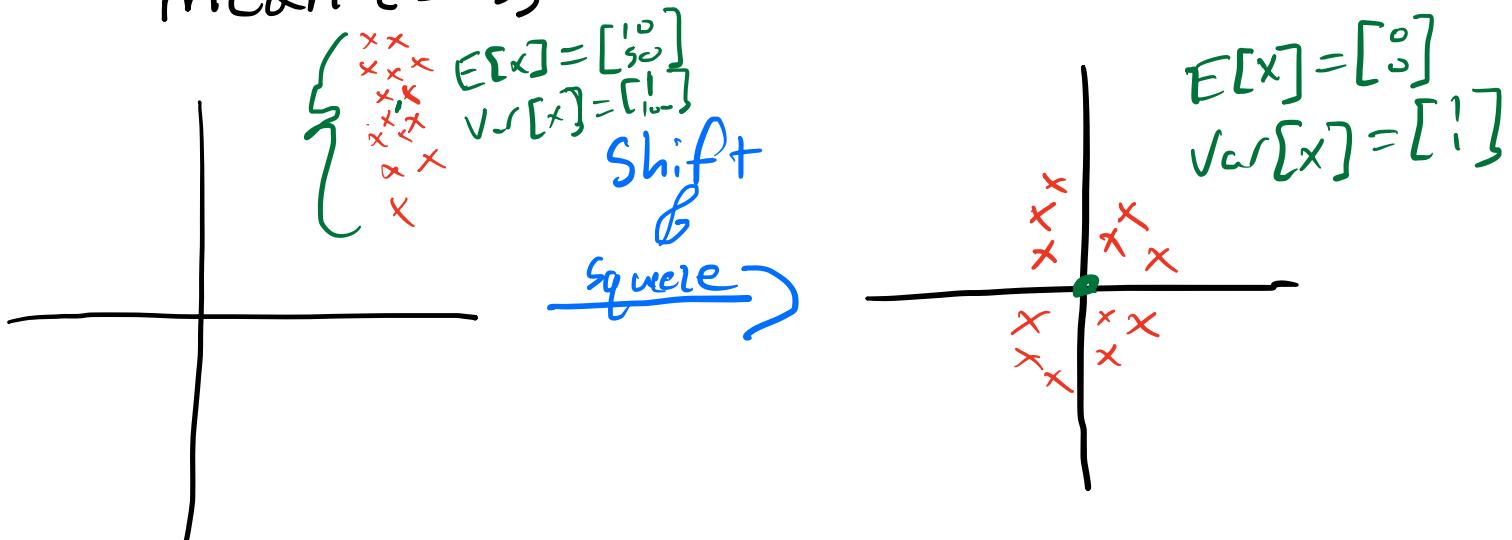


Last time

Normalization: keep data at consistent mean ($E[x]$) and variance ($\text{Var}[x]$).



Motivation

1) Gradient Scale

$$\left| \frac{dL}{dw_i} \right| = \|x\| \left| \sigma'(xw_i + b_i) \right| \prod_{j=2}^l \left| \frac{d\phi_j}{d\phi_{j-1}} \right| \dots$$

(Prevent scale of data from affecting gradient)

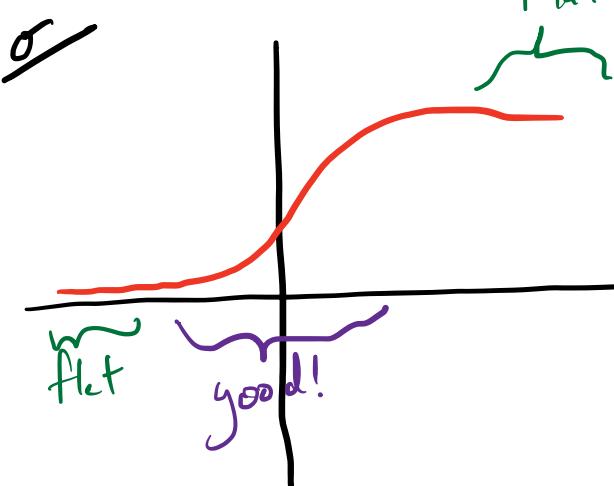
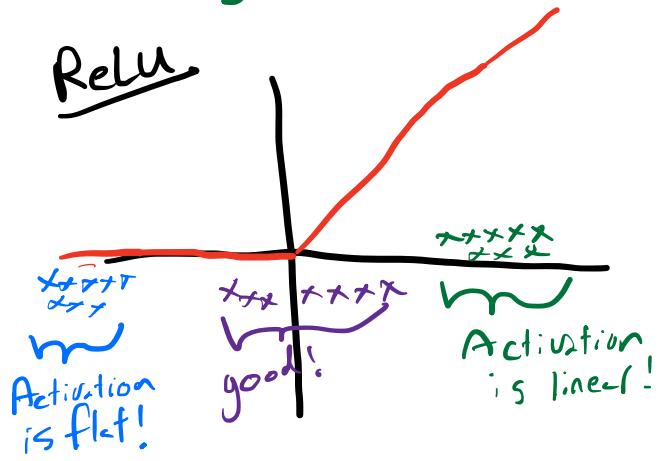
2) Scale mismatch

e.g. $x_1 = \begin{bmatrix} 2152 \\ 0.02 \end{bmatrix}$ $x_2 = \begin{bmatrix} 1.092 \\ 0.017 \end{bmatrix} \dots$

We saw why this is a problem!

Motivation

3) Activations only make sense near 0



Normalization

$$\text{Norm}(x) = \frac{x - E[x]}{\sqrt{\text{Var}[x] + \epsilon}}$$

Prevent / div. by 0!
 $\epsilon \ll 1$

Estimate $E[x]$, $\text{Var}[x]$ from dataset or Batch

$$E[x] \approx \bar{x} = \frac{1}{N} \sum_{i=1}^N x_i$$

$$\text{Var}[x] \approx s^2 = \frac{1}{N-1} \sum_{i=1}^N \left(x_i - \bar{x} \right)^2$$

\nwarrow Unbiased estimator

$$\begin{aligned} E[Norm(x)] &= E\left[\frac{x - E[x]}{\sqrt{Var(x) + \epsilon}}\right] \\ &= \frac{1}{\sqrt{Var(x) + \epsilon}} [E[x] - E[x]] = 0 \end{aligned}$$

Batch

Normalization

$$\text{Batch Norm}(x) = \frac{x - E[x]}{\sqrt{\text{Var}[x] + \epsilon}}$$

Prevent
div. by 0!

Estimate $E[x]$, $\text{Var}[x]$ from dataset or Batch

$$E[x] \approx \bar{x} = \frac{1}{B} \sum_{i=1}^B x_i$$

Batch sampled for SGD

$$\text{Var}[x] \approx s^2 = \frac{1}{B-1} \sum_{i=1}^B (x_i - \bar{x})^2$$

Unbiased estimator

Deep Network

$$f(x) = \phi_\ell(\phi_{\ell-1}(\dots \phi_1(x))) w_o + b_o$$

e.g. *Composing feature functions*

$$f(x) = \sigma(\sigma(\dots \sigma(x^T w_1 + b_1)^T w_2 + b_2) \dots)^T w_\ell + b_\ell) w_o + b_o$$

or

$$\phi_1 = \sigma(x^T w_1 + b_1)$$

$$f = \phi_\ell^T w_o + b_o$$

$$\phi_2 = \sigma(\phi_1^T w_2 + b_2)$$

$$L = \text{Loss}(f, y)$$

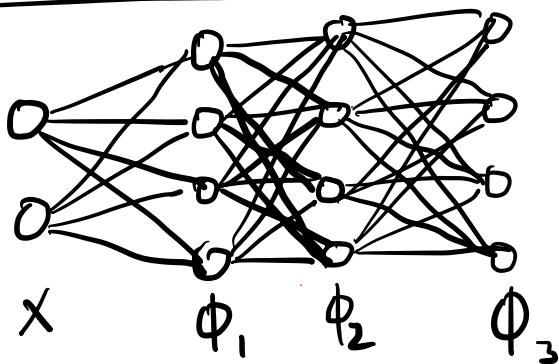
\vdots

$$\phi_\ell = \sigma(\underline{\phi_{\ell-1}^T w_\ell + b_\ell})$$

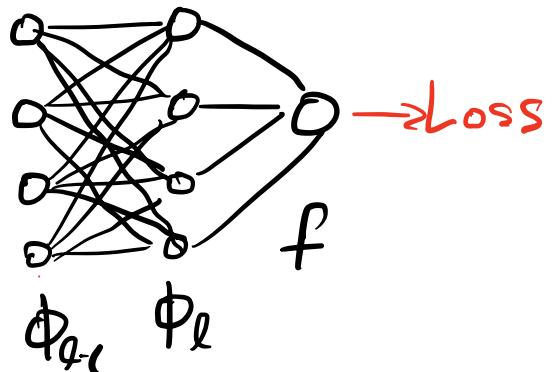
e.g. $\downarrow (f - y)^2$ MSE

Deep Network

l -layer network



...
...



or

$$\phi_1 = \sigma(x^T W_1 + b_1)$$

$$\phi_2 = \sigma(\phi_1^T W_2 + b_2)$$

⋮

$$\phi_q = \sigma(\underline{\phi_{q-1}^T W_q + b_q})$$

$$f = \phi_q^T w_o + b_o$$

$$L = \text{Loss}(f, y)$$

$$\text{e.g. } \downarrow (f - y)^2$$

$$E[x] = 0 \quad \text{Var}[x] = 1 \quad \text{w/ Norm.}$$

$$E[\phi_{q-1}] = ? \quad \text{Var}[\phi_{q-1}] = ?$$

Depend on $W_1 \dots W_{q-1}$ which could change.

Distribution Shift

Φ_{l-1}

$\begin{matrix} \times & \times \\ \times & \times \\ \times & \times \end{matrix}$
 $\begin{matrix} \times & \times & \times \\ \times & \times & \times \\ \times & \times & \times \end{matrix}$
 after k+2

$w_1 \dots w_{l-1}$
 change again

$\begin{matrix} \times & \times & \times \\ \times & \times & \times \\ \times & \times & \times \end{matrix}$
 Φ_{l-1} after
step K

$w_1 \dots w_l$

change
after
k+1

Layer l

$$\phi_l = \sigma(\phi_{l-1} w_l + b_l)$$

$$w_l^{(k+1)} \leftarrow w_l^{(k)} - \alpha \phi_{l-1}^{(k)} \sigma'(\phi_{l-1}^{(k)} w_l^{(k)} + b_l)$$

update using old $\phi_{l-1}^{(k)}$

$$f(x) = \sigma(\phi_{l-1}^{(k+1)} w_l + b_{l+1}) \dots$$

Predict using new $\phi_{l-1}^{(k+1)}$

Batch Norm

Φ_{l-1}

$\begin{matrix} \times & \times \\ \times & \times \\ \times & \times \end{matrix}$
 $\begin{matrix} \times & \times & \times \\ \times & \times & \times \\ \times & \times & \times \end{matrix}$
 after k+2

$w_1 \dots w_{l-1}$
 change again

$\begin{matrix} \times & \times & \times \\ \times & \times & \times \\ \times & \times & \times \end{matrix}$

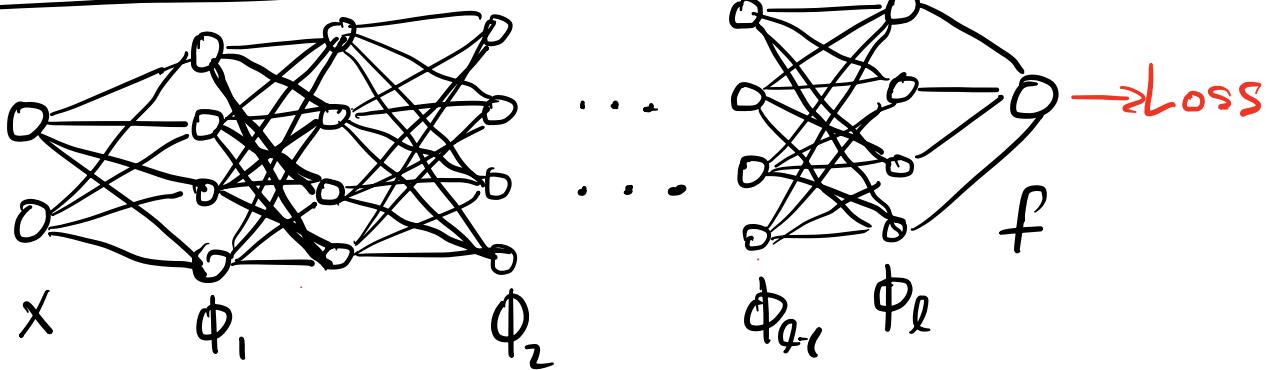
w

BatchNorm(ϕ_{l-1})

$\begin{matrix} \times & \times & \times \\ \times & \times & \times \\ \times & \times & \times \end{matrix}$
 $\begin{matrix} \times & \times & \times & \times \\ \times & \times & \times & \times \\ \times & \times & \times & \times \end{matrix}$

Can apply BatchNorm in the network

Deep Network w/ Batch Norm!



$$\phi_1 = \sigma(BN(x)^T w_1 + b_1) \quad f = BN(\phi_q)^T w_o + b_o$$

$$\phi_2 = \sigma(BN(\phi_1)^T w_2 + b_2) \quad L = \text{Loss}(f, y)$$

\vdots

$$\phi_q = \sigma(BN(\phi_{q-1})^T w_q + b_q)$$

e.g. $\downarrow (f - y)^2$

$$E[x] = 0 \quad \text{Var}[x] = 1 \quad \text{w/ Norm.}$$

$$E[BN(\phi_{q-1})] = 0 \quad \text{Var}[BN(\phi_{q-1})] = 1$$

Deep Network w/ Batch Norm

$$E[\text{BN}(\phi_{l-1})] = 0 \quad \text{Var}[\text{BN}(\phi_{l-1})] = 1$$

$E[\phi_{l-1}]$, $\text{Var}[\phi_{l-1}]$ still change every step

\therefore We need to use Batch Norm and update $\bar{x} \approx E[x]$ and $s^2 \approx \text{Var}[x]$ at every step
 $(E[\phi_{l-1}] \text{ for layer } l)$

Test time

Problem: At test time we might want to make a prediction on a single obs.

Single obs. $\rightarrow B = 1$

$$\bar{x} = \frac{1}{1} x_1 \quad s^2 = \frac{1}{1-1} (x_1 - \bar{x})^2$$

$$x_1 - \bar{x} = 0! \quad = 0 \rightarrow \text{divide by 0!}$$

Batch Norm at test time

In practice: Maintain Running Average of mean and Variance estimates as we go. (using EMA!)

$$\text{At step } k: \bar{X}_{\text{Avg.}}^{(k)} \leftarrow \beta \bar{X}_{\text{Avg.}}^{(k-1)} + (1-\beta) \bar{X}^{(k)}$$

$\beta \in [0, 1]$
typically 0.9
Est. using current batch

$$S_{\text{Avg.}}^2 \leftarrow \beta S_{\text{Avg.}}^2 + (1-\beta) S^2$$

At test:

$$\text{Batch Norm}_{\text{test}}(x) = \frac{x - \bar{X}_{\text{Avg.}}}{\sqrt{S_{\text{Avg.}}^2 + \epsilon}}$$

Use the same estimates every time we evaluate

Layer Norm

What if we don't want to do separate things at train and test time?

Batch norm

$$\begin{bmatrix} x_{11} \\ x_{12} \\ x_{13} \end{bmatrix}, \begin{bmatrix} x_{21} \\ x_{22} \\ x_{23} \end{bmatrix}, \begin{bmatrix} x_{31} \\ x_{32} \\ x_{33} \end{bmatrix}, \dots$$

$E[x_{ni}] = 0$
 $\text{Var}[x_{ni}] = 1$

Find $\bar{x} \approx E[x]$, $s^2 \approx \text{Var}[x]$
by Averaging over observations
for each dimension

Layer Norm

$$\begin{bmatrix} x_{11} \\ x_{12} \\ x_{13} \end{bmatrix}, \begin{bmatrix} x_{21} \\ x_{22} \\ x_{23} \end{bmatrix}, \begin{bmatrix} x_{31} \\ x_{32} \\ x_{33} \end{bmatrix}, \dots$$

$E[x_{i*}] = 0$
 $\text{Var}[x_{i*}] = 1$

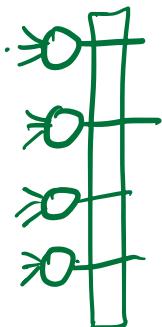
Alt. Average over dimensions within an observation

Layer Norm

$$\text{LayerNorm}(x_i) = \frac{x_i - \bar{x}_i}{\sqrt{s_i^2 + \epsilon}}$$

$$\bar{x}_i = \frac{1}{d} \sum_{j=1}^d x_{ij} \quad s^2 = \frac{1}{d-1} \sum_{j=1}^d (x_{ij} - \bar{x}_i)^2$$

In network



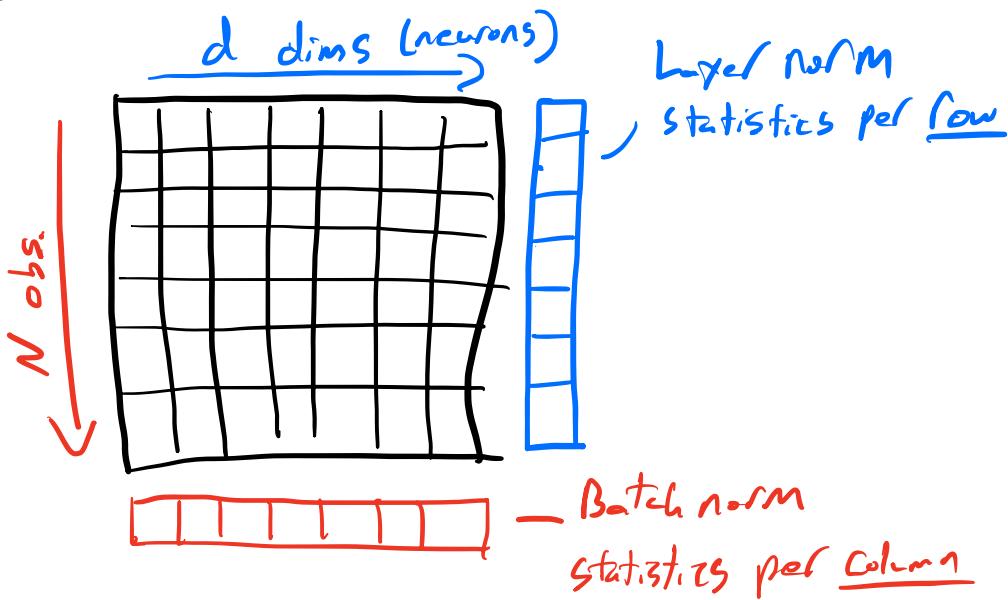
$$E[x] = 0, \text{var}[x] = 1$$

Can be applied w/
Batch size 1! As long
as $d \geq 1$

Batch Norm Vs. Layer Norm

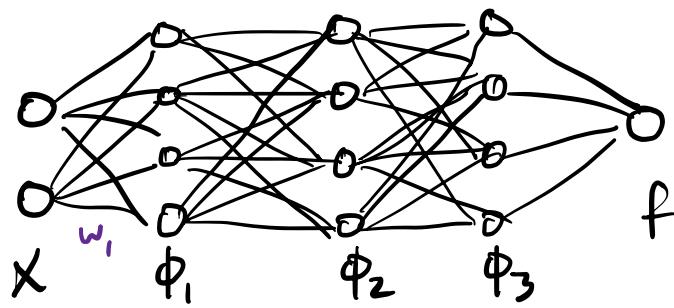
Input matrix

$$\underline{\underline{X}} \rightarrow N \times d$$



Batch Norm Usually preferred, but not always possible
[preserves more info about each observation]

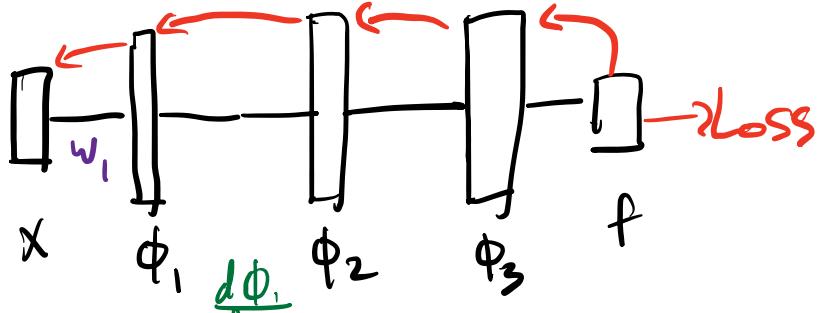
Residual Networks



- Loss

Exploding
 g_c

Simplified



$$\frac{dL}{dw_1} = x \sigma'(x^T w_1) \underbrace{\frac{d\phi_2}{d\phi_1} \cdot \frac{d\phi_3}{d\phi_2} \cdots \frac{df}{d\phi_k}}_{\text{Info from Loss gets diluted}} \cdot \frac{df}{df} \cdot \frac{d\text{Loss}}{df}$$

Info from Loss gets diluted
(Vanishing gradients)

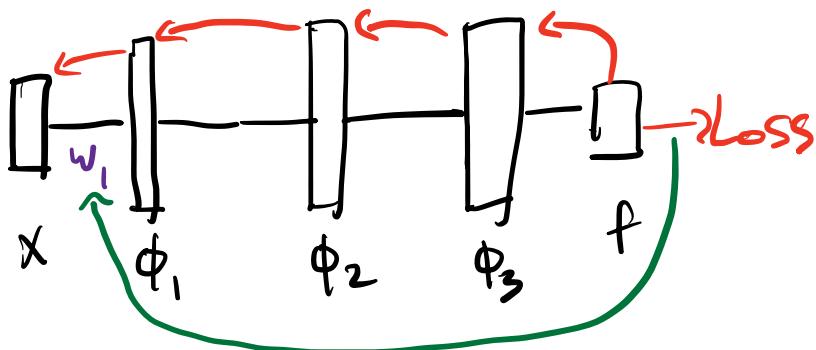
$$\rightarrow \frac{dL}{dw_1} = \frac{d\phi_1}{dw_1} \left(\frac{d\hat{\phi}_2}{d\phi_1} + 1 \right) \left(\frac{d\hat{\phi}_3}{d\phi_2} + 1 \right) \cdots \frac{df}{d\phi_k} \cdot \frac{d\text{Loss}}{df}$$

$$\frac{dL}{dw_1} = \frac{dp_1}{dw_1} \cdot \left[\dots + 1 \right] \frac{df}{d\phi_k} \cdot \frac{d\text{Loss}}{df} \prod_{i=1}^k \left(\frac{d\hat{\phi}_i}{d\phi_{i-1}} + 1 \right) \cdot \frac{df}{d\phi_1} \frac{d\text{Loss}}{df}$$

$$\frac{d\phi_1}{dw_1} \left(1 \right) \left(\frac{df}{d\phi_k} \cdot \frac{d\text{Loss}}{df} \right)$$

$$\left(\frac{d\hat{\phi}_2}{d\phi_1} + 1 \right) \left(\frac{d\hat{\phi}_3}{d\phi_2} + 1 \right)$$

Residual Networks



What if we had a direct connection?

$$\frac{dL}{dw_i} = x \sigma'(x^T w_i) \underbrace{\frac{d\phi_2}{d\phi_1} \cdot \frac{d\phi_3}{d\phi_2} \cdots}_{\left| \frac{d\phi_i}{d\phi_{i-1}} \right| \leq 1} \underbrace{\frac{df}{d\phi_i} \cdot \frac{d\text{Loss}}{df}}_{\text{Loss depends directly on } \phi_i} + \underbrace{\frac{dL}{d\phi_i} \cdot \frac{d\phi_i}{dw_i}}_{\rightarrow 0}$$

Residual function

$$\hat{f}(x) \rightarrow f(x) = \hat{f}(x) + x$$

Add input back to output!

Residual Layer

$$\text{Layer: } \phi_i = \sigma(\underbrace{\phi_{i-1}^T w_i + b_i}_{\text{Next Value}})$$

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

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

$$= \underbrace{\sigma'(\phi_{i-1}^T w_i + b_i)}_{w_i + 1} w_i + 1$$

$$\frac{d\phi_i}{d\phi_{i-1}} = \frac{d\hat{\phi}_i}{d\phi_{i-1}} + 1$$