$f(x)$

$x \quad \boxed{\phi_1} \quad \phi_2 \quad \cdots \quad \phi_L \quad f(x) \to Loss$

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

$\phi_L = \sigma(\phi_{L-1} v_L + b) \quad f = \phi_L w + b$

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

$L = Loss(f)$

$\dfrac{d\phi_1}{dw_1}$

$$\dfrac{dL}{dw_1} = x\,\sigma'(xw_1 + b) \underbrace{\quad}_{\frac{d\phi_1}{dw_1}} \dfrac{d\phi_2}{d\phi_1} \cdot \dfrac{d\phi_3}{d\phi_2} \cdots \dfrac{d\phi_L}{d\phi_{L-1}} \cdot \dfrac{df}{d\phi_L} \cdot \dfrac{dL}{df}$$

$$\dfrac{dL}{dw_1} = x\,\sigma'(xw_1 + b)\left(\prod_{l=1}^{L} \dfrac{d\phi_l}{d\phi_{l-1}}\right) \cdot \dfrac{df}{d\phi_L} \cdot \dfrac{dL}{df}$$
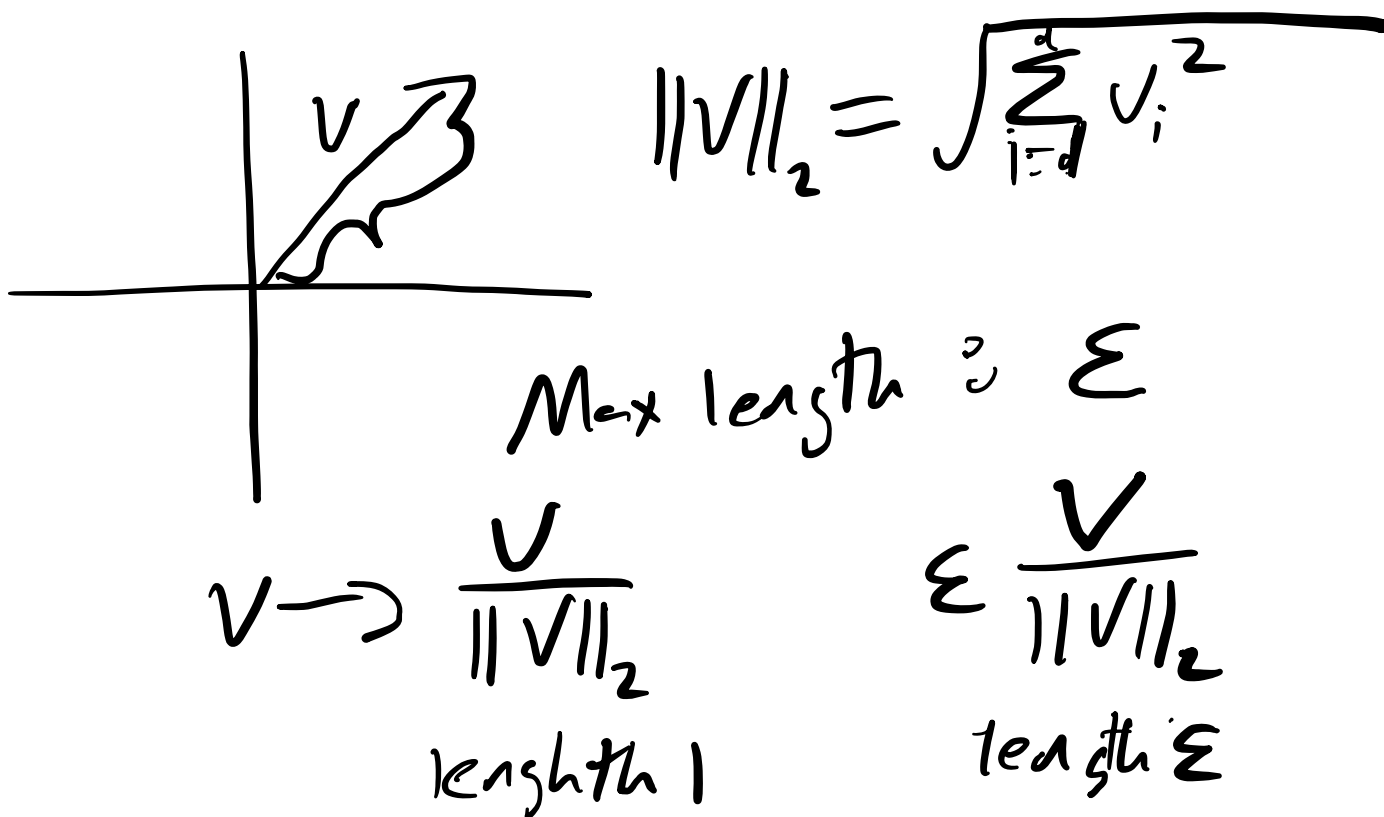
$$\left|\dfrac{dL}{dw_1}\right| = |x|\,|\sigma'(xw_1 + b)|\prod_{l=1}^{L}\left|\dfrac{d\phi_l}{d\phi_{l-1}}\right| \cdots$$

$$\left|\dfrac{d\phi_L}{d\phi_{L-1}}\right| \sim M \longrightarrow \left|\dfrac{dL}{dw_1}\right| \approx |x|\,|\sigma'|\,M^L$$

If $M > 1 \longrightarrow |\frac{dL}{dw}| >> 1$

gradient exploding

If $M < 1 \longrightarrow |\frac{dL}{dw_1}| \approx 0$

Vanishing gradients

## Gradient Clipping



$$\|V\|_2 = \sqrt{\sum_{i=1}^{d} v_i^2}$$

Max length : $\varepsilon$

$V \longrightarrow \frac{V}{\|V\|_2}$

length 1

$\varepsilon \frac{V}{\|V\|_2}$

length $\varepsilon$

Clipping by norm

Clip $(\frac{dL}{dw})$,     $\frac{dL}$

$$\longrightarrow \text{if } \left\| \frac{dL}{dw} \right\| > \varepsilon \longrightarrow \frac{\frac{dw}{dw}}{\left\| \frac{dL}{dw} \right\|} \cdot \varepsilon$$

$$\text{Otherwise} \longrightarrow \frac{dL}{dw}$$

## Clipping by Value

$$Clip_{v\text{-}loc}\left(\frac{dL}{dw}\right) = \begin{bmatrix} Max\left(Min\left(\frac{dL}{dw_i}, \varepsilon\right), -\varepsilon\right) \\ \vdots \\ \vdots \end{bmatrix}$$



## Input normalization

$$x \sim Data \qquad \underset{\text{\color{red}Issues}}{\mathbb{E}[x] >> 0}$$

$$Var[x] >> 1$$

Normalize(X)

$$X \rightarrow \frac{X - E[X]}{\sqrt{Var[X] + \varepsilon}} \qquad \varepsilon \ll 1$$



$$E[X] = ? \qquad Var[X] = ?$$

$$\bar{X} = \frac{1}{N} \sum_{i=1}^{N} X_i$$

$$S^2 = \frac{1}{N-1} \sum_{i=1}^{N} (X_i - \bar{X})^2$$

(unbiased est. of Var)

$$S^2 = \frac{1}{N} \sum_{i=1}^{N} (X_i - \bar{X})^2$$

$$\text{Normalize}(x) = \frac{X - \bar{X}}{\sqrt{s^2 + \varepsilon}} \quad \text{(biased est. of Var)}$$

# Batch Normalization

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

$$\phi_i = \sigma\left(\phi_{i-1} W + b\right) \quad \text{w/o BatchNorm}$$

$$\phi_i = \sigma\left(BN(\phi_{i-1}) W + b\right) \quad \text{w/ BatchNorm}$$

or

$$\phi_i = \sigma\left(BN(\phi_{i-1} W + b)\right)$$

Est.

$$BN(x) = \frac{X - \bar{X}}{\sqrt{s^2 + \varepsilon}}$$

## Gradient descent

$$\bar{x} = \frac{1}{N} \sum_{i=1}^{N} x_i$$

$$s^2 = \frac{1}{N-1} \sum_{i=1}^{N} (x_i - \bar{x})^2$$

## Stochastic gradient descent

Minibatch of $B$ obs.

$$\bar{x} = \frac{1}{B} \sum_{i=1}^{B} x_i, \quad s^2 = \frac{1}{B-1} \sum_{i=1}^{B} (x_i - \bar{x})^2$$

$$B > 1 \qquad BN(x) = \frac{x - \bar{x}}{\sqrt{s^2 + \epsilon}}$$

Training

## Test time

while training

$$\underbrace{\bar{\mu}^{(k+1)}}_{\substack{SS \\ E[x]}} \longleftarrow B\bar{\mu}^{(k)} + (1-B)\bar{x}^{(k)}$$

$$\overline{\sigma}^{2(k+1)} \leftarrow \beta\overline{\sigma}^{2(k)} + (1-\beta)S^{2(k)}$$
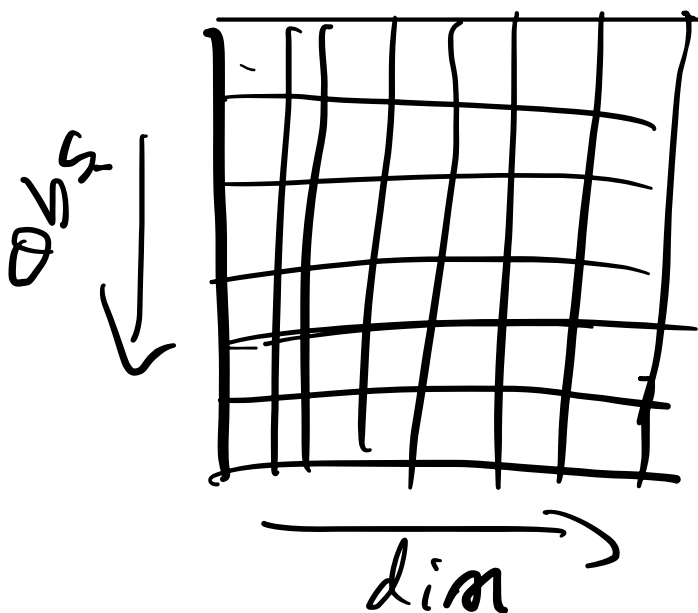
$$\underset{\substack{SS \\ Var[x]}}{} $$

$$\underset{test}{BN(x)} = \frac{x - \overline{\mu}}{\sqrt{\overline{\sigma}^2 + \varepsilon}}$$

## Layer normalization

$$LN(\underline{x}) = \frac{x - \frac{1}{d}\sum_{i=1}^{d}x_i}{\sqrt{\frac{1}{d-1}\sum_{i=1}^{d}(x_i - \overline{x}) + \varepsilon}} = \frac{\overline{x} - \overline{x}}{\sqrt{s^2 + \varepsilon}}$$

$$\overline{x} = \frac{1}{d}\sum_{i=1}^{d}x_i$$



Batch Norm
Normalize columns

obs

dim

LayerNorm Normalize rows