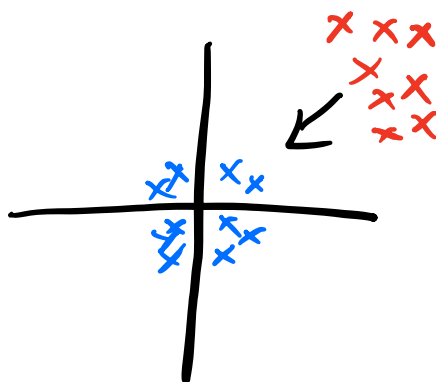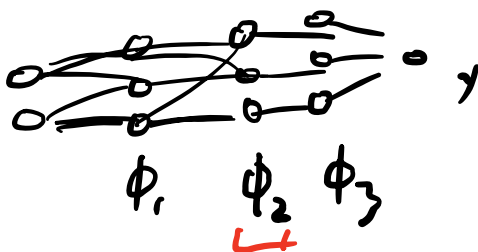# Last time

Normalization: keep data at consistent mean and variance

## Advantages

$$BN(x) = \frac{x - E[x]}{\sqrt{Var[x]}}$$

$$\left|\frac{dL}{dw_1}\right| = |x| \, |\sigma'(xw_1)| \underbrace{\left|\prod_{i=1}^{L}\right] \frac{d\phi_i}{d\phi_{i-1}}\right|}$$

$$\underbrace{m}$$



$$\phi_1 \quad \phi_2 \quad \phi_3$$



$$\left|\frac{dL}{dw_i}\right| = \left|\phi_{i-1}\right| \left|\sigma'(\phi_{i-1}w_i)\right| \left|\prod_{j=i+1}^{L}\right| \left|\frac{d\phi_j}{d\phi_{j-1}}\right|$$

$$BN(\phi_{i-1}) = \frac{\phi_{i-1} - E[\phi_{i-1}]}{\sqrt{Var(\phi_{i-1}) - \varepsilon}}$$
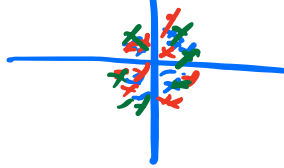
w/o Normalization
$$\phi_{i-1}$$



$$\phi_{i-1} = \sigma(xw)$$

changes

# w/ normalization

# Batch Norm

$E[\phi_{i,}] = 0 \quad Var = 1$

$E[\ ] = 0 \quad Var$

$\vdots$

# Layer Norm

$E[x] = 0$
$Var = 1$

$= X$

$N \times d$

\# obj \quad \# neurons

obs

# Batch norm

## Residual Networks



$$\phi_1 \qquad \phi_2 \qquad \phi_3$$

$$X \qquad \phi_1 \qquad \phi_2 \qquad \phi_3 \qquad \rightarrow \text{Loss}$$

$$\phi_1 = \sigma(xw_1) \qquad \phi_2 = \sigma(\phi_1 w_2) \dots$$

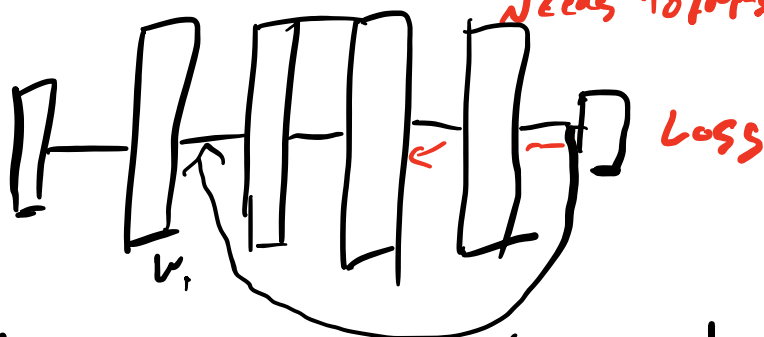$$\frac{dL}{dw_1} = x\sigma'(x^T w_1) \frac{d\phi_2}{d\phi_1} \cdot \frac{d\phi_3}{d\phi_2} \dots \frac{dL}{d\phi_L}$$

Needs to forget Info about Loss



$$\frac{dL}{dw_1} = \underbrace{x\sigma'(xw_1)} \left[ \left( \prod_{i=1}^{L} \frac{d\phi_i}{d\phi_{i-1}} \right) \frac{dL}{d\phi_L} + \frac{dL}{d\phi_1} \frac{d\phi_1}{dw_1} \right]$$

## Approach: Residual function

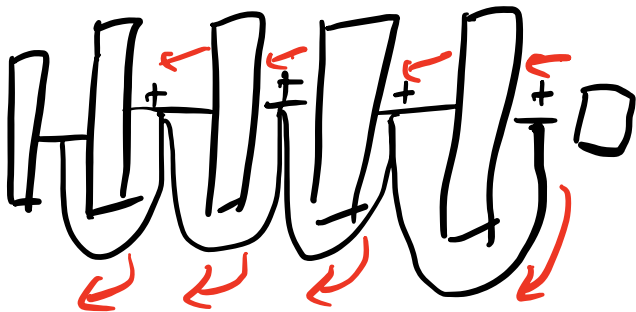$$\hat{f}(x) \longrightarrow f(x) = \underline{\hat{f}(x)} + x$$

# Residual Layer

Layer: $\phi_i = \sigma(\phi_{i-1}^T W_i)$

<span style="color:red">Computes next Value</span>

Residual Layer: $\phi_i = \sigma(\phi_{i-1}^T W) + \phi_{i-1}$

$\hat{\phi}_i$  <span style="color:red">$\frac{d}{d\phi_{i-1}} = 1$</span>

<span style="color:red">$\hat{\phi}_i = \sigma(\phi_{i-1} w) = \phi_i - \phi_{i-1}$</span>

<span style="color:red">residual</span>



## Backprop step

$$\frac{dL}{d\phi_{i-1}} = \frac{dL}{d\phi_i} \cdot \frac{d\phi_i}{d_{i-1}} = \frac{dL}{d\phi_i}\left[\frac{d}{d\phi_{i-1}}\sigma(\phi_i w_i)\right] + \frac{dL}{d\phi_s}$$

$$= \frac{dL}{d\phi_i}\left[\frac{d}{d\phi_{i-1}}\left(\sigma(\phi_{i-1} w_i) + \phi_{i-1}\right)\right]$$

$$= \frac{dL}{d\phi_i}\left[w_i \sigma'(\phi_{i-1}) + 1\right]$$

w/o residual

$$\frac{dL}{dw_i} = x\sigma'(xw_i)\prod_{i=1}^{L}\left(\frac{d\phi_i}{d\phi_{i-1}}\right)\frac{dL}{d\phi_L}$$

w/ residual

$$\frac{dL}{dw_1} = x\sigma'(xw_1) \cdot \prod_{i=1}^{2}\left(\frac{d\hat{\phi}_i}{d\phi_{i-1}} + 1\right)\frac{dL}{d\phi_2}$$