

$$J(w) = \underbrace{R(w)}_{\text{avg loss}} + \underbrace{2\lambda \|w\|^2}_{\text{Regularizer}}$$

$$\nabla J(w) = \nabla R(w) + 2\lambda w$$

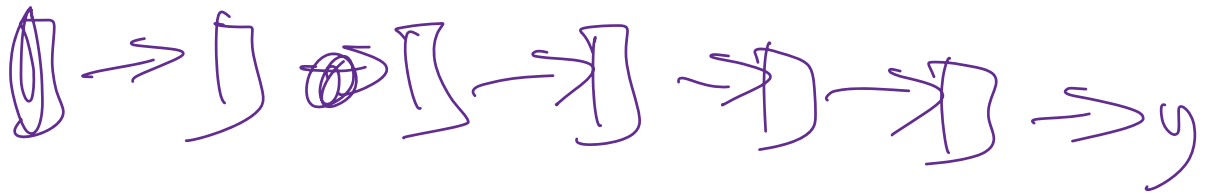
Gradient Step:

$$\begin{aligned} W^{\text{new}} &= W^{\text{old}} - \alpha (\nabla J(w^{\text{old}})) \\ &= W^{\text{old}} - \alpha \nabla R(w) + 2\alpha \lambda W^{\text{old}} \end{aligned}$$

Say you gradient step only on regularizer

$$\begin{aligned} W^{\text{new}} &= W^{\text{old}} - 2\alpha \lambda W^{\text{old}} \\ &= W^{\text{old}} \underbrace{(1 - 2\alpha \lambda)}_{\text{shrinkage}} \end{aligned}$$

Many layers



$W^{(2)}$

$$\frac{\partial R}{\partial W^{(2)}} = \frac{\partial W^{(2)}}{\partial f(a)} \frac{\partial f(a)}{\partial W^{(2)}}$$

All small numbers, makes
final gradient small

Layer Normalization

~~No~~

Normally:

$$h^1 = f(w^1 x)$$

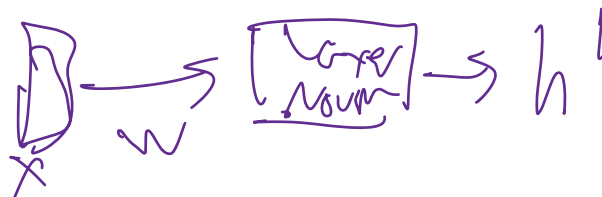
$$h^2 = f(w^2 h^1)$$

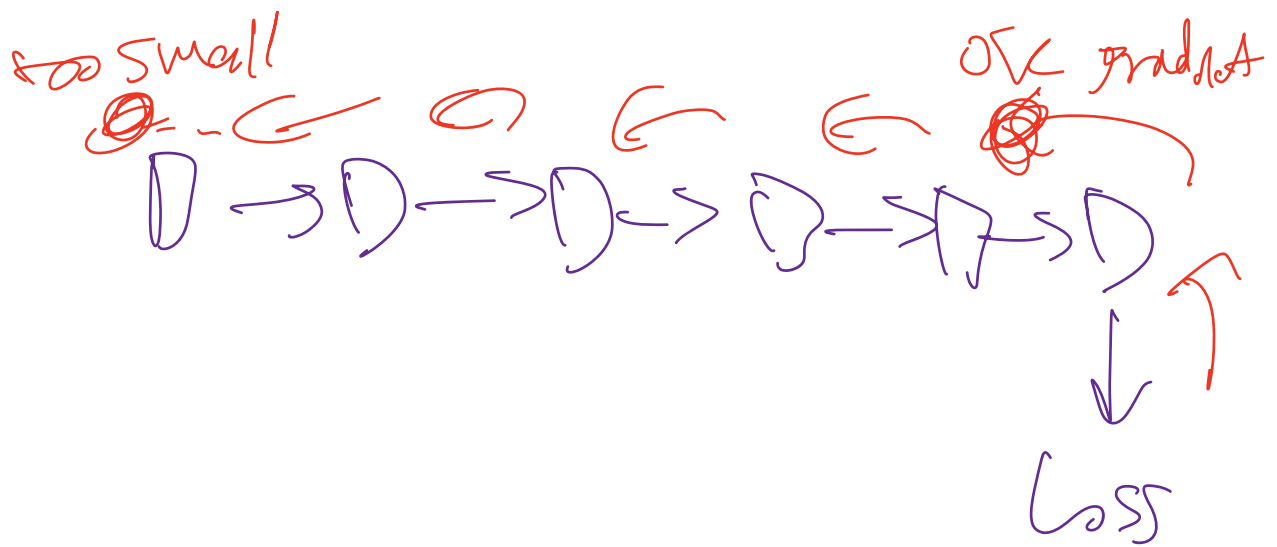
with layer norm

$$h^1 = \frac{f(w^1 x)}{\|f(w^1 x)\|}$$

so values
are
nearly
scal

etc. ✖





→ forward pass

→ backward pass = chain rule