \mathbf{x} \sim \mathcal{N}\left(\mu_a \sum_i \mathbf{x}_i, \|\mathbf{x}\|^2 \sigma_a^2\right) \quad \text{for } \mathbf{a} \sim \mathcal{N}(\mu_a, \sigma_a^2 \mathbf{I})$$

Random matrix multiplication

$$\mathbf{z}_i = \mathbf{W}_{i-1} \mathbf{z}_{i-1} \sim \mathcal{N}(0, \|\mathbf{z}_{i-1}\|^2 \sigma_{W_{i-1}}^2 \mathbf{I}) \quad \text{for} \quad \mathbf{W}_{i-1} \sim \mathcal{N}(0, \sigma_{W_{i-1}}^2 \mathbf{I})$$

Random ReLU

$$\mathbf{z}_{i+1} = \max(\mathbf{z}_i, 0) \quad \text{for} \quad \mathbf{z}_i \sim \mathcal{N}(0, \sigma_i^2 \mathbf{I})$$

$$\mathbb{E}[\|\mathbf{z}_{i+1}\|^2] = \frac{1}{2} n_{\mathbf{z}_i} \sigma_i^2$$

Putting things together

$$\mathbf{z}_i \sim \mathcal{N}(0, \sigma_i^2 \mathbf{I})$$

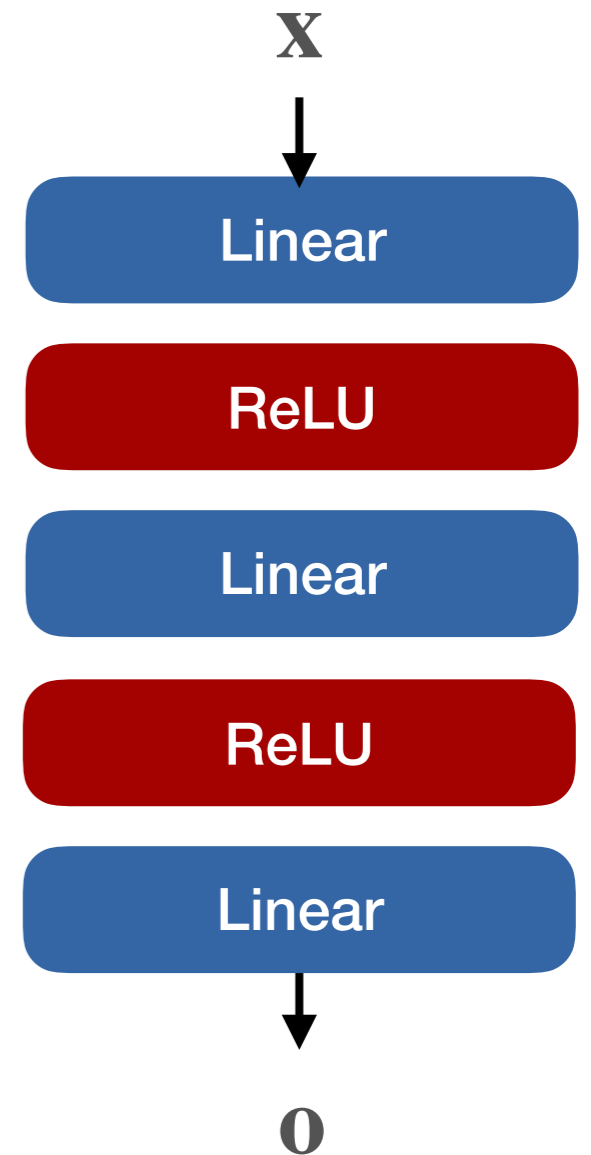
$$\|\mathbf{z}_{i+1}\|^2 \approx \mathbb{E}[\|\mathbf{z}_{i+1}\|^2] = \frac{1}{2} n_{\mathbf{z}_i} \sigma_i^2 \quad \mathbf{z}_{i+2} \sim \mathcal{N}(0, \underbrace{\|\mathbf{z}_{i+1}\|^2 \sigma_{W_{i+1}}^2}_{\sigma_{i+2}^2} \mathbf{I})$$

$$\sigma_{i+2} = \frac{1}{\sqrt{2}} \sigma_{W_{i+1}} \sigma_i \sqrt{n_{\mathbf{z}_i}}$$

$$\sigma_i = \prod_{k=0}^{(i-1)/2} \left(\frac{1}{\sqrt{2}} \sigma_{W_{2k+1}} \sqrt{n_{\mathbf{z}_{2k}}} \right) \sigma_x$$

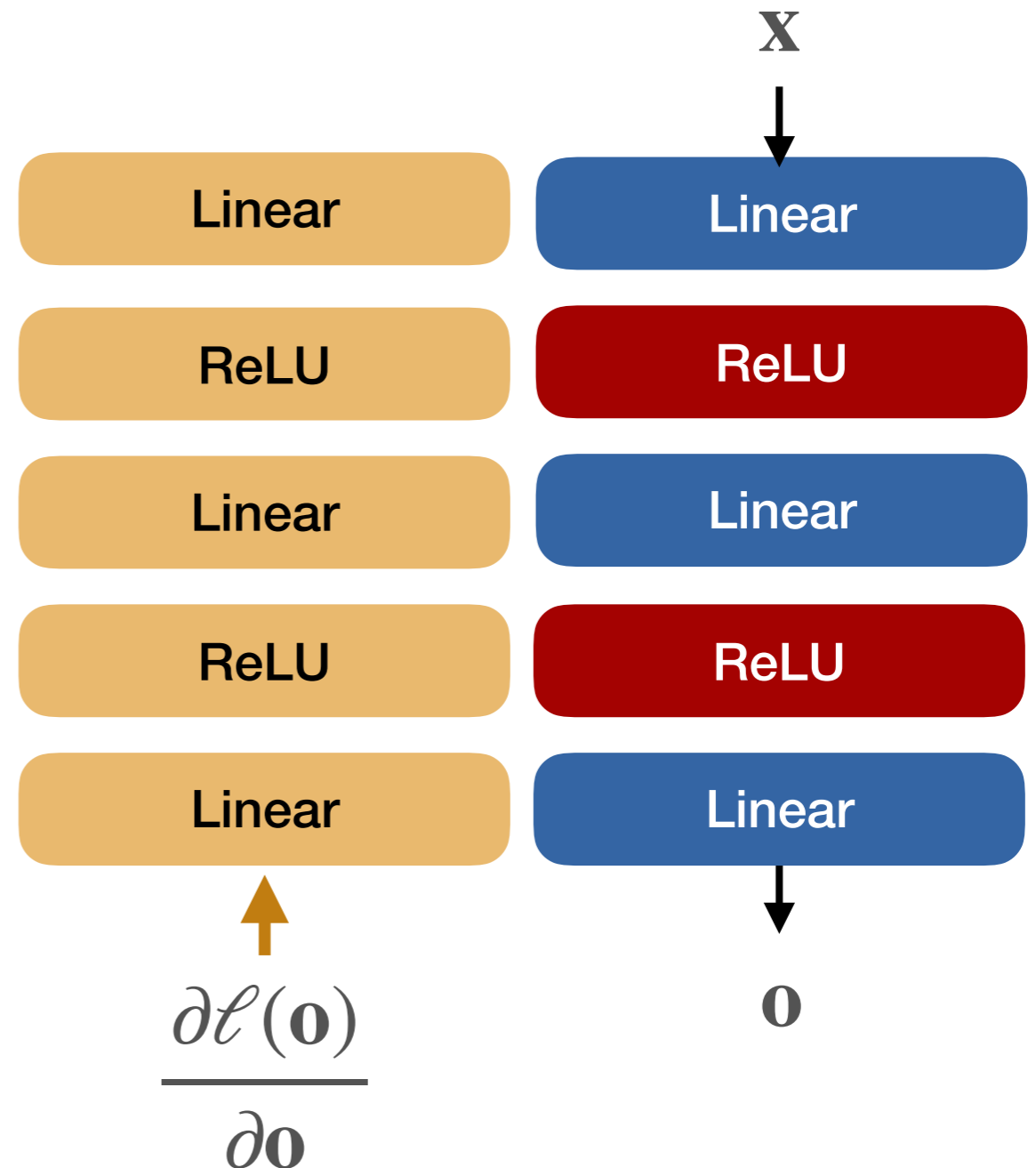
Randomly initialized network

$$\sigma_i = \prod_{k=0}^{(i-1)/2} \left(\frac{1}{\sqrt{2}} \sigma_{W_{2k+1}} \sqrt{n_{z_{2k}}} \right) \sigma_x$$



Variance of back-propagation graph

$$\hat{\sigma}_i = \prod_{k=i/2}^{(N-1)/2} \left(\frac{1}{\sqrt{2}} \sigma_{W_{2k+1}} \sqrt{n_{z_{2k+2}}} \right) \hat{\sigma}_N$$



Xavier initialization

$$\sigma_i = \prod_{k=0}^{(i-1)/2} \left(\frac{1}{\sqrt{2}} \sigma_{W_{2k+1}} \sqrt{n_{z_{2k}}} \right) \sigma_x$$

- Try to keep both activations and gradient magnitude constant

$$\hat{\sigma}_i = \prod_{k=i/2}^{(N-1)/2} \left(\frac{1}{\sqrt{2}} \sigma_{W_{2k+1}} \sqrt{n_{z_{2k+2}}} \right) \hat{\sigma}_N$$

- $$\sigma_W = \sqrt{2} \sqrt{\frac{2}{n_{z_i} + n_{z_{i+1}}}}$$

Kaiming initialization

- Try to keep either activation or gradient magnitude constant

$$\sigma_i = \prod_{k=0}^{(i-1)/2} \left(\frac{1}{\sqrt{2}} \sigma_{W_{2k+1}} \sqrt{n_{z_{2k}}} \right) \sigma_x$$

$$\hat{\sigma}_i = \prod_{k=i/2}^{(N-1)/2} \left(\frac{1}{\sqrt{2}} \sigma_{W_{2k+1}} \sqrt{n_{z_{2k+2}}} \right) \hat{\sigma}_N$$

- $\sigma_W = \sqrt{2} \frac{1}{\sqrt{n_{z_i}}}$

- $\sigma_W = \sqrt{2} \frac{1}{\sqrt{n_{z_{i+1}}}}$

Initialization in practice

- Xavier (default) is often good enough
- Initialize last layer to zero