

31.3.2020

$$\delta^{(l-1)} = \text{diag}(\varphi'(z^{(l-1)})) \cdot W^{(l)} \cdot \delta^{(l)}$$

$$\nabla_{W^{(l)}} l(\cdot) = \delta^{(l)} \cdot v^{(l-1)} \leftarrow \text{activations from layer } l-1$$

$$v^{(l)} = \varphi(W^{(l)} v^{(l-1)})$$

Failure mode:

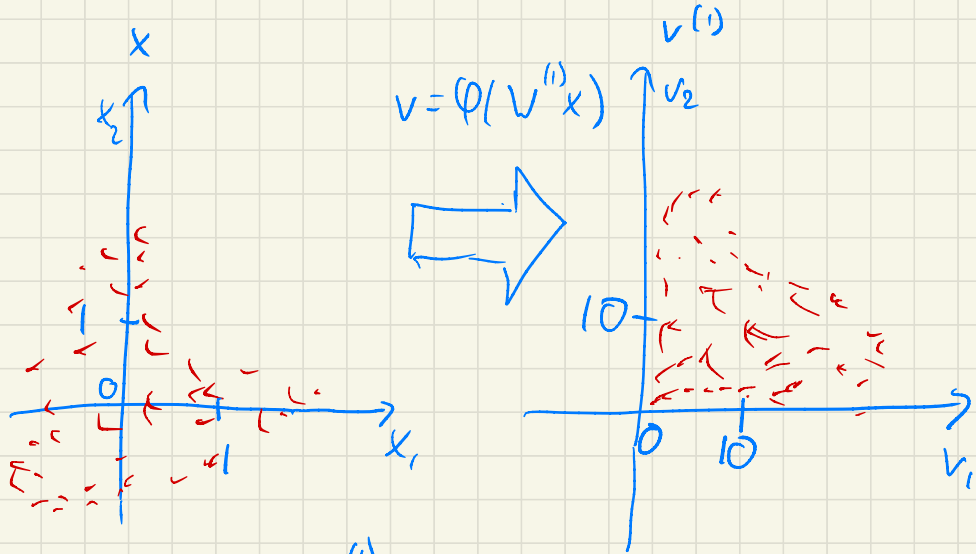
(1) vanishing gradients $\nabla_{W^{(l)}} l \rightarrow 0$

(2) exploding gradients $\nabla_{W^{(l)}} l \rightarrow \infty$

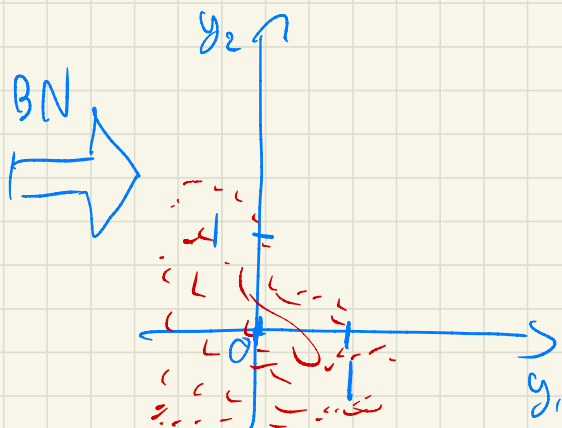
One possible reason for (1)/(2) is $v^{(l)} \rightarrow 0 / \infty$

He-initialization tries to ensure that $\text{Var}(v^{(l)}) = \text{const}$
at initialization

Batch-normalization aims to ensure standardized
activations throughout training



$$y^{(1)} = \text{BN}(v^{(1)})$$



Effect of batch normalization